# End-to-End Delay Analyses
# in IP Networks

Rossitza. Iv. Goleva[1], Dimitar K. Atamian[2]

*Abstract -* **Quality of Service analyses in IP networks and especially end-to-end management in real time services is dynamic investigation area. The transmission of voice, audio, and video is sensitive to the delay and delay variation. Wired and wireless technologies set different requirements to the quality parameters. They also influence end-to-end delay in different ways. This paper proposes approximate analytical/simulation solution to end-to-end analyses in packet switched transmission. The analyses are made taking into account the class of the services, their delay bounds and codec implemented. End-to-end resource management is estimated using Network Signalling (NSIS) protocol. The derived results are applicable to IP scalable network planning, optimization, and congestion management.**

*Keywords -* **Packet network, IP, Quality of Service, shaping.**

## I. INTRODUCTION

Mixed IPv4 and IPv6 networks, wired and wireless solutions and hybrid networks with multimedia traffic carried are interesting investigation area. The technological circumstances are changing continuously and this requires dynamic Quality of Service estimation and management. IPv4 and IPv6 headers have different size. This influences the servicing rate at the router and switch interfaces. The minimal IPv4 header is 20 bytes. Minimal IPv6 header is 40 bytes. The next additional headers can be added in IPv6 datagram. The maximal length of IPv6 is 1280 bytes. The maximal length in IPv4 is 1500 bytes. The average packet length in both versions can vary significantly depending on the traffic nature and applications [1], [2].

Some of the interfaces apply traffic shaping and policing at packet level [3], [4], [5]. The policing are implemented at access points of the networks. Shaping is applied on the network and technology boundaries. In a typical end-to-end connection there are at least one policing point and few shaping points. The policing technique is capable to reject packets. Shaping technique applies additional delays to some of the packets and the gained capacity is assigned to the packets without enough delay reserve.

End-to-end delay and delay variation phenomena in TCP and UDP services under bursty traffic depends strongly on the traffic distribution, policing and shaping applied [6], [7]. The access points behave as worst case delay points under heavy and bursty load traffic.

[1]Rossitza Iv. Goleva is with the Department of Telecommunications, Technical University of Sofia, Bulgaria, Kl. Ohridski blvd. 8, Sofia. 1756, Bulgaria, E-mail: rig@tu-sofia.bg

[2]Dimitar K. Atamian is with Department of Telecommunications, Technical University of Sofia, Bulgaria, Kl. Ohridski blvd. 8, Sofia. 1756, Bulgaria, E-mail: dka@tu-sofia.bg

In this paper we investigate shaping influence on the end-to-end delay and delay variation. The work is based on the packet stream with Poisson arrival at session level and Deterministic arrival at packet level already observed during simulation [3], [8], [9]. Distribution of the packets at the queue entry and at the output of the router interface is mixed. The probability of the packet to wait and the probability of place and waiting losses are investigated further on [10], [11], [12].

The aim of the paper is to demonstrate the capability of the priorities and shaping techniques under typical interface load of 40 to 50%. The derived results are applicable to the routers that are capable to keep state parameters per session, estimate them and manage dynamically the queue and priority parameters.

We investigate the four mostly used techniques for traffic and Quality of Service (QoS) management – IntServ, DiffServ, RSVP, and NSIS. The analytical/ simulation approach is used for this purpose. It is based on the FIFO queue with priorities and limited waiting bound per priority.

## II. QoS MANAGEMENT

Integrated Services (IntServ) is a complex technique often called protocol that ensures Quality of Service in IP networks. It is applied usually in access routers or gateways. It tries to serve packets from different services in a different ways depending on the quality requirements. IntServ classifies services into three main classes depending on the traffic requirements as elastic, tolerant real-time and intolerant real-time:

- Elastic applications are served in a "best effort" discipline. The quality parameters cannot be guaranteed. It is applied for not time critical applications like email

- Tolerant real-time applications are delay sensitive and usually require high bandwidth. LAN-to-LAN connectivity is usually modelled this way

- Intolerant real-time applications require low delay and almost guaranteed bandwidth. VoIP service is intolerant to the delay

Differentiated Services (DiffServ) is another quality management technique that is more applicable for core networks. Due to its nature DiffServ applies its rules on aggregated traffic. After appropriate marking of the aggregated packets they are gathered in the way that is defined for their class. There are three main types of services we try to highlight in this paper:

- Premium service with low delay, low loss, guaranteed bandwidth applied for VoIP

- Assured service with less requirements to the delay and loss in comparison to the premium service for LAN-to-LAN connectivity

- Olympic service with no time requirements for email

Dynamic traffic management and especially Quality of Service management requires signalling protocol that is capable to confirm traffic contract end-to-end. The two protocols investigated here are Resource Reservation Protocol (RSVP) and Network Signalling (NSIS).

Resource Reservation Protocol (RSVP) is a technique useful for delay sensitive traffic like VoIP. Three types of services are identified for RSVP like:

- Wildcard filter with maximal requirements for given interface applied for LAN-to-LAN connectivity

- Shared explicit with maximal requirements for the interface taking into account called address. It is applied for email

- Fixed filter with full reservation for quality sensitive services like VoIP

Network Signalling (NSIS) protocol is a new generation of RSVP/ IntServ protocols that is capable to confirm end-to-end Quality of Service parameters. The new phenomenon in NSIS is in the distinction between signalling transport and signalling application in different layers. The signalling information is transported using TCP session or UDP protocol. The analyses of the signalling information can be transparent to some of the network nodes. It also can be analysed in those nodes where there is a need of quality estimation and management. The same packet filters like those in RSVP are applied. They are modified by means of reservation, traffic measurement and reconfiguration parameters. NSIS also keeps state parameters per session like IntServ. The protocol can update flow parameters, support multihoming, tunnelling and IPv4/ IPv6 traverse.

The transport part of the NSIS supports both datagram and virtual connection modes for signalling transport. It also associates security protocols. NSIS is the unique dynamic QoS protocol nowadays. The protocol node is capable also to implement shaping by adding variable delay in the packet flow depending on their quality requirements. Whenever packets are delayed this means adding additional delay and delay jitter. Therefore, shaping is bounded by end-to-end delay and delay jitter constrains [13], [14]. The appliance of NSIS in a network enable traffic contracts at all network interfaces – access, edge or core. NSIS is applicable for Service Level Agreement (SLA). It is applicable for customer profile specification and management.

## III. TRAFFIC SOURCES

Three types of traffic sources are assumed in an example wide area network – Voice over IP, LAN-to-LAN connectivity, email. LAN traffic is lower priority in comparison to the VoIP traffic and with higher priority in comparison to the email. The there services are mixed together with some assumptions. In Voice over IP (VoIP) service silence and talk intervals are exponentially distributed. On-off model is applied.

LAN emulation is specific with its sessions. Sessions are established for any Internet connections. Packet rate is higher in comparison to the VoIP. Session duration is low. The traffic source is behaving as on-off model with exponential duration of the silence and transmission intervals [5]. Emails are specific with packet exchange mostly in one direction. The service is not time demanding. Number of traffic sources is taken from the typical image in a business area. Packets are taken to be long. In VoIP traffic 200 bytes carry up to 20 milliseconds voice. This means that quality voice can be transmitted only in the area using up to 20-30 and even more hops.

The limits for waiting times are calculated under consideration of end-to-end delay for every service. Servicing times per packets are fixed on 100 Mbps line interface. Table I represents all the parameters for traffic sources in the model.

TABLE I
TRAFFIC SOURCES PARAMETERS

| Parameter | VoIP | LAN-to-LAN | Email |
|---|---|---|---|
| Pear rate, packets per second | 10-30 | 164-250 | 1-5 |
| Mean call/ session duration, sec | 180 | 20-50 | 10-50 |
| Mean duration between beginning of calls/ sessions, sec | 360 | 10 | 15 |
| Mean talk/ silence duration, sec | 10/20 | 50/10 | 20/10 |
| Distribution of call/series duration | Exp. | Exp. | Exp. |
| Traffic sources | 5000 | 500 | 1500 |
| Priorities | High | Medium | Low |

## IV. ANALYTICAL AND SIMULATION MODEL

Simulation is performed on C++ language. The pseudo exponential pseudo deterministic characteristics of the traffic sources are reached after usage of combination between many random generators. The queue behaviour is complex due to the priorities and limits on waiting times. Waiting times limits are calculated taking into account specific requirements of the four QoS techniques - IntServ, DiffServ, RSVP, NSIS. The bounds are calculated analytically, the statistical results are derived via simulation. Many parameters have been derived from the model like time and space loss probabilities, probabilities to wait for different types of traffic, queue lengths, waiting times at many interface points in the model like output of the traffic sources, input and output of the queue. Statistical accuracy of the derived results is proven by Student criterion.

The overall load of the interface is calculated with Eq. 1.

$$InterfaceLoad = \frac{T_{Occupancy,}}{T_{Modelling}},$$ (1)

Where InterfaceLoad is the overall occupancy of the interface

$T_{Occupancy}$ is the duration in seconds when the servicing module is occupied

$T_{Modelling}$ is the overall modelling time

The probability of packet loss is estimated with Eq. 2.

$$P_{PacketLoss} = \frac{PacketLoss}{TotalNoTransmitted}, \qquad (2)$$

Where PacketLoss is total number of lost packets

TotalNoTransmitted is total number of transmitted packets

Packet losses are divided into waiting bound losses and place losses. It is demonstrated further in this paper that the place losses dominate on the overall losses.

Mean waiting time is calculated by Eq. 3 dividing total duration of waiting packets and total number of waited packets.

$$MeanWaitinTime = \frac{WaitingTime}{NumberOfWaitedPackets} \qquad (3)$$

Most of the LAN traffic is considered to be TCP. The waiting time limits for such traffic depend on the round trip time of the TCP segments.

TCP applies many different mechanisms that allow retransmission and slow start in the session. The limits for slow start are different (Eq. 4). In typical TCP session of up to 15 hops waiting time limits for queues is function of the slow start limit. Otherwise, the interface will force all LAN sessions to decrease the transmission rate.

$$SlowStart = 2E_{RTT}, \qquad (4)$$

Where SlowStart is the value of the timer and $E_{RTT}$ is the estimated round trip time of the packet. Many authors propose also formulae Eq. 5.

$$E_{RTT} = aE_{RTT} + bS_{RTT} \quad, \quad \text{where} \quad 0,8 \le a \le 0,9 \quad,$$
$$0,1 \le b \le 0,2, \ a+b=1 \qquad (5)$$

$S_{RTT}$ is a slow start round trip time. This formula enables small adjustment of the timer and more precise calculation of the waiting time limits in queues.

Therefore we propose that end-to-end delay limit calculation to use Eq. 6, where $N_{hops}$ is the number of hops in the end-to-end connection.

$$W_{\max LAN} \le \frac{SlowStart - 100ms}{N_{hops}} \qquad (6)$$

Typical number of hops $N_{hops}$ is up to 15. In case of 100 milliseconds of segment/fragmentation delay than the limit for waiting time in the queues is divided between hops. This limit is doubled or increased in different ways.

Because of the difference in service activity and distribution we apply in the simulation model the following bound for the peak traffic from all sources in Eq. 7. For three types of services we propose for more accuracy Eq. 8. The priority and delay requirements for email are low and they do not interference the overall behaviour in the interface.

$$R_{Peak} \le \sum_{i=1}^{3} R_{Pi} N_i \lambda_i \text{, where} \qquad (7)$$

$R_{peak}$ – total peak rate for all sources to the given interface

i – stands for VoIP, LAN and email parameters as follows:

$R_{PVoIP}$ – peak rate of 1 VoIP traffic source in packets

$\lambda_{VoIP}$ – VoIP source intensity per call

$N_{VoIP}$ – number of VoIP sources

$R_{PLAN}$ – peak rate of 1 LAN traffic source in packets

$\lambda_{LAN}$ – LAN source intensity per session

$N_{LAN}$ – number of LAN sources

$R_{Pemails}$ – peak rate of 1 email traffic source in packets

$\lambda_{email}$ – Email source intensity per session if any

$N_{email}$ – number of email sources

$$R_{Peak} \le \sum_{i=1}^{2} R_{Pi} N_i \lambda_i + 0,1 R_{Pemail} N_{email} \lambda_{email} \qquad (8)$$

The packet service time is estimated on the 100 Mbps interface rate to 0.00001732 seconds. This is the time interval for 200 bytes packet. The length of the queue fraction per service type $Q_{LenVoIP}$ is made equal to the series length $S_{VoIP}$, i.e. VoIP service. This is done to avoid series loss Eq. 9.

$$S_{VoIP} = Q_{LenVoIP}$$
$$S_{LAN} = Q_{LenLAN} \qquad (9)$$
$$S_{email} = Q_{Lenemail}$$

We also denote with $P_{VoIP}$, $P_{LAN}$ and $P_{email}$ payload per packet per service. Thus we derive the length of the series in Eq. 10.

$$S_{VoIP} = P_{VoIP} R_{PVoIP} \text{, packets} \qquad (10)$$

The overall queue length of the interface in packets is $Q_{Len}$ Eq. 11 or Eq. 12 depending on whether we allow series loss or not.

$$Q_{Len} \le S_{VoIP} + S_{LAN} + S_{Trans} \qquad (11)$$

$$Q_{Len} \ge S_{VoIP} + S_{LAN} + S_{Trans} \qquad (12)$$

Waiting time limits depends in the type of service. For example for VoIP end-to-end delay should be below 150 milliseconds. A typical number of hops are up to 25. Minimal delay on fragmentation at both ends is equal to the voice buffer, i.e. 20-30 milliseconds. Maximal waiting time per queue $W_{maxVoIP}$ can be calculated with Eq. 13.

$$W_{\max VoIP} \le Q_{LenVoIP} T_{Serv} \qquad (13)$$

Maximal waiting time limit for LAN packet $W_{\max LAN}$ depends on VoIP packets because they are of higher priority. The formula Eq. 14 is applied.

$$W_{\max LAN} \le n W_{\max VoIP} + Q_{LenLAN} T_{Serv}, \qquad (14)$$

Where n can be any number. For reasonable waiting time limits we choose n=2. The same rule is applied for maximal waiting time limit for the third queue fraction with the lowest priority $W_{\max email}$ shown in Eq. 15.

$$W_{\max email} \le n W_{\max VoIP} + m W_{\max LAN} + Q_{Lenemail} T_{Serv} \qquad (15)$$

Where n and m are any numbers but are chosen to be 2.

The number of parallel VoIP sessions $NS_{VoIP}$ can be calculated from Eq. 16, where AVoIP is traffic per VoIP source and $N_{VoIP}$ is the number of traffic sources.

$$NS_{VoIP} = N_{VoIP} A_{VoIP} \qquad (16)$$

The same is applied for other types of services.

## V. RESULTS

Simulation is performed on C++ language. The pseudo exponential pseudo deterministic characteristics of the traffic sources are reached after usage of combination between many random generators. The queue behaviour is complex due to the priorities and limits on waiting times. Many parameters have been derived from the model like probability of packet loss due to the lack of place in the queue, probability packet to be dropped due to the waiting limit exceed, probability to wait for different types of traffic, observations on of the packets intervals, queue lengths, delay, delay jitter, waiting times at many interface points in the model. Statistical accuracy of the derived results is proven by Student criterion. The presented results are in the 90% confidence interval from statistical point of view. IntServ, DiffServ, RSVP and NSIS have different way to gather with packets and this influences the way they police, drop and shape them.

Interesting results that influence directly interfaces and queue management are derived on the basis of queue length per service type. The queue fraction of the three services is observed. For services with highest priority like VoIP IntServ it is the most proper mechanism. DiffServ offers good overall utilization. RSVP and NSIS demonstrate the excellent quality for VoIP service. NSIS is the most flexible technique.

The results after the investigation of the priority queue with different waiting and place bounds are shown on Table II. Under almost the same utilization factor the utilisation of the fractions of the queue per service is changeable. Waiting losses are quite small for the fast interfaces and can be considered negligible. Table III represents probabilities of place losses on different utilization factor.

### TABLE II
### NUMERICAL RESULTS ON UTILIZATION

| Parameter | IntServ | DiffServ | RSVP | NSIS |
|---|---|---|---|---|
| Utilization | 0.4732 | 0.45847 | 0.44051 | 0.46791 |
| VoIP Utilization | 0.04533 | 0.04367 | 0.04238 | 0.04475 |
| LAN-to-LAN Utilization | 0.42334 | 0.40996 | 0.3936 | 0.41855 |
| Email Utilization | 0.00453 | 0.00484 | 0.00452 | 0.00461 |

### TABLE III
### NUMERICAL RESULTS ON PROBABILITY OF PLACE LOSSES

| | Medium traffic | Above medium traffic | Close to heavy Traffic | Heavy traffic |
|---|---|---|---|---|
| Utilization | 0.46184 | 0.48059 | 0.64364 | 0.71751 |
| Probability of Place Loss | 0.00001 | 0.00059 | 0.02651 | 0.08808 |
| Traffic sources | 0.00008 | 0 | 0.00011 | 0.00564 |
| Priorities | 0 | 0.00067 | 0.03087 | 0.09781 |
| | 0 | 0 | 0.00173 | 0.13891 |

The overall shaping effect is seen for the service with highest priority. The delay for the service with lowest priority becomes bigger. Therefore, the delay bound for real time services is kept on the favour of the non real time services.

## VI. CONCLUSION

The four QoS techniques distribute the queue resource in a different way. This is the reason to see different mean values and shaping effects. The acceleration effect of the services with highest priority is demonstrated.

The deterministic nature of the packets streams suppress shaping and increase losses. The statistical multiplexing effect is very limited due to the deterministic streams.

NSIS is the most flexible and tunable resource management and utilization technique. The authors refine the simulation model with more traffic sources and more precise generation of the packets from these sources based on the observation of the real traffic. NSIS protocol as well as non real time services has to be simulated as pure TCP traffic. Limits criteria for queue management and especially its tune adjustments criteria are under evaluation.

## REFERENCES

[1] S. Jha, M. Hassan, "Engineering Internet QoS", Artech Hose, 2002.
[2] T. Janevski, "Traffic Analysis and Design of Wireless IP Networks", Artech House, 2003.
[3] R. Goleva, M. Goleva, D. Atamian, T. Nikolov, K. Golev, "Quality of Service System Approximation in IP Networks", Serdica Journal of Computing, 2008, pp. 101-112
[4] J. Pitts, J. Schormans, "Introduction to IP and ATM Design and Performance", John Wiley&Sons, Ltd., 2000.
[5] V. Ralsanen, "Implementing Service Quality in IP Networks", John Wiley & Sons, Ltd., 2003.
[6] RFC 3742, "Limited Slow-Start for TCP with Large Congestion Windows", S. Floyd, 2004.
[7] RFC 2581, "TCP Congestion Control", M. Allman, 2001.
[8] L. Kleinrock, "Queueing Systems", Volumes I and II, John Wiley and Sons, 1976.
[9] V. Iversen, "Teletraffic Engineering Handbook", ITU-D, 2005.
[10] S. Mirtchev, "Study of Queueing Behavior in IP Buffers", iTech'07 Conference, Varna, St. Kostantin and Elena, June, 2007, pp. 187-193.
[11] B. Tsankov, R. Pachamanov, K. Kasev, "Traffic Management in IP Networks with Priorities", Telecom 2004, Varna, Bulgaria, pp. 283-289.
[12] L.P.Khadjiivanov, B.T.Taskov, A.A.Aliazidi, B.P.Tsankov, "Application of Priority Queueing Mechanisms to ATM Multiplexing and Traffic Control" Proc. of Integrated Broadband Communications Networks and Services, Copenhagen, Denmark, April 20 – 23, 1993, pp. 33.3.1 – 33.3.11.
[13] Fu, X., H. Schulzrinne, A. Bader, D. Hogrefe, C. Kappler, G. Karagiannis, H. Tschofenig, S. Van den Bosch, "NSIS: A New Extensible IP Signaling Protocol Suite", IEEE Communications Magazine, Oct., 2005, pp. 133-141
[14] "Next Steps in Signaling: Framework", RFC 4080, H. Tschofenig, D. Kroeselberg, June 2005.