

A New Framework and Burst Assembly for IP DiffServ over Optical Burst Switching Networks

Keping Long *, Rodney S. Tucker **, Chonggang Wang ***

* Special Research Centre for Optical Internet & Wireless Information Networks (COIWIN), School of Communication and Information Engineering, ChongQing University of Posts & Telecommunications (CQUPT), ChongQing 400065, P.R. of China, Email: longkp@cqupt.edu.cn

** ARC Special Research Centre for Ultra Broadband Information Networks (CUBIN), Dept of Electrical and Electronic Engineering, University of Melbourne, Victoria 3010, Australia, Email: k.long@ee.mu.oz.au , r.tucker@ee.mu.oz.au

*** HongKong University of Science and Technology, cgwang@cs.ust.hk

Abstract—IP Differentiated Services (DiffServ) has been standardized by the IETF and is considered as a promising IP QoS solution due to its scalability and ease of implementation. In this paper, we present a novel framework for IP Differentiated Services (DiffServ) over optical burst switching (OBS), namely, DS-OBS. We present the network architecture, functional model of edge nodes and core nodes, the control packet format, a novel burst assembly scheme at ingress nodes and scheduling algorithm of core nodes. The basic idea is to apply DiffServ capable burst assembly at ingress nodes and perform different per hop behavior (PHB) electronic treatment for control packets of different QoS classes service at core nodes. Simulation results show that the proposed schemes can provide the best differentiated service for expedited forwarding (EF), Assured forwarding (AF) and best effort (BE) service in terms of end-to-end delay, throughput and IP packet loss probability.

Index Terms—IP Quality of Service, differentiated services, per hop behavior, burst assembly, optical burst switching, DS-OBS

I. INTRODUCTION

IP over WDM has been envisioned as the dominant network architecture for the next generation optical Internet since it eliminates the intermediate layers and can make better use of advanced optical technologies [1, 2]. Currently, optical burst switching (OBS) is under study as a promising switching paradigm for the optical backbone network in IP over WDM networks [3-5]. The basic idea of optical burst switching (OBS) is to decouple the data channel from the control channel and to combine the best of the coarse-grained optical circuit switching and the fine-grained optical packet switching in order to facilitate the efficient integration of IP and WDM.

Due to the increasing deployment of Internet applications requiring QoS, such as voice and video application, a challenge issue for IP over WDM networks is how to support IP quality of service (QoS) at the WDM layer. Among the different IP QoS solutions, the Differentiated Services (DiffServ) architecture [6] is currently considered as a promising IP QoS solution due to its scalability and ease of implementation in electronic switching and routing techniques. An end-to-end differentiated service is

achieved by concatenation of per-domain services and service level agreements (SLAs) between adjoining domains along the route that the traffic crosses from the source to the destination. The per-domain service is realized by providing qualitative per hop behavior (PHB) using active queue management and scheduling algorithms at the core router for packets with different marking at DiffServ codepoints (DSCP) marked by the ingress router. Two PHBs are also standardized by the IETF, there are: Expedited Forwarding Per Hop Behavior (EF PHB) [7] to implement services requiring low loss and low delay ; and Assured Forwarding Per hop behavior (AF PHBs) [8] to implement services requiring assured bandwidth. However, these buffer-based schemes are difficult to apply to the WDM layer because optical RAM is not yet available even though a limited delay can be provided to optical packets by fiber delay lines (FDLs).

Recently, an offset-time-based QoS scheme for optical burst switching networks has been proposed by Qiao [9] to provide basic QoS by only considering the metric of burst blocking probability (BBP). Using this scheme, Qiao has analyzed the degree of class isolation using the offset times, upper/lower bounds on burst blocking probability, and QoS performance in terms of burst blocking probability and queuing delay [9]. Dolzer [10] has extended model of the offset-time-based QoS scheme by considering the interaction of offered load between the higher priority traffic and lower priority traffic. Their analysis and simulation results show that the offset-time-based QoS scheme can achieve lower burst blocking probability for higher priority traffic while degrading the performance of lower priority traffic.

However, there are several key issues are unclear for the offset-time-based QoS scheme. First, there have not been any reported studies on how the offset-time-based OBS QoS scheme supports the DiffServ or IntServ standards of the IETF. Secondly, the end-to-end (at least edge-to-edge) IP Packet QoS performance for the offset-time-based QoS scheme does not appear to have been analyzed or evaluated. The QoS performance of bursts is different to the QoS performance of IP packets. Chen [11] argues that the end-to-end delay for data

bursts will not represent the end-to-end delay of IP packets since packets of a real-time application may be separated by different data bursts at the burst assembly stage. Third, all existing studies assume that no contention occurs between control packets (reservation requests). As discussed in [12], this assumption is incorrect when multiple control packets (reservation requests) from different edge routers arrive at a core node synchronously, especially in the case of a fixed timer-based burst assembly algorithm and constant offset-time setup.

In this paper, we address the above key issues and proposed a novel framework for IP DiffServ over OBS networks, namely, DS-OBS. This paper is organized as follows. In section II, we define a new format for burst control packets and present a DiffServ capable OBS network architecture and functional model of edge/core nodes. In section III, we propose a novel DiffServ-capable burst assembly scheme at ingress nodes. Extensive simulation results in terms of end-to-end delay, throughput and IP packet loss probability for proposed architecture and schemes are given in section IV. Section V concludes this paper.

II. NETWORK ARCHITECTURE AND NODE MODEL FOR IP DIFFSERV OVER OBS

A. Network Architecture for IP DiffServ over OBS

The proposed network architecture and the basic function of edge routers and core routers for IP DiffServ over OBS (DS-OBS) are shown in Fig. 1, which consists of ingress nodes, core nodes and egress nodes. The basic idea of the proposed architecture is to provide IP DiffServ over OBS networks (DS-OBS) while avoiding the deployment of extra offset and optical buffers. In contrast to the traditional OBS network architecture. In our proposed architecture, when traditional IP packets enter the ingress router that integrates the function of a DiffServ edge router with the function of an OBS edge router, the functions of both traffic conditioning and burst assembly are performed for the IP packets. Then, a data burst and a proposed control packet carrying the DSCP information as well as other information are generated and transmitted into the core of the OBS network separately on different wavelengths with the Burst Head Packets (BHPs) sent ahead by an offset time. At core nodes, BHPs are processed electronically to provide different per-hop behaviors (PHBs) for different classes of service. As a result of this processing, the differentiated service for corresponding data bursts was implemented. At the egress node of a domain, the data bursts are disassembled back into IP packets to be forwarded to their next hops (e.g. traditional IP routers).

B. Functional model of ingress Nodes and control packet

The proposed functional model of an ingress router and the control packet format for the DS-OBS are shown in Fig.2. After being marked as EF, AF and BE service, the IP packets are enqueued into separate queues. The burst assembler performs

DiffServ-capable assembly for EF, AF and BE service by adaptively adjusting the assembly parameters (timer and burst size) based on the QoS requirements of different classes of service and the metering results of actual traffic arrivals. A detailed burst assembly algorithm will be given in section III.

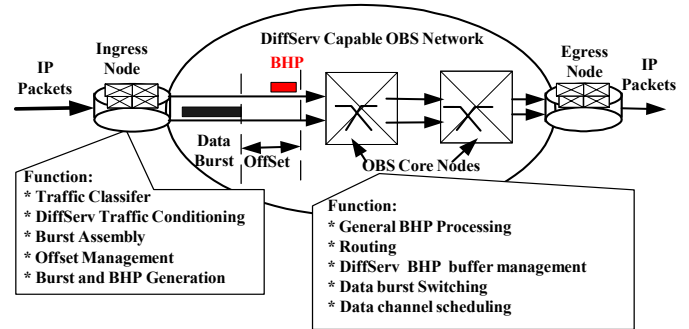


Fig.1 Network architecture for IP DiffServ over OBS

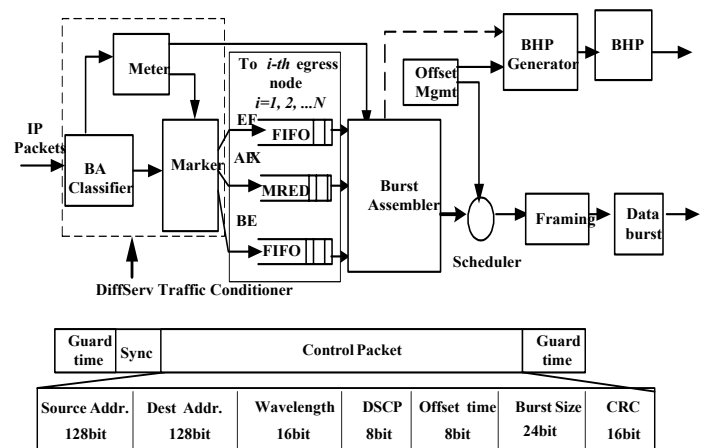


Fig.2 Model of ingress node and control packet format for DS-OBS Network

Once a burst of EF, AF or BE service has been created, based on the information of burst assembly and offset management, control packet generator will generate a burst header packet (BHP), which is defined as Fig.2. This includes information such as DiffServ codepoints (DSCP), offset time, source/destination address and burst size. From Fig.2, we can see the maximum length of a BHP is less than 45bytes and 128bits of Source/Destination address is reserved for the future IPv6 protocol. So far, no standardization of OBS control packets has been reported in the existing literatures. We define a new BHP not only for our scheme but also for being extended to other OBS networks, e.g. if (G)MPLS is deployed for traffic engineering in the OBS network, we can only change the of source/destination address fields to a "Label" field with 20bits length.

In this paper, both the constant offset and random offset settings are considered. As discussed in [12], in order to prevent excessive loss due to the synchronized arrivals of data bursts at the same output of OBS core nodes, a random offset generation (e.g. Possion token bucket) is deployed for data burst shaping.

C. Functional model and scheduling algorithm of core nodes

The control part of a core router is shown in Fig. 3, which keeps the optical data processing as simple as possible. The priority treatment of data bursts is performed by the different per-hop behavior treatment for their BHPs of EF, AF and BE services. This is achieved by priority scheduling for BHPs and appropriately configuring the transmission rate of control channel as follows.

$$Num_{hop} \times Num_{dataChannel} \times \frac{Avg_{BHPSize}}{Rate_{controlChannel}} = \frac{Avg_{BurstSize}}{Rate_{DataChannel}} \quad (1)$$

In contrast to existing schedulers in OBS core nodes, the scheduler presented in this paper as shown in Fig. 3 has two functions. First, it performs priority scheduling (or Weighted Round Robin scheduling, WRR) for multiple BHPs queues of EF, AF and BE service according to the information carried in the BHPs. As a result of differentiated service for BHPs of different service classes, the differentiated service for the corresponding data bursts was achieved. Secondly, it schedules data bursts on outgoing data channels using data channel scheduling algorithms, e.g. the *latest available unused channel with void filling* (LAUC-VF) algorithm [13] and our proposed the *Least Gap with Void Filling* (LGVF) algorithm [15]. Since the main differentiated treatment has already been conducted for the BHPs of different service classes, the performance enhancement of our LGVF over LAUC-VF in the DS-OBS architecture is not obvious as the enhancement in other OBS network [15].

The scheduler is bidirectionally connected to the switching configuration module. The aim is to exchange configuration information between them, to make the scheduler update the state information of control channel and data channels, to rewrite fields of BHPs and transmit the BHPs to the next hop.

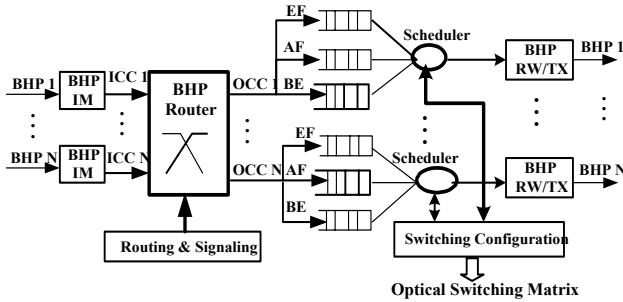


Fig.3 Control part of a core node for DS-OBS network

III. DIFFSERV-CAPABLE BURST ASSEMBLY SCHEME

Several burst assembly algorithms have been proposed in [13][16]. Xiong [13] proposed a bucket-based burst assembly scheme based on the egress router addresses, assembly time intervals and maximum data burst size. There are several open issues in this scheme, such as: how to dimension the parameters of buckets according to the traffic load and number of QoS classes and destinations, burst assembly time T_c and maximum burst size L . To achieve a better trade-off between the burst assembly delay and burst assembly efficiency when input traffic

load is temporarily low, An Ge [16] presented a simple burst assembly mechanism based on a time-window and minimum data burst size. However, both of these two mechanisms cannot support IP DiffServ and cannot adjust the parameters of burst assembly according to the QoS characteristic of different service classes.

In this section, we propose a novel DiffServ-based burst assembly scheme. The proposed scheme can support the DiffServ services (EF PHB, AF PHB and BE) standardized by the IETF, and can adjust the parameters adaptively according to the actual traffic arrival rate (output of meter as shown in Fig.2) and the differentiated QoS requirements of different service classes. The detailed burst assembly algorithm is given as follows.

Definition:

T_A^{EF}, T_A^{AF} : Maximum burst assembly time of EF and AF service, configured by the delay requirement of EF and AF.

T_A : Allowable maximum burst assembly time in an OBS network. Dimensioned based on the number of destinations, the number of QoS classes and the end-to-end delay.

MAX_{BS}, MIN_{BS} : Allowable maximum/minimum burst size in an OBS network. Decided by the allowable burst assembly delay, burst assembly efficiency and link utilization [12, 13].

$BS_{min}^{EF}, BS_{min}^{AF}, BS_{min}^{BE}$: Allowable minimum burst size of EF, AF and BE service. BS_{min}^{AF} is configured to be more than MIN_{BS} and is decided by the Committed Information Rate (CIR) in the metering using trTCM [14]. BS_{min}^{BE} is configured to be equal to MAX_{BS} .

$BS_{max}^{EF}, BS_{max}^{AF}, BS_{max}^{BE}$: Allowable maximum burst size of EF, AF and BE service. Considering the delay characteristic of EF PHB service, set the $BS_{min}^{EF} = BS_{max}^{EF} = MIN_{BS}$. BS_{max}^{AF} is less than MAX_{BS} and is decided by the Peak Information Rate (PIR) in the metering using trTCM [14].

$T_k^{EF}, T_k^{AF}, T_k^{BE}$ is the timer for assembling EF, AF, BE, respectively.

$L_{EF}(k), L_{AF}(k), L_{BE}(k)$ is the assembly meter for computing the assembled burst length of EF, AF, BE, respectively.

Algorithm:

1) When a packet with length of pkt_size arrives at a queue associated with a destination at the ingress node of the OBS domain.

If it is marked as EF PHB and $L_{EF}(k)=0$, start timer T_k^{EF} and $L_{EF}(k) = L_{EF}(k) + pkt_size$

Else if it is marked as AF PHB and $L_{AF}(k)=0$, start timer T_k^{AF} and $L_{AF}(k) = L_{AF}(k) + pkt_size$

Else it is marked as BE and $L_{BE}(k)=0$, start timer T_k^{BE} and $L_{BE}(k)=L_{BE}(k)+\text{pkt_size}$

2) If $\{T_k^{EF} > T_A^{EF}$ or $L_{EF}(k) > BS_{\min}^{EF} = MIN_{BS}\}$, the EF burst is created and reported with the length of burst $L_{EF}(k)$ or MIN_{BS} if $L_{EF}(k)$ is less than MIN_{BS} when the timer T_k^{EF} reaches the T_A^{EF} , reset $T_k^{EF}=0$ and $L_{EF}(k)=0$.

Else if $\{T_k^{AF} > T_A^{AF}$ or $BS_{\min}^{AF} < L_{AF}(k) < BS_{\max}^{AF}\}$, the AF burst is created and is reported with the length of burst $L_{AF}(k)$ or BS_{\min}^{AF} if $L_{AF}(k)$ is less than BS_{\min}^{AF} when the timer T_k^{AF} reaches T_A^{AF} , reset $T_k^{AF}=0$ and $L_{AF}(k)=0$.

Else if $\{T_k^{BE} > T_A$ or $L_{BE}(k) \geq MAX_{BS}\}$, the BE burst is created and reported with the length of burst $L_{BE}(k)$ or MAX_{BS} if $L_{BE}(k)$ is less than MAX_{BS} when the timer T_k^{BE} reaches the T_A , reset $T_k^{BE}=0$ and $L_{BE}(k)=0$.

(3) Otherwise, keep all the values unchanged and waiting for the next IP packet arrival and go to (1).

IV. PERFORMANCE EVALUATION

Using OPNET, we have developed simulation models of our proposed architecture, protocol, burst assembly scheme at ingress nodes and scheduling algorithm at core nodes. We have simulated the end-to-end IP performance and IP packet loss rate at core nodes using the simulation network shown in Fig.4. Each of four Possion traffic sources (one EF, one AF and two BE) with the same traffic arrival rate and same average packet length (500bytes) was accessed by each of four ingress routers via 2.5Gbit/s wavelengths and were multiplexed into core router 0. The multiplexed stream was transmitted to each of four egress nodes via core routers over 4-hops with $4 \times 2.5\text{Gbit/s}$ data wavelengths. The BHPs are transmitted on a separate channel in the core of network.

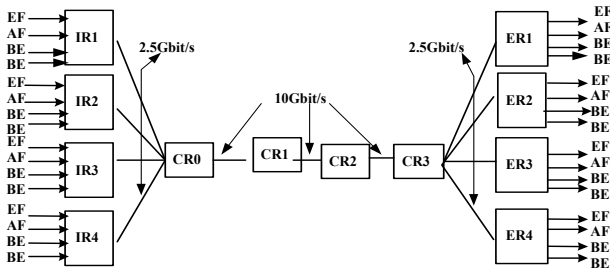


Fig.4 Simulation network topology

Three scenarios are simulated to evaluate the performance of our schemes: Scenario 1, constant offset at ingress nodes, FIFO for BHPs at core nodes; Scenario 2, constant offset at ingress nodes, WRR at core nodes; Scenario 3, random offset at ingress nodes, WRR at Core nodes. The parameters of the burst assembler are shown in Table 1.

TABLE 1 PARAMETERS OF BURST ASSEMBLER

Service classes	Min burst size (Bytes)	Max burst size (Bytes)	Timer (ms)
EF	5k	5k	4.8
AF	30k	50k	55
BE	125k	125k	600

Fig.5 shows the IP packet loss rate at core0 as a function of the source traffic load, which was caused by both the congestion of control packets and the scheduling failure of data bursts. The IP packet loss rate in all three scenarios is low due to the proposed burst assembly scheme at ingress nodes. As the load increases, the packet loss rate increases. However, The IP packet loss rate is acceptable even if the load is 1.0. Moreover, the best performance is achieved by using random offset setting at ingress nodes and WRR at core nodes. When a constant offset is applied, the performance of WRR scheduling is better than the performance of FIFO.

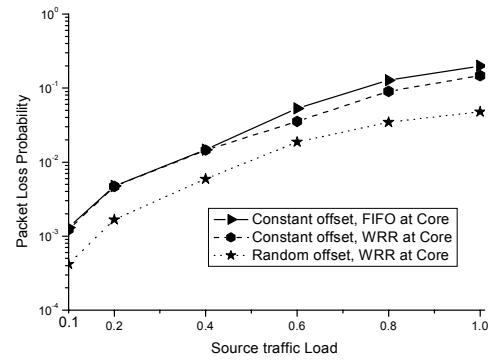


Fig.5 IP packet loss rate at core0 in the three scenarios

Fig. 6 shows the actual offered load at the input of core0 and the link utilization at the output of Core0 as a function of source traffic load. The link utilization of all three scenarios is high due to the low IP packet loss rate. The link utilization for the scenario with random offset at ingress nodes and WRR at core nodes is just slightly lower than the actual offered load at core0 due to the very low packet loss rate even if the source traffic load is 1.0. This is because DiffServ traffic conditioning (metering, marking, shaping, dropping) and admission control have been applied for bursty IP traffic to assure their performance and to avoid the waste of resources at the core of the DS-OBS network. So, the actual offered load at input of Core0 is less than source traffic load.

Furthermore, in order to illustrate the service differentiation of our proposed architecture, the packet loss rate for EF, AF and BE service under all three scenarios when the traffic load is 1.0 is shown in Fig.7. The scenario with random offset at ingress nodes and WRR at core nodes can achieve the best service differentiation and achieve the lowest packet loss rate for AF and EF as well as the lowest total packet loss rate. Scenario 2 (constant offset at ingress nodes, WRR at core nodes) has achieved better service differentiation and total packet loss rate

than scenario 1 (constant offset at ingress nodes, FIFO at core nodes).

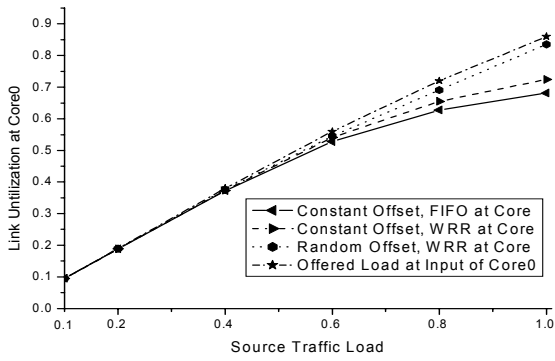


Fig. 6 Link Utilization at Core0

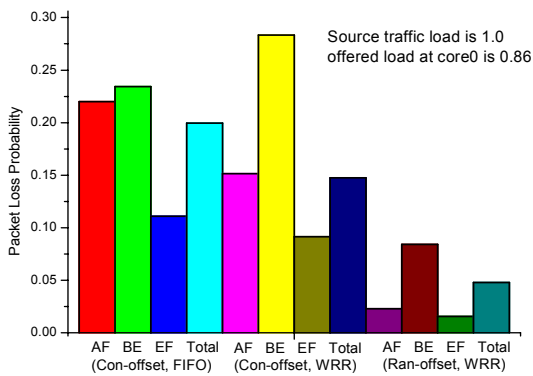


Fig.7 The packet loss rate for EF, AF, BE in three scenarios (load 1.0)

The actual achieved end-to-end (E2E) throughput normalized by reserved bandwidth for EF, AF and BE service as a function of source traffic load is given in Fig.8. When the load is below 0.4, three types of service can achieve over 0.92 times the expected bandwidth and there is no obvious difference among the three service classes. As load increases from 0.6, the service differentiation becomes obvious following the idea of DiffServ architecture for all three scenarios: AF service achieved the most assured bandwidth and EF service achieved as much bandwidth as they can due to the strict delay requirement. The bandwidth assurance of AF and EF service is achieved by decreasing the end-to-end throughput of BE service. The performance of setting random offset at ingress nodes and WRR at core nodes is the best compared to the other two scenarios. When the constant offset is applied, the performance of WRR scheduling is better than FIFO. Results show that our DS-OBS is very efficient to support the IP DiffServ architecture.

For each of AF, EF and BE service, the end-to-end (E2E) delay under the three scenarios is almost the same and is shown in Fig.9. Due to our DiffServ-capable burst assembly and since no extra offset is deployed for QoS, E2E Delay of EF service is less than 6.67ms even if the load is 0.1. So, the possible worst

E2E delay of EF service in 100-hops is about 160ms and meets the delay demand of voice and video service. However, the E2E delay of AF is tens of milliseconds and E2E delay of Best effort (BE) is hundreds of milliseconds

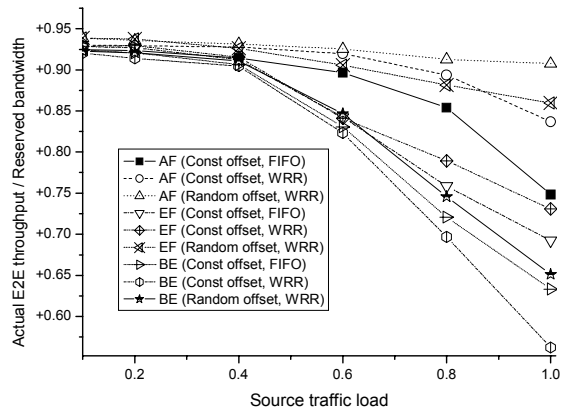


Fig. 8 Achieved end-to-end throughput normalized by reserved bandwidth

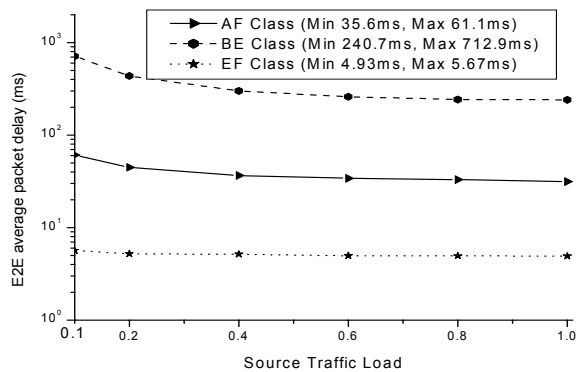


Fig.9 End-to-end delay for EF, AF and BE for 4-hops network

V. CONCLUSION

In this paper, we propose a novel framework for IP DiffServ over OBS, namely DS-OBS, including the network architecture, novel functional model of ingress nodes and core nodes, a control packet structure, a DiffServ-capable burst assembly scheme at ingress nodes, and scheduling scheme at core nodes. Using a 4-hop network, we have simulated IP packet loss probability at core nodes, and the end-to-end (E2E) throughput and end-to-end (E2E) delay for EF, AF and BE service. The results show that the proposed schemes can not only provide the best service differentiation for EF, AF and BE in terms of both the E2E performance and IP packet loss probability at core nodes, but also enhance the whole performance of network.

REFERENCES

[1] P.Bonenfant, A.Rodriguez-Moral and J. Manchester, "IP over WDM: the missing Link", White Paper, Lucent Technologies, 1999.

- [2] N. Gnani, S. Dixit and T. Wang, "On IP-over-WDM integration", IEEE Communication Magazine, pp.72-84, March 2000.
- [3] C. Qiao and M. Yoo, "Optical burst switching (OBS) - a new paradigm for an optical internet", J. High Speed Networks, vol. 8, pp. 69-84. Jan. 1999.
- [4] C. Qiao, "Labeled Optical Burst Switching for IP-over-WDM Integration", IEEE Communication Magazine, Sept. 2000.
- [5] J. Turner, "Terabit burst switching", J. High Speed Networks, vol. 8, pp. 3-16, 1999.
- [6] S. Blake, D. Black, M. Carlson, "An Architecture for Differentiated Services", IETF RFC 2475, Oct. 1998.
- [7] V. Jacobson, K. Nichols and K. Poduri, "An Expedited Forwarding PHB", IETF RFC 2598, June 1999
- [8] J. Heinanen, F. Baker, W. Weiss and J. Wroclawski, "Assured Forwarding PHB group", IETF RFC 2597, June 1999
- [9] M. Yoo, C. Qiao and Sudhir Dixit, "QoS performance of optical burst Switching in IP-over-WDM networks", IEEE J. Select. Areas Commun., vol. 18, No. 10, pp. 2062-2070, Oct. 2000.
- [10] K. Dolzer, C. Gauger, J. Spith and S. Bodamer, "Evaluation of reservation mechanisms for optical burst switching", AEU International Journal of Electronics and Communications, vol. 55, No. 1, Jan. 2001.
- [11] Y. Chen, M. Hamdi, D. Tsang and C. Qiao, "Proportional QoS over OBS networks", Proc. IEEE Global Telecommunications Conference (Globecom' 01), San Antonio, Nov 2001.
- [12] H. M. Chaskar, S. Verma and R. Ravikanth, "Optical burst switching: a viable solution for terabit IP backbone", IEEE Network, pp.48-53, Nov./Dec. 2000.
- [13] Yijun Xiong and Marc Vandenhoute, "Control Architecture in Optical Burst-Switched WDM Networks", IEEE J. Select. Areas Commun., vol. 18, No. 10, pp. 1838-1851, Oct. 2000.
- [14] J. Heinanen and R. Guerin, "A Two Rate Three Color Marker", IETF RFC 2698, Sept. 1999.
- [15] Keping Long, R. S. Tucker and Seen-Yoon OH, "Fairness scheduling algorithms for supporting QoS in optical burst switching networks", Proceedings of the SPIE's International Symposium: Asia-Pacific Optical and Wireless Communications (APOC), Shanghai, CHINA, 2002.
- [16] An Ge, Franco Callegati and Lakshman S. Tamil, "On optical burst switching and self-similar traffic", IEEE Communication Letters, vol. 4, No. 3, pp. 98-100, March 2000.