

July 2000

Differentiated service performance analysis

L. V. Nguyen
University of Wollongong

A. Eyers
University of Wollongong, teyers@uow.edu.au

Joe F. Chicharo
University of Wollongong, chicharo@uow.edu.au

Follow this and additional works at: <https://ro.uow.edu.au/infopapers>



Part of the [Physical Sciences and Mathematics Commons](#)

Recommended Citation

Nguyen, L. V.; Eyers, A.; and Chicharo, Joe F.: Differentiated service performance analysis 2000.
<https://ro.uow.edu.au/infopapers/221>

Research Online is the open access institutional repository for the University of Wollongong. For further information contact the UOW Library: research-pubs@uow.edu.au

Differentiated service performance analysis

Abstract

Differentiated service (DiffServ) has been proposed as an alternative for Integrated service. It aims to provide the same service to a group of flows that have similar quality of service requirements. Assured forwarding (AF) and expedited forwarding (EF) are two proposals for DiffServ provision. We present a performance analysis of an N drop-precedences threshold dropping (TD) queue, which is one of the proposed mechanisms for AF. In this analysis, traffic flows are assumed Poisson with exponentially distributed service time. We present simulation results that verify the analysis. This paper is an extension of the work attempted by Bolot et al. (see Proceedings of IEEE INFOCOM, 1999) and Sahu (see Umass CMPCSI Technical Report 99-09, University of Massachusetts) since it considers the general case with multiple classes of flow. We also show that the Poisson based analysis can be shown to hold for aggregation of bursty Markov sources in some cases and not to hold in others.

Keywords

Markov processes, Poisson distribution, exponential distribution, quality of service, queueing theory, telecommunication traffic

Disciplines

Physical Sciences and Mathematics

Publication Details

This paper originally appeared as: Nguyen, LV, Eysers, A & Chicharo, JF, Differentiated service performance analysis, Proceedings. ISCC 2000. Fifth IEEE Symposium on Computers and Communications, 3-6 July 2000, 328-333. Copyright IEEE 2000.

Differentiated Service Performance Analysis

Long V. Nguyen, Tony Eyers, Joe F. Chicharo
The Switch Network Research Centre
University of Wollongong, Australia
long@snrc.uow.edu.au

Abstract

Differentiated service (DiffServ) has been proposed as an alternative for Integrated Service. It aims to provide the same service to a group of flows that have similar Quality of Service requirements. Assured Forwarding (AF) and Expedited Forwarding (EF) are two proposals for DiffServ provision. We present a performance analysis of an N drop-precedences Threshold Dropping (TD) queue, which is one of the proposed mechanisms for AF. In this analysis, traffic flows are assumed Poisson with exponentially distributed service time. We present simulation results that verify the analysis. This paper is an extension of the work attempted by Bolot et al [8] and Sahu [9] since it considers the general case with multiple classes of flow. We also show that the Poisson base analysis can be shown to hold for aggregation of bursty Markov sources in some cases and not to hold in others.

Keywords— Threshold Dropping, Differentiated Service, Quality of Service, drop precedence, loss probability, expected delay, poissonian hypothesis and On-Off traffic

1. Introduction

The current Internet provides Best-Effort service with no specific performance guarantees for individual application. The IETF Integrated Service (IntServ) working group, formed to address this issue, has produced RSVP and service classes such as guaranteed-QoS, Controlled-Load. IntServ uses RSVP to provide network resources for individual flow [1], [2]. At each network node, the assigned bandwidth is maintained by a priority queueing algorithm, such as weighted Fair Queuing (WFQ) or Worst-case Fair Weighted Fair Queuing (WF²Q) [3]. WFQ and WF²Q are approximations of the idealised Generalised Process Sharing (GPS) mechanism [4]. The order of complexity for WFQ and WF²Q is $O(N^2)$

where N is the number of connections supported by the network node.

There are two related issues that arise from the implementation of RSVP and IntServ: the amount of overhead traffic and the scalability of this mechanism. IntServ aims to provide network resources for individual flows, thereby producing significant overhead traffic. It is impossible to implement RSVP and IntServ in wide area networks due to its poor scalability. In such networks, a router will have to support thousands or even millions QoS connections and becomes more and more complicated as the number of connections increase.

Differentiated service was proposed as an alternative to Intserv. It aims to provide the same service to a group of connections that have similar QoS requirements (whilst IntServ guarantees service requirement for individual connections by using RSVP). This helps lower the overhead, as network nodes have to handle only a small number of aggregations. Hence, it improves the efficiency of the network and DiffServ should also scale well in a larger network.

Recently, there have been two proposals for DiffServ provision: Assured Forwarding (AF) and Expedited Forwarding (EF). The AF schemes offer different levels of forwarding assurance for data packets received from a customer Diffserv domain [5]. In the current definition of AF, there are 4 traffic classes and within each traffic class, there are 3 drop precedences [6]. Packets of different application are given different drop precedence. More AF classes or levels of drop precedence may be defined for local use. Moreover, a DiffServ node must allocate a configurable, minimum amount of forwarding resources (buffer space and bandwidth to each implemented AF class [6]. Examples of AF mechanisms are Threshold Dropping, Random Early Detection In-profile/Out-profile (RIO) [8].

Meanwhile, in EF schemes higher priority packets receive preferential link access over lower priority packets [7]. During congestion periods, bandwidth is reallocated from low priority flows to high priority flows to minimise the delay and delay jitter [5]. Examples of EF mechanisms are Class Based Queuing (CBQ) [7] and

Priority Queuing [8]. In comparison, AF is a simpler mechanism to implement than EF since AF's buffer management is simpler than EF's packet scheduler. Moreover, low priority flows in AF are not significantly affected by higher priority flows.

Threshold Dropping (TD) is a queuing mechanism proposed to implement AF DiffServ. In a TD node, there is a buffer threshold assigned to each level of drop precedence [8]. IP Packets with higher drop precedence are more likely to be dropped during congestion. Within a class, flows of similar QoS requirements are given the same drop precedence (i.e. the same buffer threshold). A packet is discarded when the buffer exceeds the threshold corresponding to its drop precedence at its arrival.

Traffic characteristics are important parameters to determine the performance (loss probability and mean delay of packets) of a network. There have been a number of traffic models (eg. Poisson, MMPP, Gamma, etc) proposed to capture the characteristics of IP packets in a network. Hence, it is important to analyse the effect of traffic models on a network's performance so that Internet Service Provider can dimension and design DiffServ networks accordingly. In this paper, we will present an analytical approach to estimate packet loss probability and mean delay for poisson traffic (an well known model) when applied to the Threshold Dropping associated with DiffServ. The TD queue can be considered as an AF class with a configurable amount of forwarding resources and a number of drop precedences. This paper is an extension of the work attempted by Bolot *at al* [8] and Sahu [9] since it considers the general case with multiple classes of flow. We also show that the Poisson base analysis can be shown to hold for aggregation of bursty Markov sources in some cases and not to hold in others.

This paper is organised as follows: in Section II we present our analytical approach to calculate the packet loss probability and expected delay for an N drop-precedences TD queue (extension from the 2 drop-precedences TD queue). Section III presents simulation results that confirm the analytical results presented in Section II. Section III also presents simulation results for aggregation of MMPP traffic sources hence highlighting the validity of the analytical results. Section IV concludes the paper.

2. Analysis

In [8] the authors suggested a general approach for loss probability and expected delay calculations for AF mechanisms. In this Section, we extend this analytical approach for a TD queue with Poisson arrivals to the N drop-precedences case. Some adjustments for the mean delay calculation are also added.

In an N drop-precedences TD queue (Figure 1), there are N flows (each flow corresponds to a level of drop precedence) arriving at the queue. A packet is discarded at

its arrival when its corresponding buffer threshold has been reached or exceeded.

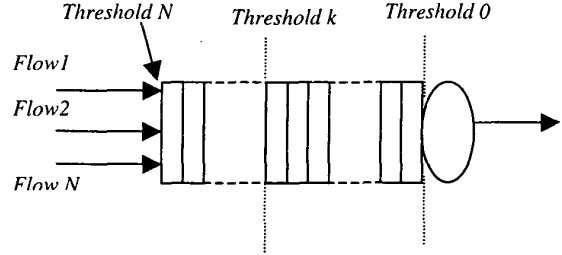


Figure 1. A Threshold Dropping Queue with N drop-precedences

This paper presents our analysis with the assumption that the incoming traffic flows are Poisson. We introduce the following terms:

- The arrival rate of the i^{th} priority flow is λ_i .
- The packet service times are exponentially distributed service times with mean $1/\mu$.
- The loads of the i^{th} priority flow and the aggregation are ρ_i and ρ respectively
- The buffer threshold of the i^{th} priority flow is L_i packets (L_0 is 0)
- At steady-state, the probability that there are n packets in the system is $\Pi(n)$
- $\alpha(n)$ is the acceptance probability of a packet which arrives to the queue seeing n other packets already in the system
- $\alpha_i(n)$ is the acceptance probability of an i^{th} priority packet which arrives to the queue seeing n other packets already in the system. For a TD queue, this probability can be determined as

$$\alpha_k(n) = \begin{cases} 1 & \text{if } n < L_k \\ 0 & \text{if } L_k \leq n \end{cases} \quad (1)$$

- p_i is the ratio of the i^{th} priority flow's load to the overall load. Hence, p_i is the ratio of λ_i over the sum of all arrival rates.

It is important to notice that the lower the drop precedence of a flow, the higher the priority of the flow (eg. the 1st priority flow has the lowest drop precedence and a buffer threshold of L_N , which is the buffer size of the queue). From the definition of $\alpha(n)$ and $\alpha_i(n)$ we have

$$\alpha(n) = \sum_{i=1}^N p_i \alpha_i(n) \quad (2)$$

$$\text{and } \alpha(n) = \begin{cases} p_1 + \dots + p_N & \text{if } n < L_1 \\ p_2 + \dots + p_N & \text{if } L_1 \leq n < L_2 \\ \dots & \dots \\ p_k + \dots + p_N & \text{if } L_{k-1} \leq n < L_k \\ \dots & \dots \\ p_N & \text{if } L_{N-1} \leq n < L_N \\ 0 & \text{if } L_N = n \end{cases} \quad (3)$$

It can be seen that this TD queue can be modelled as a birth-death process. For a state n , the birth rate is $\rho * \mu * \alpha(n)$ while the death rate is μ . The steady-state distribution of buffer content is:

$$\Pi(n) = \Pi(0) \rho^n \prod_{i=0}^{n-1} \alpha(i) \quad (4)$$

with the probability that the buffer is empty $\Pi(0)$

$$\Pi(0) = \left[\sum_{n=0}^{L_N} \rho^n \prod_{i=0}^{n-1} \alpha(i) \right]^{-1} \quad (5)$$

or

$$\Pi(0) = \left[1 + \sum_{i=1}^N \left(\prod_{j=1}^{i-1} (\rho_j + \dots + \rho_N)^{L_j - L_{j-1}} \sum_{k=1}^{L_i - L_{i-1}} (\rho_i + \dots + \rho_N)^k \right) \right]^{-1} \quad (6)$$

From (3) and (4), we obtain:

$$\Pi(n) = \Pi(0) \prod_{j=1}^{k-1} \left(\sum_{i=j}^N \rho_i \right)^{L_j - L_{j-1}} \left(\sum_{i=k}^N \rho_i \right)^{n - L_{k-1}} \quad \text{if } L_{k-1} < n \leq L_k \quad (7)$$

The loss probability of the i th priority flow is determined as:

$$\begin{aligned} \text{Loss}_i &= 1 - \sum_{n=0}^{L_N} \alpha_i(n) \Pi(n) \\ \Rightarrow \text{Loss}_i &= 1 - \sum_{n=0}^{L_i-1} \Pi(n) \quad (8) \end{aligned}$$

Using (6), (7) and (8), loss probability of the i th priority flow is:

$$\begin{aligned} \text{Loss}_i &= \Pi(0) \left[\prod_{j=1}^i (\rho_j + \dots + \rho_N)^{L_j - L_{j-1}} + \right. \\ &\quad \left. \sum_{j=i+1}^N \left(\prod_{k=1}^{j-1} (\rho_k + \dots + \rho_N)^{L_k - L_{k-1}} \sum_{k=1}^{L_j - L_{j-1}} (\rho_j + \dots + \rho_N)^k \right) \right] \quad (9) \end{aligned}$$

Clearly, when a packet arrives at the queue which already has n packets, it has a delay of n packets service times plus its own service time. Therefore, the mean delay of the i th priority flow (excluding rejected packets) is:

$$\text{Delay}_i = \frac{1}{\mu} \frac{\sum_{n=0}^{L_N-1} (n+1) \Pi(n) \alpha_i(n)}{\sum_{n=0}^{L_N-1} \Pi(n) \alpha_i(n)}$$

Substitute the values of $\alpha_i(n)$ into the above equation, we have

$$\Rightarrow \text{Delay}_i = \frac{1}{\mu} \frac{\sum_{n=0}^{L_i-1} (n+1) \Pi(n)}{\sum_{n=0}^{L_i-1} \Pi(n)} \quad (10)$$

Using (6), (7) and (10), the mean delay of the i th priority flow is:

$$\text{Delay}_i = \frac{1}{\mu} \frac{A_i}{B_i} \quad (11)$$

with

$$A_i = 1 + \sum_{k=1}^{i-1} \left[\prod_{j=0}^{k-1} (\rho_j + \dots + \rho_N)^{L_j - L_{j-1}} \sum_{n=1}^{L_k - L_{k-1}} (1+n+L_{k-1}) (\rho_k + \dots + \rho_N)^n \right] - (1+L_i) \prod_{j=1}^i (\rho_j + \dots + \rho_N)^{L_j - L_{j-1}} \quad (12)$$

and

$$B_i = 1 + \sum_{k=1}^{i-1} \left[\prod_{j=1}^{k-1} (\rho_j + \dots + \rho_N)^{L_j - L_{j-1}} \sum_{n=1}^{L_k - L_{k-1}} (\rho_k + \dots + \rho_N)^n \right] - \prod_{j=1}^i (\rho_j + \dots + \rho_N)^{L_j - L_{j-1}} \quad (13)$$

In the case of a 2 drop-precedences TD queue, we denote the loads of high and low priority flows and the aggregation are ρ^h , ρ^l and ρ respectively. The buffer size and the threshold of the TD queue is K and L packets.

The loss probability and mean delay can be determined from (9) and (11) with $N=2$ and $L_1=L$ and $L_2=K$. Hence the loss probabilities are:

$$\begin{aligned} \text{Loss Probability of} \\ \text{High Priority Packet} &= \Pi(0) (\rho_h + \rho_l)^L \rho_h^{K-L} \quad (14) \end{aligned}$$

and

$$\begin{aligned} \text{Loss Probability of} \\ \text{Low Priority Packet} &= \Pi(0) (\rho_h + \rho_l)^L \frac{1 - \rho_h^{K-L+1}}{1 - \rho_h} \quad (15) \end{aligned}$$

The probability that there is no packet in the system $\Pi(0)$ can be determined as:

$$\Pi(0) = \left[\sum_{n=0}^L (\rho_h + \rho_l)^n + (\rho_h + \rho_l)^L \sum_{n=1}^{K-L} \rho_h^n \right]^{-1} \quad (16)$$

The mean delay for high and low priority flows are:

$$\text{Delay}_{\text{high}} = \frac{1}{\mu} \frac{\sum_{n=0}^L (n+1) (\rho_h + \rho_l)^n + (\rho_h + \rho_l)^L \sum_{n=1}^{K-L} (n+1+L) \rho_h^n}{\sum_{n=0}^L (\rho_h + \rho_l)^n + (\rho_h + \rho_l)^L \sum_{n=1}^{K-L} \rho_h^n} \quad (17)$$

and

$$\text{Delay}_{\text{low}} = \frac{1}{\mu} \frac{\sum_{n=0}^{L-1} (n+1) (\rho_h + \rho_l)^n}{\sum_{n=0}^{L-1} (\rho_h + \rho_l)^n} \quad (18)$$

In these mean delay calculations, discarded packets are not included since the retransmission mechanism is not defined while in [8], the authors accounted for discarded packets. However, in [9], the authors provided delay calculations that include the probability for the system to be empty. The calculations for mean delays are:

$$\text{Delay}_{\text{high}} = \frac{\Pi(0)}{\mu} \left(1 + \sum_{n=1}^L (n+1)(\rho_h + \rho_l)^n + (\rho_h + \rho_l)^L \sum_{n=1}^{K-L} (n+1+L)\rho_h^n \right) \quad (19)$$

and

$$\text{Delay}_{\text{low}} = \frac{\Pi(0)}{\mu} \frac{\left(1 + \sum_{n=1}^L (n+1)(\rho_h + \rho_l)^n \right)}{\sum_{n=1}^L (\rho_h + \rho_l)^n} \quad (20)$$

These calculations imply the proportionality of mean delay to $\Pi(0)$. Hence as load increases, $\Pi(0)$ approaches zero and so does the expected delay. This contradicts the observation that the expected delay approaches K/μ for high priority packets and L/μ for low priority packets (Little's theory). In our analysis, the term $\Pi(0)$ is cancelled since it appears in both the numerator and denominator of the calculation.

In the next Section, we will present simulation results that verify our theoretical analysis. We will also compare our mean delay calculations with [8]'s delay calculations. Also, simulation results will be shown to raise the question if the analysis developed for Poisson traffic can be applied to aggregates of bursty Markov sources.

3. Simulation Results

In this section, we present the results obtained from simulation, compare them with analytical calculations. Also, we show our observation that the Poisson based analysis in some cases hold for bursty input traffic. In our simulation, data packets (packet size and time stamp) are generated based on the traffic model (Poisson and 2-state On-Off) to form traffic flows. These flows are fed into a TD queue model and simulation results (loss probability and mean packet delay) are measured and compared with estimated ones.

3.1. Verification of the analysis

A simulation was developed to obtain experimental results and verify our analytical approach. The results are used to compare with the calculation provided by Bolot *et al* [8]. This simulation was developed for the 2-class case and expanded to a multiple priority case. We repeat the experiment in [8] by introducing a high priority and a low priority flow to the TD queue. High priority packets account for 95% while low priority packets account for

5% of the load. The buffer size and threshold was set to 16 and 6 packets respectively since the traffic flows consist of a single source.

Figure 2 shows simulation results in comparison with our analysis while Figure 3 compares simulation results with [8]. The expected delay presented in our figures is normalised with respect to the mean service time of the queue ($1/\mu$). Figure 2 shows the loss probability and expected delay (for both priorities) obtained from the simulation and our analysis. It can be seen that our analyses closely match with simulation results (the points are on top of the analytical graphs)

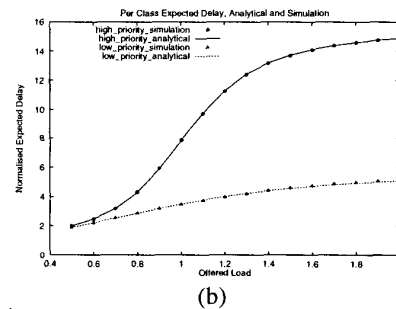
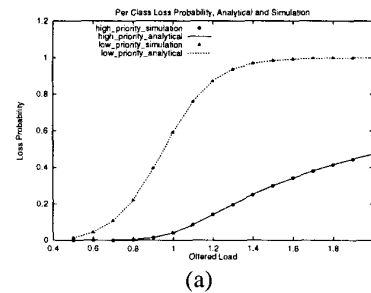


Figure 2. (a) Loss Probability and (b) Expected Delay (normalised with respect to Packet Mean Service Time) of High and Low Priority Packets as a Function of the Total Load – Analytical and Simulation Results. High Priority Packets Contribute 95% of the Load ($K=16$; $L=6$)

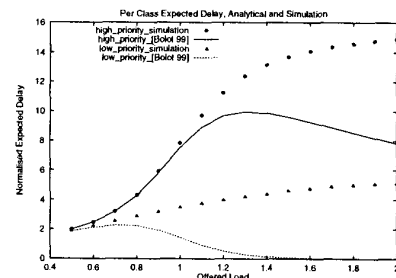


Figure 3. Expected Delay (normalised with respect to Packet Mean Service Time) of High and Low Priority Packets as a Function of the Total Load – [8]'s Approach

and Simulation Results. High Priority Packets Contribute 95% of the Load ($K=16; L=6$)

We have established that the delay calculation presented by Bolot *et al* [8] matches with the results obtained when discarded packets are included. Figure 3 shows that during the period while the queue is not heavily loaded, the calculated results are close to simulated ones since the number of discarded packets are insignificant. However, as the queue becomes heavily congested the delay calculation of Bolot *et al* [8] approaches zero since there are less and less packets that get accepted. Meanwhile, the normalised expected delay (excluding discarded packets) approaches the buffer threshold value.

The simulation was also developed to verify our analytical approach to determine loss and delay in a multiple priority TD queue. In this simulation, a 3-priority TD queue was implemented with 3 single-source flows corresponding to the 3-drop priorities. Medium and low priority flows are set with load of 0.7 and 0.4 accordingly while the load of the high priority flow is varied from 0.1 to 0.9. The buffer thresholds were set at 16, 12 and 8 packets for high, medium and low priority flow respectively. The loss probability and expected delay were measured and plotted as a function of the high priority load in Figure 4.

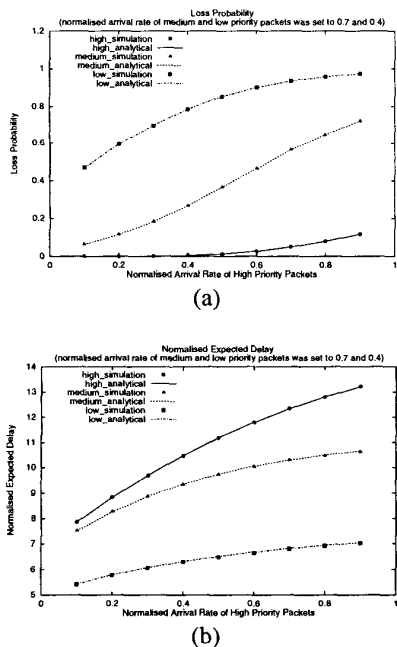


Figure 4. (a) Loss Probability and (b) Expected Delay (normalised with respect to Packet Mean Service Time) of High, Medium and Low Priority Packets as a Function of the Load of High Priority Flow—Analytical and Simulation Results ($L_3=16; L_2=12; L_1=8$)

It can be seen that the analytical results matches those obtained from simulation. Hence, this validates our analysis for a TD queue with a generalised number of drop precedences. In a DiffServ environment, this can be useful to estimate loss and delay where for each service class data packets are associated to more than two levels of discarding priority.

3.2. Poisson based analysis with bursty input sources

The analysis presented in Section II was developed with the assumption that input traffic are Poisson. It is clear that this Poisson based analysis does not cover the case when the input is a single bursty flow such as MMPP. However, if we alter the experiment by replacing the Poisson input sources by aggregations of bursty sources such as MMPP or On-Off, calculation and simulation results are shown to match in some cases and not in others. The Poisson parameter of the traffic model is calculated based on the On-Off parameters. Figure 5 presents loss and delay obtained from simulation with 2 single-source flows (high and low priority). The bursty sources are On-Off with duty cycle of 50%. These two sources are selected such that the high priority source contributes 95% of the total load while the low priority source contributes 5% of the total load.

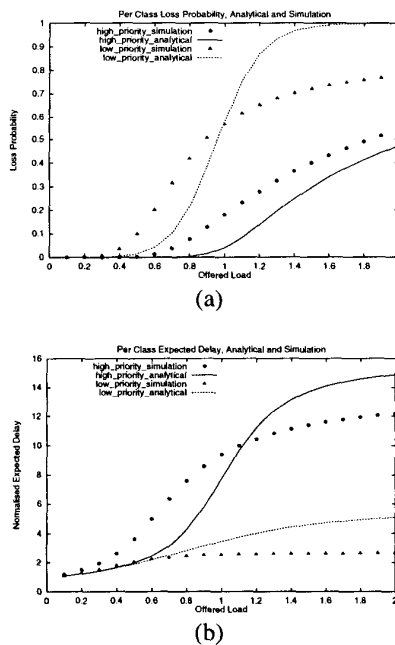


Figure 5. (a) Loss Probability and (b) Expected Delay (normalised with respect to Packet Mean Service Time) of High and Low Priority Packets as a Function of the Total Load – Analytical and Simulation Results. Both Streams

are On-Off with High Priority Packets contribute 95% of the Load ($K=16$; $L=6$)

In this simulation the overall load of the TD queue ranges from 0.1 (lightly loaded) to 2 (heavily loaded). It can be seen that there is a difference between analytical and simulation results for both loss and delay. However, if the single bursty sources are replaced by aggregation of bursty sources (in the simulation, 20 identical On-Off sources are treated as a flow), the obtained simulation results are close to our predicted results. Figure 6 shows the match between simulation results and analysis. Each On-Off source is the same as in the previous example. The high priority flow still contributes 95% of the packets while the low priority flow contributes only 5%.

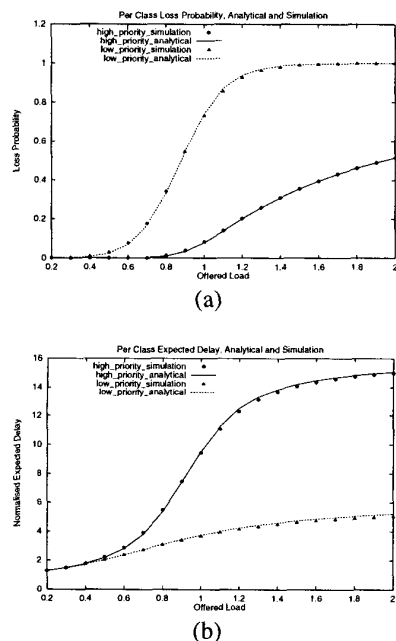


Figure 6. (a) Loss Probability and (b) Expected Delay (normalised with respect to Packet Mean Service Time) of High and Low Priority Packets as a Function of the Total Load – Analytical and Simulation Results. Both Streams are Aggregates of 20 identical On-Off sources. High Priority Packets contribute 95% of the Load ($K=16$; $L=6$)

Hence, it can be seen that for some cases (with large aggregation of bursty sources) the Poisson analysis holds while in other cases (small aggregations), it does not. As a result it is necessary to investigate the performance of a TD queue with other traffic models

4. Conclusions

This paper has presented an analytical approach to determine QoS metrics of a TD queue (loss probability and expected delay). This method can be applied to provide the solution for a TD queue with multiple

discarding priority TD queues. We have also corrected the discrepancies of the delay calculations provided in [8] and [9]. Our analytical calculations were verified with simulation results in Section III.

It is noted that our analytical approach allows a service provider to determine what performance is expected based on the traffic parameters as well as network resources in a multiple priority situation. These calculations can be used to help dimension the network to satisfy QoS requirements: loss and delay.

Moreover, we showed that the poisson hypothesis can not hold for single bursty traffic flow yet it can hold for large aggregation of bursty sources. The question arises is to determine the sufficient size of the aggregation so that the poisson analysis can be applied. Also, it emphasises the need to investigate DiffServ performances with different traffic models.

5. References

- [1]. White P., "RSVP and Integrated Services in the Internet: A Tutorial", *IEEE Communications Magazine*, May 1997.
- [2]. Zhang L., Deering S., Estrin D., Shenker S., Zappala D., "RSVP a new Resource Reservation Protocol", *IEEE Communications Magazine*, September 1993.
- [3]. Zhang H., "Service Disciplines for Guaranteed Performance Service in Packet-Switching Networks", *Proceeding of the IEEE*, October 1995.
- [4]. Parekh A. K., Gallager R. G., "A generalised Processor Sharing Approach to Flow Control in Integrated Services Networks: The Single-Node Case", *IEEE/ACM Transactions on Networking*, June 1993.
- [5]. Weiss W., "QoS with Differentiated Services", *Bell Labs Technical Journal*, Oct-Dec 1998.
- [6]. Heinanen J., Baker F., Weiss W., "Assured Forwarding PHB Group", Internet Draft, Feb. 1999, <http://www.ietf.org/internet-drafts/draft-ietf-diffserv-phb-af-06.txt>.
- [7]. Jacobson V., Nichols K., Poduri K., "An Expedited Forwarding PHB", Internet Draft, Feb. 1999, <http://www.ietf.org/internet-drafts/draft-ietf-diffserv-phb-ef-02.txt>.
- [8]. Bolot J.C., May M., Jean-Marie A., Diot C., "Simple Performance Models of Differentiated Services Schemes for the Internet", *Proceedings of IEEE INFOCOM*, March 1999.
- [9]. Sahu S., Towsley D., Kurose J., "A Qualitative Study of Differentiated Services for the Internet", *Umass CMPCSI Technical Report 99-09*, University of Massachusetts.