

# PRNN/ERLS-based predictive QoS-promoted DBA scheme for upstream transmission in EPON

Jan-Wen Peng · Chung-Ju Chang · Po-Lung Tien

Received: 2 November 2008 / Accepted: 25 January 2010 / Published online: 11 February 2010  
© Springer Science+Business Media, LLC 2010

**Abstract** This article proposes a PRNN/ERLS-based predictive QoS-promoted dynamic bandwidth allocation (PQ-DBA) scheme for upstream transmission in Ethernet passive optical network (EPON) systems. The proposed PQ-DBA scheme originally divides incoming packets of voice, video, data service traffic into six priorities, where packets having less room before QoS requirements violation or being in starvation situation will be dynamically promoted to high priority cycle-by-cycle. It predicts packets arriving at prediction interval for ONUs using pipeline recurrent neural network (PRNN)/extended recursive least squares (ERLS) so that the bandwidth allocation can be more up-to-date and then accurate. Simulation results show that the proposed PQ-DBA scheme achieves higher system utilization and lower average voice, video, data packet delay time than the DBAM scheme [Luo and Ansari, *OSA J Opt Netw* 4(9):561–572] by 4, and 21, 90, 43%, respectively, and the PQ-DBA scheme but without prediction by 2, and 26, 29, 34%, respectively.

**Keywords** Ethernet passive optical network (EPON) · Pipelinerecurrent neural network (PRNN) · Extended recursive least squares (ERLS) · Predictive QoS-promoted dynamic bandwidth allocation (PQ-DBA) · Optical line terminal (OLT) · Optical network unit (ONU)

## 1 Introduction

The Ethernet passive optical network (EPON), which represents the combinations of low-cost Ethernet equipment and low-cost fiber infrastructure, is considered to be the cost-effective solution of the optical access network [1]. The EPON system has to support triple-play services in optical access networks, and must fulfill the quality-of-services (QoS) requirements of each service. Kramer et al. [2] proposed an IPACT upstream multiple access mechanism for EPON. This work dynamically assigned bandwidth to all ONUs by polling ONUs' demands, and was selected as IEEE 802.3ah standard [3]. However, quality of service (QoS) requirements, such as voice/video packet delay and packet dropping probability, were not considered.

Several polling-based upstream scheduling algorithms were proposed to deal with the QoS problem in [4–7]. Cheng et al. [4] proposed a DBA-high priority (DBA-HP) scheme, which focused on the high-priority traffic. The DBA-HP minimized the packet delay time and delay variation of high-priority packets, but sacrificed the packet delay, packet dropping probability, and throughput of low-priority packets. An intra-ONU priority scheduling scheme [5] and a two layer bandwidth allocation (TLBA) scheme [6] were proposed to resolve the unfairness for the low priority packets by designing a maximum cycle time to each traffic class. However, it produced increment of delay time of high-priority traffic and decrement of system throughput, due to that the available bandwidth cannot meet all demands resulted by the burst or the heavy traffic load. A traffic-class burst-polling-based delta dynamic bandwidth allocation (TCBP-DDBA) scheme was studied in [7]. The scheme not only reallocated extra bandwidth to the heavily loaded ONUs to improve under utilization problems but also provided QoS guarantee to delay sensitive services. However, it did not individually allocate

---

J.-W. Peng (✉)  
Chunghwa Telecom Laboratories, Taipei 100, Taiwan, ROC  
e-mail: janwen@cht.com.tw

C.-J. Chang · P.-L. Tien · J.-W. Peng  
National Chiao Tung University, Hsinchu 300, Taiwan, ROC  
e-mail: cjchang@mail.nctu.edu.tw

P.-L. Tien  
e-mail: tbl@cm.nctu.edu.tw

bandwidth to each class and resulted in the longest delay to non-delay sensitive services if the traffic load was heavy and the ONU did not arrange the transmission well.

Moreover, Wu et al. [8] proposed a prediction-based longest queue first (PLQF) scheduling algorithm, which is based on the information of user queue length and the incoming traffic. The system resource is allocated to the user who has the largest probability to overflow in the near future to achieve a minimal cell loss rate. Luo and Ansari proposed a dynamic bandwidth allocation with multiple services (DBAM) scheme in [9] and a limited sharing with traffic prediction (LSTP) scheme in [10]. The DBAM scheme adopted a limited bandwidth allocation (LBA) and employed a linear *estimation credit* for class-based traffic prediction to estimate queue's traffic arriving during cycle time. The ONUs apply *priority queueing* to buffer the three types of service frames, and a maximum time slot length of a specific class of traffic in each ONU is pre-assigned to a service level agreement (SLA) between the end user and the service provider. The DBAM can reduce the packet delay and queue length. The LSTP scheme used the same upper bounded maximum time slot length in bytes of ONUs as the DBAM scheme, but with a linear predictor and a least-mean square (LMS) adaptive update algorithm in ONUs to predict the data traffic arrived during the cycle time. An early-DBA mechanism with prediction-based fair excessive bandwidth reallocation (PFEBR) scheme was introduced in [11], where the modified *line estimation credit* from DBAM [9] was used to ensure fairness of all ONUs. Yin et al. [12] proposed a nonlinear prediction-based dynamic bandwidth allocation (NLPDBA) scheme to improve the upstream transmission efficiency and prediction accuracy. The NLPDBA scheme used the same LMS update algorithm as the LSTP scheme [10] and controlled the estimated result for stability. However, when the predictor underestimates or overestimates the traffic, it may drop packets or waste bandwidth and cannot show how soon it will converge. Thus, the NLPDBA scheme is hard to decide the precise granted bandwidth for ONUs in the next cycle.

In this article, we propose a PRNN/ERLS-based predictive QoS-promoted dynamic bandwidth allocation (PQ-DBA) scheme for upstream transmission in EPON. The PQ-DBA not only commits the QoS requirements for real-time services, but also promotes the fairness for non-real-time packets, especially when the traffic intensity is larger than 0.8. The PQ-DBA originally divides voice, video, and data types of service traffic into six priorities. The proposed PQ-DBA has different treatments to packet's transmission. Priority of real-time video packets is raised if the packets will violate the QoS requirements at the beginning of the next cycle. Similarly, priority of non-real-time data packets is promoted if the packets suffer a long delay and will get into a starvation situation. Besides, a pipeline recurrent neural network (PRNN)/extended recursive least square (ERLS)

update predictor [13] is adopted to precisely predict arrival packets at ONU during the next prediction interval cycle-by-cycle. The predictor has good nonlinear prediction capability and fast convergent time [14], and we had successfully employed nonlinear pipelined RNN to predict the interference variation of DS-CDMA/PRMA system [15]. Simulation results show that proposed PQ-DBA improves the system utilization, the average voice, video, and data delay time by 4, and 21, 90, 43%, respectively, over the DBAM [9]. The PQ-DBA scheme also achieves lower average voice, video, and data delay time by 26, 29, and 34%, respectively, and higher system utilization by 2% than the PQ-DBA without prediction. Besides, the proposed PQ-DBA can fulfill the video packet dropping probability requirement whereas the DBAM scheme fails.

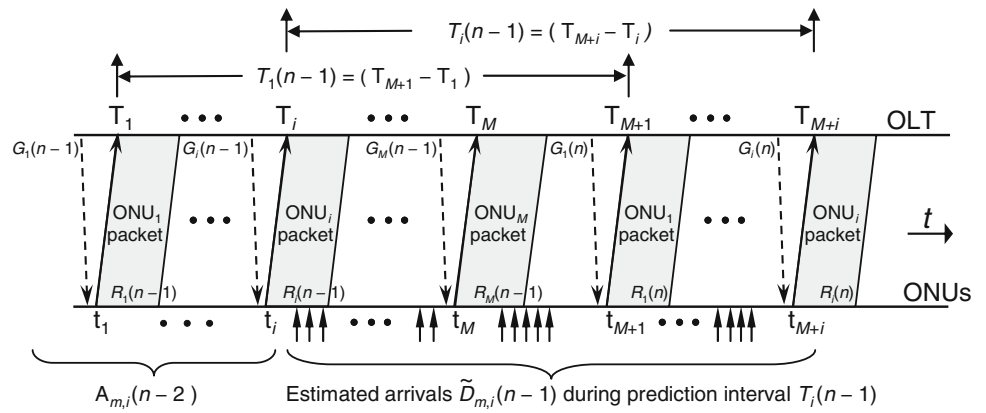
The rest of the article is organized as follows. Section 2 describes the system model. Section 3 introduces the PQ-DBA scheme, includes the PRNN/ERLS predictor, the QoS-promoted operation and the PQ-DBA assignment procedure. Simulation results and discussions are presented in Sect. 4. Finally, concluding remarks are given in Sect. 5.

## 2 System model

Assume that there is one optical line terminal (OLT) and  $M$  optical network units (ONUs) split by 1:  $M$  splitter in an EPON system. The line rate is  $R_E$  bps between OLT and each ONU, and the line rate  $R_U$  bps between ONU and its own end users. Two wavelengths are used to serve downstream and upstream traffic individually. The system supports three types of services, real-time voice, real-time video, and non-real-time data. In  $ONU_i$ , three types of queues are provided to store real-time voice, real-time video, and non-real-time data packets, which are denoted by  $Q_{0,i}$ ,  $Q_{1,i}$ , and  $Q_{2,i}$ , respectively,  $1 \leq i \leq M$ . The incoming packets from users will be put into the corresponding queues at ONU according to their service types. The arriving packet will be dropped if its queue is full, and the packet will be discarded if its QoS requirement is violated.

There is a GATE (REPORT) message sent from OLT (ONU) to ONU (OLT) to manage the transmission between OLT and ONUs. Figure 1 shows the upstream transmission between the PQ-DBA at OLT and ONUs in an EPON system. Assume that the EPON system is at the present cycle  $(n - 1)$ , which starts from when OLT sends the GATE message for  $ONU_1$  at cycle  $(n - 1)$ , denoted by  $G_1(n - 1)$ , to when the OLT sends the  $G_1(n)$  for the next cycle  $n$ . The PQ-DBA assumes a prediction interval for  $ONU_i$  at cycle  $(n - 1)$ , denoted by  $T_i(n - 1)$ ,  $1 \leq i \leq M$ , which starts from when OLT receives the first packet from  $ONU_i$  with the  $(n - 1)$ -th REPORT message (at time  $T_i$  of Fig. 1) to when OLT receives the first packet of  $ONU_i$  with the  $n$ -th

**Fig. 1** Upstream transmission between OLT and ONUs in the EPON system



REPORT message (at time  $T_{M+i}$  of Fig. 1). The  $T_i(n-1)$  is not necessarily equal to the  $T_j(n-1)$ ,  $\forall n$ , when  $i \neq j$ , and  $1 \leq i, j \leq M$ . Assume that the PQ-DBA can know when the  $T_{M+i}$  will be.

The GATE message for  $ONU_i$  at cycle  $(n-1)$ ,  $G_i(n-1)$ , is a set of  $\{G_{0,i}(n-1), G_{1,i}(n-1), G_{2,i}(n-1)\}$ , where  $G_{m,i}(n-1)$  denotes the granted bandwidth for the type  $m$  of service traffic in  $ONU_i$ ,  $m \in \{0, 1, 2\}$ ,  $1 \leq i \leq M$ . When the ONU receives the GATE message from the OLT, it transmits packets at its assigned timeslot and piggybacks a REPORT message at the end of the packet. The REPORT message from  $ONU_i$  to OLT at cycle  $(n-1)$ , denoted by  $R_i(n-1)$ , is a set of  $\{L_{0,i}(n-1), L_{1,i}(n-1), L_{2,i}(n-1), L_{dp,i}(n-1), L_{d,i}(n-1), L_{w,i}(n-1)\}$ . The  $L_{m,i}(n-1)$  is the occupancy of queue  $Q_{m,i}$  at cycle  $(n-1)$ ,  $m \in \{0, 1, 2\}$ , and the  $L_{dp,i}(n-1)$ ,  $L_{d,i}(n-1)$ , and  $L_{w,i}(n-1)$  are numbers of bytes that had better be transmitted at the next cycle otherwise these packets will be dropped and/or the QoS requirement will be violated. The QoS requirements  $L_{dp,i}(n-1)$ ,  $L_{d,i}(n-1)$ , and  $L_{w,i}(n-1)$  are derived in the following.

The  $L_{dp,i}(n-1)$  indicates the total amount of bytes of real-time video packets that will violate the video packet delay requirement, denoted by  $T_d^*$ , if they are not transmitted at the next cycle. Denote  $T_{d,k}$  to be the delay time of the  $k$ -th video packet in  $Q_{1,i}$  of  $ONU_i$ ,  $1 \leq i \leq M$ , at the present cycle  $(n-1)$ , where  $k=1$  means the first packet in  $Q_{1,i}$ . Denote  $x$  to be the  $x$ -th packet with the least delay time, which will violate the delay requirement at the beginning of the next prediction interval, such that any video packet queued before  $x$ -th packet will be dropped. Then,  $x$  can be calculated by

$$x = \arg \min_k \{T_{d,k} + T_i(n-1), \forall n, k, i, \text{ and } T_{d,k} + T_i(n-1) > T_d^*\} \tag{1}$$

Then, the  $L_{dp,i}(n-1)$  can be obtained by

$$L_{dp,i}(n-1) = \sum_{k=1}^x S_{k,1}, \tag{2}$$

where  $S_{k,1}$  is the number of bytes of the  $k$ -th packet's size in  $Q_{1,i}$ ,  $1 \leq i \leq M$ .

The  $L_{d,i}(n-1)$  represents the total amount of bytes of real-time video packets, which should be transmitted at the next cycle, otherwise the video packet delay requirement will be violated, and the requirement of video packet dropping probability, denoted by  $P_d^*$ , cannot be kept. Denote  $P_{d,i}$  the dropping probability of video packets at  $Q_{1,i}$ , measured by  $ONU_i$ ,  $1 \leq i \leq M$ . A moving time window is adopted to calculate the video packet dropping probability. It contains the latest  $N$  output video packets of  $ONU_i$ , which have been dropped and transmitted, or are going to be dropped or transmitted at the next cycle. Assume that there are  $N_d$  video packets, among the  $N$  video packets, which have been dropped so far. Assumed there are  $x$  video packets waiting in the queue  $Q_{1,i}$ , and being dropped if they are not transmitted at the next cycle. The  $x$  is given in Eq. 1. Thus, a number of packets among these  $x$  packets, denoted it by  $y$ , must be transmitted otherwise the requirement of video packet dropping probability,  $P_d^*$ , will be violated. Then,  $y$  can be obtained by

$$y = (N_d + x - \lceil N \times P_d^* \rceil)^+, \tag{3}$$

where  $(a)^+ = a$  if  $a \geq 0$ ,  $(a)^+ = 0$  if  $a < 0$ ; and  $\lceil b \rceil$  denotes the smallest integer greater than  $b$ . Then, the  $L_{d,i}(n-1)$  can be derived by

$$L_{d,i}(n-1) = \sum_{k=1}^y S_{k,1}, \tag{4}$$

where  $S_{k,1}$  means the number of bytes of the  $k$ -th packet's size in  $Q_{1,i}$ ,  $1 \leq i \leq M$ .

The  $L_{w,i}(n-1)$  is the total amount of bytes of non-real-time data packets whose waiting time will be larger than a starvation-threshold time, denoted by  $T_w^*$ , at the next cycle. It tells OLT how much bandwidth is required for non-real-time data packets in  $Q_{2,i}$  to prevent starvation. Denote  $T_{w,k}$  to be the waiting time of the  $k$ -th data packet in  $Q_{2,i}$  of  $ONU_i$ ,  $1 \leq i \leq M$ , at the present cycle  $(n-1)$ . A number of packets with a waiting time larger than the starvation-threshold time,

$T_w^*$ , should be served at the next cycle. Denote the number of packets to be  $z$ , and the  $z$  is given by

$$z = \arg \min_k \{T_{w,k} - T_w^*, \forall k, i, \text{ and } T_{w,k} - T_w^* > 0\}. \quad (5)$$

Then, the  $L_{w,i}(n - 1)$  can be obtained by

$$L_{w,i}(n - 1) = \sum_{k=1}^z S_{k,2}, \quad (6)$$

where  $S_{k,2}$  means the number of bytes of the  $k$ -th packet's size in  $Q_{2,i}$ ,  $1 \leq i \leq M$ . The non-real-time data packets do not have strict delay criterion, but they ought to be protected from such an unfair condition by an inappropriate bandwidth assignment. Similar to the random early detection (RED) scheme [16], the starvation–threshold time  $T_w^*$  starts a mechanism to keep data packets from reaching a starvation situation.

### 3 PQ-DBA scheme

The PQ-DBA begins the prediction procedure for the  $i$ -th ONU when it receives the  $i$ -th ONU REPORT message at the present cycle  $(n - 1)$ ,  $R_i(n - 1)$ ,  $1 \leq i \leq M$ . It adopts a pipeline recurrent neural network/extended recursive least square (PRNN/ERLS) for prediction, which has good nonlinear prediction capability and fast convergent time. When the predictions for all ONUs have accomplished, the PQ-DBA performs the bandwidth allocation. The PQ-DBA scheme classifies the three types of services traffic from ONUs into six priorities and provides a QoS promotion capability at ONUs, where service packets, which will violate the QoS requirement will be promoted to higher priority. The bandwidth allocation is according to the REPORT messages and the predicted arrival packets of service traffic of all ONUs. Then, the PQ-DBA at OLT sends the  $i$ -th GATE message to  $ONU_i$  for the next cycle,  $G_i(n)$ ,  $1 \leq i \leq M$ .

#### 3.1 PRNN/ERLS predictor

For prediction of the arrival packets during the prediction interval of  $ONU_i$  at present cycle  $(n - 1)$ ,  $T_i(n - 1)$ ,  $1 \leq i \leq M$ , the PRNN/ERLS predictor would have to predict the packet arrival rate of service type  $m$  traffic to  $ONU_i$  during  $T_i(n - 1)$ , which is defined as the ratio of actual arrival packets to the updating period, denoted by  $\tilde{\lambda}_{m,i}(n - 1)$ ,  $m \in \{0, 1, 2\}$ . From Fig. 1, the actual packet arrival rate at prediction interval  $T_i(n - 2)$ ,  $\lambda_{m,i}(n - 2)$ , can be obtained from the reported queue occupancy and the granted bandwidth. Let  $A_{m,i}(n - 2)$  be the amount of actual arrival packets at  $Q_{m,i}$  in bytes at the prediction interval  $T_i(n - 2)$ . At prediction interval  $T_i(n - 1)$ , the PQ-DBA can obtain the amount of actual arrival packets  $A_{m,i}(n - 2)$  according to the reported

queue occupancy  $L_{m,i}(n - 1)$ ,  $L_{m,i}(n - 2)$  and the granted bandwidth  $G_{m,i}(n - 1)$ . If  $L_{m,i}(n - 1) > 0$ , then the amount of actual arrival packets  $A_{m,i}(n - 2)$  can be determined from  $G_{m,i}(n - 1)$ ,  $L_{m,i}(n - 1)$  and  $L_{m,i}(n - 2)$ . If  $L_{m,i}(n - 1) = 0$ , assume that the actual arrival packets arrive in uniform distribution. Therefore the  $A_{m,i}(n - 2)$  can be obtained by

$$A_{m,i}(n - 2) = \begin{cases} G_{m,i}(n - 1) - L_{m,i}(n - 2) + L_{m,i}(n - 1), & \text{if } L_{m,i}(n - 1) > 0, \\ \frac{G_{m,i}(n - 1) - L_{m,i}(n - 2)}{2}, & \text{if } L_{m,i}(n - 1) = 0. \end{cases} \quad (7)$$

Then ,the  $\lambda_{m,i}(n - 2)$  can be calculated by

$$\lambda_{m,i}(n - 2) = \frac{A_{m,i}(n - 2)}{T_i(n - 2)}. \quad (8)$$

In the implementation of the PRNN predictor, a fully connected recurrent neural network (RNN) structure with  $R$  neurons and  $p + q + R$  input nodes can be adopted [17]. However, the computational complexities are pretty high. Instead the PRNN structure is here considered for its computation efficiency. In our design, the prediction value of the  $\tilde{\lambda}_{m,i}(n - 1)$  can be determined from  $p$  previously measured samples  $\lambda_{m,i}(k)$ ,  $n - p - 1 \leq k \leq n - 2$ , and  $q$  predicted errors,  $e_{m,i}(j)$ ,  $n - q - 1 \leq j \leq n - 2$ , and  $e_{m,i}(j) = \lambda_{m,i}(j) - \tilde{\lambda}_{m,i}(j)$ . The  $\tilde{\lambda}_{m,i}(n - 1)$  can be expressed as

$$\tilde{\lambda}_{m,i}(n - 1) = H(\lambda_{m,i}(n - 2), \dots, \lambda_{m,i}(n - p - 1); \tilde{\lambda}_{m,i}(n - 2), \dots, \tilde{\lambda}_{m,i}(n - q - 1)) \quad (9)$$

where  $H(\bullet)$  is an unknown nonlinear function and the PRNN is adopted to approximate it. The PRNN refers to the RNN predictor with a pipelined structure, which consists of  $q$  levels of processing, and each level has a RNN module with  $N$  neurons,  $(d + N + 1)$  input nodes and a comparator, where  $d = p - q + 1$  and  $q \times N = R$ . For the detailed description of the PRNN predictor, please refer to [14].

Here, the extended recursive least square (ERLS) is applied as the learning algorithm for PRNN. In order to reduce the complexity, all the modules of the PRNN are designed to have exactly the same synaptic weight matrix. Hence, each level of the PRNN is a sub-prediction and each RNN module has a prediction error, and the total prediction errors of the PRNN predictor must be combined for weight adjustment. The prediction errors for the  $k$ -th RNN module around prediction interval  $T_i(n - 1)$ , denoted by  $e_k(n - 2)$ ,  $1 \leq k \leq q$ , and for the PRNN predictor, denoted by  $E_{m,i}(n - 2)$  are, respectively, defined as

$$e_k(n - 2) = \lambda_{m,i}(n - k - 1) - y_k(n - 2), \quad (10)$$

and

$$E_{m,i}(n - 2) = \sum_{k=1}^q \xi^{k-1} e_k^2(n - 2), \quad (11)$$



where  $y_k(n - 2)$  is the output of  $k$ -th RNN module, and  $\xi \in (0, 1]$  is the forgetting factor. The term  $\xi^{k-1}$  is an approximate measure of the memory of the individual modules in the PRNN. The cost function of ERLS, denoted by  $\varepsilon_{\text{ERLS}}(n - 2)$ , is defined as

$$\varepsilon_{\text{ERLS}}(n - 2) = \sum_{k=1}^{n-2} \xi^{n-k-2} E_{m,i}(k). \tag{12}$$

The ERLS algorithm minimizes the cost function in Eq. 12 and then updates the weights of the neurons in the modules accordingly. According to analysis in [13], ERLS considers present and previous errors, so it can minimize the margin of errors. When the PRNN/ERLS predictor has finished the prediction of the packet arrival rate during  $T_i(n - 1)$ ,  $\tilde{\lambda}_{m,i}(n - 1)$ , the amount of estimated arrival packets of  $Q_{m,i}$  in bytes at prediction interval  $T_i(n - 1)$ ,  $m \in \{0, 1, 2\}$ ,  $1 \leq i \leq M$ , denoted by  $\tilde{D}_{m,i}(n - 1)$ , can be calculated by

$$\tilde{D}_{m,i}(n - 1) = \tilde{\lambda}_{m,i}(n - 1) \times T_i(n - 1). \tag{13}$$

The predicted queue occupancy of  $Q_{m,i}$  at the end of  $T_i(n - 1)$ , denote it by  $P_{m,i}(n - 1)$ , can be obtained as

$$P_{m,i}(n - 1) = L_{m,i}(n - 1) + \tilde{D}_{m,i}(n - 1). \tag{14}$$

### 3.2 PQ-DBA scheme

The PQ-DBA scheme classifies the three types of services traffic from ONUs into six priorities. The first priority is real-time voice packets; the second priority is real-time video packets, which have violation problem of packet dropping probability requirement if they are not served at the next cycle; the third priority is video packets, which have violation problem of delay requirement if they are not served at the next cycle; then data packets have violation problem of starvation–threshold requirement as fourth priority; then, real-time video packets as fifth priority, and non-real-time data packets as the least priority. The PQ-DBA provides a QoS promotion capability to ONUs. Service packets are transmitted in accordance with their priorities, and these packets, which

packets first, and estimates the newly arrival packets during prediction interval, and then allocates the available bandwidth to each ONUs in a proportional method. The PQ-DBA scheme allocates bandwidth in unit of bytes to all ONUs from the service of the highest priority to that of the lowest one successively until all the bandwidths are used up. The granted bandwidth to the service of any priority for OUN<sub>*i*</sub> is based on the predicted queue occupancy of  $Q_{m,i}$ ,  $P_{m,i}$ , given in Eq. 14,  $m \in \{0, 1, 2\}$ ,  $1 \leq i \leq M$ , and  $L_{dp,i}$ ,  $L_{d,i}$ , and  $L_{w,i}$  given in Eqs. 2, 4, and 6, respectively. Denote  $B$  the total bandwidth of fiber link. The bandwidth allocation of PQ-DBA scheme is described in detail as follows:

#### Step 1 : Bandwidth allocation to voice

The voice packet plays the highest priority of service since it is strictly delay sensitive. Denote the allocated bandwidth to voice packets in  $Q_{0,i}$  by  $G'_{0,i}$ . Based on the predicted occupancy of  $Q_{0,i}$ ,  $P_{0,i}$ , and the total bandwidth  $B$ ,  $G'_{0,i}$  is given by

$$G'_{0,i} = \begin{cases} P_{0,i}, & \text{if } \sum_{i=1}^M P_{0,i} \leq B, \\ B \times \frac{P_{0,i}}{\sum_{i=1}^M P_{0,i}}, & \text{elsewhere.} \end{cases} \tag{15}$$

#### Step 2 : Bandwidth allocation to video packets with the second and the third priorities

To ensure the QoS requirements of video packets, the PQ-DBA secondly allocates the bandwidth to the video packets, which will be dropped if they are not transmitted at the next cycle. Denote  $G'_{1,i}$  as the allocated bandwidth to video packets with the second and the third priorities in  $Q_{1,i}$ . Based on the reported video packet with problem of dropping probability and delay,  $L_{d,i}$ ,  $L_{dp,i}$ , and the residual bandwidth of fiber link,  $B - \sum_{i=1}^M G'_{0,i}$ , the allocated bandwidth to video packets with the second and the third priorities  $G'_{1,i}$  is given by

$$G'_{1,i} = \begin{cases} L_{dp,i}, & \text{if } B - \sum_{i=1}^M [G'_{0,i}] \geq \sum_{i=1}^M L_{dp,i}, \\ L_{d,i} + (B - \sum_{i=1}^M [G'_{0,i} + L_{d,i}])^+ \times \frac{L_{dp,i} - L_{d,i}}{\sum_{i=1}^M (L_{dp,i} - L_{d,i})}, & \text{if } \sum_{i=1}^M L_{d,i} < B - \sum_{i=1}^M G'_{0,i} < \sum_{i=1}^M L_{dp,i}, \\ (B - \sum_{i=1}^M G'_{0,i})^+ \times \frac{L_{d,i}}{\sum_{i=1}^M L_{d,i}}, & \text{if } B - \sum_{i=1}^M G'_{0,i} \leq \sum_{i=1}^M L_{d,i}. \end{cases} \tag{16}$$

will violate the QoS requirement, will be promoted to higher priority.

Upon receiving the REPORT messages from each ONU, the PQ-DBA at OLT calculates the amount of QoS-promoted

#### Step 3 : Bandwidth allocation to data packets with the fourth priority

The PQ-DBA continues to allocate the bandwidth to the data packets whose waiting time exceeds the starvation–threshold time  $T_d^*$ , to ensure QoS requirement of data service and avoid starvation. Denote  $G'_{2,i}$  as the allocated bandwidth in  $Q_{2,i}$ . Based on the amount of reported data packets considering starvation,  $L_{w,i}$ , and the residual bandwidth of fiber link,  $B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i}]$ ,  $G'_{2,i}$  is given by

$$G'_{2,i} = \begin{cases} L_{w,i}, & \text{if } B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i}] > \sum_{i=1}^M L_{w,i}, \\ (B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i}])^+ \times \frac{L_{w,i}}{\sum_{i=1}^M L_{w,i}}, & \text{elsewhere.} \end{cases} \tag{17}$$

*Step 4* : Bandwidth allocation to video packets with the fifth priority

The PQ-DBA assigns the bandwidth to the unallocated video packets. Denote  $G''_{1,i}$  as the allocated bandwidth to video packets with the fifth priority in  $Q_{1,i}$ . Based on the amount of unallocated video packets,  $P_{1,i} - G'_{1,i}$ , and the residual bandwidth of fiber link,  $B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i} + G'_{2,i}]$ ,  $G''_{1,i}$  is given by

$$G''_{1,i} = \begin{cases} P_{1,i} - G'_{1,i}, \\ (B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i} + G'_{2,i}])^+ \times \frac{P_{1,i} - G'_{1,i}}{\sum_{i=1}^M (P_{1,i} - G'_{1,i})}, \end{cases}$$

*Step 5* : Bandwidth allocation to data packets with the sixth priority

The PQ-DBA assigns the bandwidth to the unallocated data packets. Denote  $G''_{2,i}$  as the allocated bandwidth to data packets with the sixth priority in  $Q_{2,i}$ . Based on the amount of unallocated data packets,  $P_{2,i} - G'_{2,i}$ , and the residual bandwidth of the fiber link,  $B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i} + G'_{2,i} + G''_{1,i}]$ ,  $G''_{2,i}$  is given by

$$G''_{2,i} = \begin{cases} P_{2,i} - G'_{2,i}, & \text{if } B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i} + G'_{2,i} + G''_{1,i}] \\ > \sum_{i=1}^M (P_{2,i} - G'_{2,i}), \\ (B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i} + G'_{2,i} + G''_{1,i}])^+ \\ \times \frac{P_{2,i} - G'_{2,i}}{\sum_{i=1}^M (P_{2,i} - G'_{2,i})}, & \text{elsewhere.} \end{cases} \tag{19}$$

*Step 6* : Residual bandwidth allocation

Finally, the PQ-DBA assigns voice and video packets proportionally, sharing the residual bandwidth based on their predicted queue occupancy to make the best use of the bandwidth and guarantee QoS further. Denote  $G''_{0,i}$  and  $G'''_{1,i}$  as the allocated residual bandwidth per voice packets in  $Q_{0,i}$  and video packets in  $Q_{1,i}$ , respectively. Based on the predicted queue occupancy,  $P_{0,i}$ ,  $P_{1,i}$  and the residual

bandwidth,  $B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i} + G'_{2,i} + G''_{1,i} + G''_{2,i}]$ ,  $G''_{0,i}$  and  $G'''_{1,i}$  are given by

$$\begin{cases} G''_{0,i} = (B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i} + G'_{2,i} + G''_{1,i} \\ + G''_{2,i}])^+ \times \frac{P_{0,i}}{\sum_{i=1}^M (P_{0,i} + P_{1,i})}, \\ G'''_{1,i} = (B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i} + G'_{2,i} + G''_{1,i} \\ + G''_{2,i}])^+ \times \frac{P_{1,i}}{\sum_{i=1}^M (P_{0,i} + P_{1,i})}. \end{cases} \tag{20}$$

$$\text{if } B - \sum_{i=1}^M [G'_{0,i} + G'_{1,i} + G'_{2,i}] > \sum_{i=1}^M (P_{1,i} - G'_{1,i}), \text{ elsewhere.} \tag{18}$$

*Step 7* : GATE message generation

Based on the allocated bandwidth,  $G'_{0,i}$ ,  $G'_{1,i}$ ,  $G'_{2,i}$ ,  $G''_{0,i}$ ,  $G''_{1,i}$ ,  $G''_{2,i}$  and  $G'''_{1,i}$ , the final granted bandwidth in GATE message  $G_{0,i}$ ,  $G_{1,i}$ , and  $G_{2,i}$  for ONU<sub>*i*</sub> are given by

$$\begin{cases} G_{0,i} = G'_{0,i} + G''_{0,i}, \\ G_{1,i} = G'_{1,i} + G''_{1,i} + G'''_{1,i}, \\ G_{2,i} = G'_{2,i} + G''_{2,i}. \end{cases} \tag{21}$$

The  $\{G_{0,i}, G_{1,i}, G_{2,i}\}$  are included in the GATE message  $G_i$ , and the OLT sends the GATE message to the ONU<sub>*i*</sub>,  $1 \leq i \leq M$ . Each ONU follows the information in the GATE message to transmit its own packets. At the same ONU, if there is some bandwidth left at some queues, the rest bandwidth can be reallocated to other queues of the same ONU. The PQ-DBA scheme not only decreases packet delay and packet dropping probability but also increases bandwidth efficiency.

#### 4 Simulation results and discussions

An event-driven packet-based simulation is elaborated to show the performance comparison among three upstream

transmission schemes: the proposed PQ-DBA, the proposed PQ-DBA but without prediction, and the DBAM [9]. In the simulations, the considered EPON system is in an OLT with 32 ONUs connected configuration, and there is one 1:32 splitter located at 20 km away from OLT, and at 5 km distance to each ONU. The downstream rate  $R_E = 1$  Gbps, and the upstream rate  $R_U = 100$  Mbps. Each queue in ONU is with the same buffer space of 1 Mb. The guard time for laser ON/OFF between each ONU transmission is set to  $1 \mu\text{s}$ . For simulation convenience, the system cycle and the prediction interval of all ONUs are assumed to be fixed at 0.72 ms. Every simulation result includes 100 simulation cycles, and each of which contains 1,000 updating periods. As for the PRNN predictor, parameters are selected as:  $N = 2$ ;  $p = 4$ ;  $q = 2$ ; and  $\xi = 0.99$  [14].

The voice traffic source is modeled as two-state Markov modulated deterministic process (MMDP) with  $\alpha$  and  $\beta$  as the transition rates. The mean durations of talk spurts and silence periods are assumed to be exponentially distributed with  $1/\alpha = 1$  s and  $1/\beta = 1.35$  s, respectively. Assume the voice traffic is packetized in the ONU by placing 24 bytes of data in a packet. By adding the overhead such as Ethernet, UDP (User Data Protocol) and IP (Internet Protocol) headers in a packet, the packet results in a 70-byte of frame. The generation rate of voice packets is constant bit rate (CBR) on every  $125 \mu\text{s}$  during talk spurts (ON state), but none during silence (OFF state). On the other hand, the highly burst video and data packets are modeled by a superposition of  $K$  independent and identical ON-OFF Pareto-distributed source in order to generate self-similar and long-range dependence (LRD) traffic. The typical mean ON period of the Pareto distribution is 7.2 s with a “heaviness” of  $\rho = 1.4$  for relationship to Hurst parameter of 0.8, and the typical mean OFF period is 10.5 s with a “heaviness” of  $\rho = 1.2$  [18]. The packet sizes are uniformly distributed between 64 and 1,518 bytes. The traffic arrival rates for each ONU are set as follows: (i) voice service: 4.48 Mbps with CBR, (ii) video service: variable bit rate (VBR) at 0.55–15.55 Mbps, (iii) data service: VBR at 0.28–7.67 Mbps.

The voice delay criterion is according to the ITU-T Recommendation G.114: 1.5 ms for “one way transmission time” in access network. The voice dropping probability is set to zero. The video packet delay requirement  $T_d^*$ , the video packet dropping probability requirement  $P_d^*$ , and the data packets starvation–threshold time  $T_w^*$ , are defined as 10 ms, 1%, and 500 ms, respectively. For DBAM [9], here we set the SLA of the DBAM scheme the same QoS requirement as the PQ-DBA scheme.

Figure 2 shows the system utilization versus the traffic intensity. Under the QoS constraint, it can be found that the proposed PQ-DBA has higher system utilization than the DBAM and the PQ-DBA without prediction by an amount of 4, and 2% in average, respectively. It is because the PQ-DBA

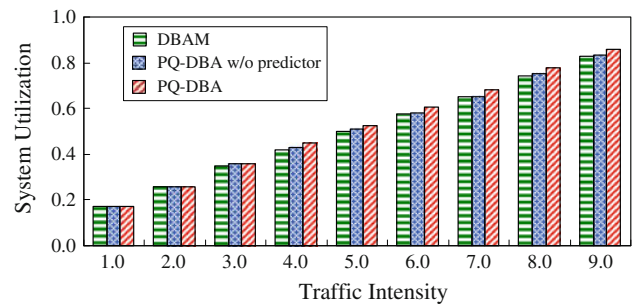


Fig. 2 System utilization versus traffic intensity

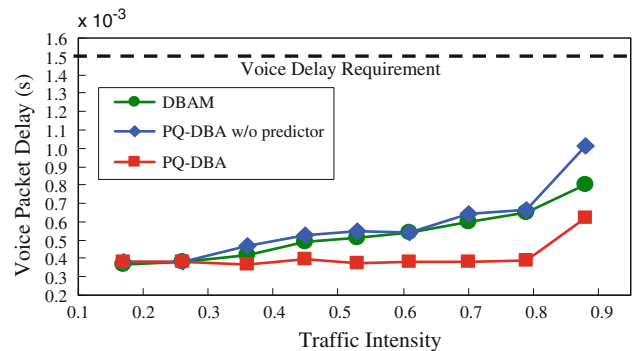
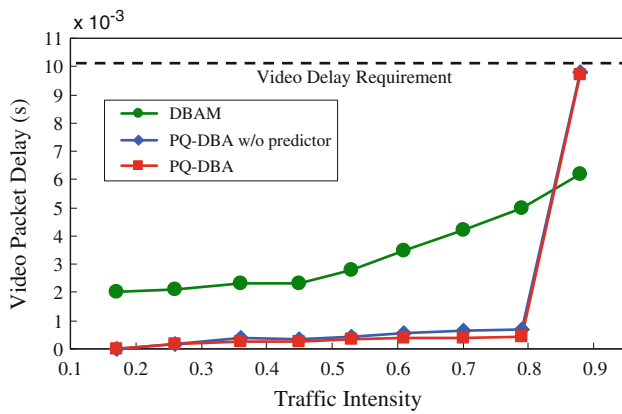


Fig. 3 The average voice packet delay versus traffic intensity

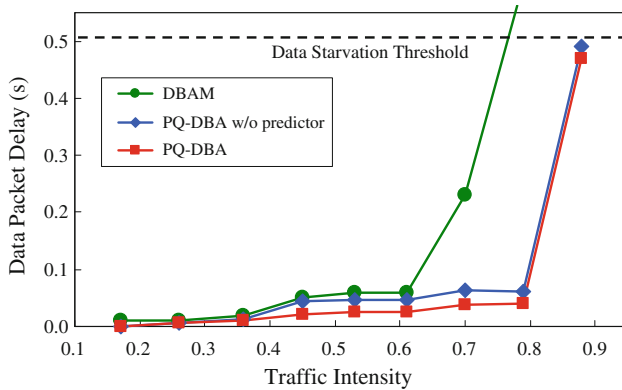
and the PQ-DBA without prediction allocate bandwidth step by step to six priorities, rather than the DBAM by a *limited bandwidth allocation* (LBA) to each class in advance. The PQ-DBA and the PQ-DBA without prediction dynamically promote some video packets cycle-by-cycle to avoid packet dropping due to the violation of delay requirement, whereas the DBAM scheme just applies *priority queueing*. Moreover, the PRNN/ERLS also helps the PQ-DBA scheme to enhance the efficiency of bandwidth allocation by providing more precise estimation and thus more accurate allocation bandwidth for ONUs than the PQ-DBA without prediction and the DBAM scheme with linear prediction.

Figure 3 shows the average voice packet delays versus traffic intensity. It can be observed that as traffic intensity is larger than 0.3, the average voice packet delay of the PQ-DBA outperforms by about 21 and 26% in average over the DBAM scheme and the PQ-DBA without prediction, respectively. It is because the PQ-DBA with PRNN/ERLS makes a precise prediction for new voice packet arrivals during each prediction interval. The DBAM also adopts prediction, but the prediction is linear, which would be less sophisticated than the nonlinear PRNN/ERLS. Notice that the three schemes have the voice packets as the first priority to serve.

Figure 4 shows the average video packet delay versus the traffic intensity. It can be observed that the video packet delays of the three schemes are within the delay requirement when the traffic intensity is below 0.9. The PQ-DBA with PRNN/ERLS improves the video packet delay time by



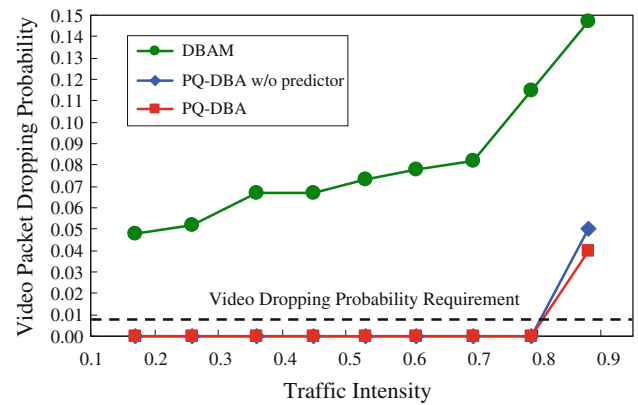
**Fig. 4** The average video packet delay versus traffic intensity



**Fig. 5** The average data packet delay versus traffic intensity

about 90% better than that in DBAM scheme (and 29% over PQ-DBA without prediction) when the traffic intensity is below 0.8. It is because the PQ-DBA makes the video packet with the problem of delay requirement be transmitted with a higher priority in order to decrease the delay and satisfy the dropping probability. When the traffic intensity exceeds 0.9, the average video packet delay for PQ-DBA and PQ-DBA without prediction is close to the video delay requirement. It is because the burst data packets with starvation condition are greatly increased and could be sent prior to normal video packets.

Figure 5 illustrates the average data packet delay versus the traffic intensity. The starvation ratio is defined as the proportion of the number of data packets with a delay time exceeding the  $T_w^*$  among the transmitted data packets. It can be found that when the traffic intensity is larger than 0.8, the data starvation in DBAM happens earlier than that in the PQ-DBA scheme and PQ-DBA without prediction. It is because the priority of data packets can be promoted by considering the starvation–threshold requirement in the PQ-DBA scheme and PQ-DBA without prediction, but is the lowest by *priority queueing* in the DBAM scheme. When the traffic intensity is larger than 0.9, the data packet delay time



**Fig. 6** The average video packet dropping probability versus traffic intensity

increases apparently in the PQ-DBA and PQ-DBA without prediction. It is because the data packets are the lowest priority, and the system does not have enough bandwidth to support the greatly increasing of the burst data arrivals, the starvation occurs.

Figure 6 illustrates the average video packet dropping probability versus the traffic intensity in EPON. The video packet dropping probability plays an index role of video service degradation; a high quality transmission must be ensured otherwise the mosaic phenomenon might be happened due to low networking QoS. Unfortunately, the video packet dropping probability in DBAM cannot be satisfied, due to it adopts LBA by setting an upper bound of grants to video packets, and inaccurate linear class-based traffic prediction. However, the average video dropping probability of the PQ-DBA scheme and PQ-DBA without prediction is almost zero when the traffic intensity is below 0.8. It is because the priority of video packets with the problem of delay requirement can be raised in both the PQ-DBA scheme and PQ-DBA without prediction scheme. When the traffic intensity is larger than 0.8, the average video dropping probability exceeds the video dropping probability requirement due to the fiber link capacity limitation.

## 5 Conclusion

In this article, a PRNN/ERLS-based predictive QoS-promoted dynamic bandwidth allocation (PQ-DBA) algorithm has been proposed for EPON upstream transmission. Specifically, the PQ-DBA promotes packets to higher priority by QoS control and then performs bandwidth allocation based on the priorities from the highest to the lowest, which guarantees the QoS requirements of real-time services while preserving the grade of service of the non-real-time packets. Moreover, thanks to the nonlinear PRNN/ERLS predictor [15] with high accuracy and fast convergence, the PQ-DBA



accurately predicts the burst traffic and thus significantly improves the system efficiency. Simulation results show that the PQ-DBA scheme achieves the highest system throughput, lowest delay and guarantees the stringent QoS requirement of video packets, which outperforms DBAM and PQ-DBA without prediction.

## References

- [1] MaGarry, M.P., Maier, M., Reisslein, M.: Ethernet PONs: a survey of dynamic bandwidth allocation (DBA) algorithms. *IEEE Opt. Commun.* **42**(8), S8–15 (2004). doi:[10.1109/MCOM.2004.1321381](https://doi.org/10.1109/MCOM.2004.1321381)
- [2] Kramer, G., Mukherjee, B., Pesavento, G.: IPACT: a dynamic protocol for an Ethernet PON (EPON). *IEEE Commun. Mag.* **40**(2), 74–80 (2002)
- [3] IEEE Standard 802.3ah-2004, IEEE P802.3ah Ethernet in the First Mile Task Force.
- [4] Cheng, H., Chen, M., Xie, S.: A dynamic bandwidth allocation scheme supporting different priority services in EPON. *Proc. Int. Soc. Opt. Eng. (SPIE)* **5626**(2), 1123–1127 ISBN 0-8194-5580-6 (2005). doi:[10.1117/12.575286](https://doi.org/10.1117/12.575286)
- [5] Assi, C.M., Ye, Y., Dixit, S., Ali, M.A.: Dynamic bandwidth allocation for quality-of-service over Ethernet PONs. *IEEE J. Sel. Areas Commun.* **21**(9), 1467–1477 (2003)
- [6] Xie, J., Jiang, S., Jiang, Y.: A dynamic bandwidth allocation scheme for differentiated services in EPONs. *IEEE Commun. Mag.* **42**(8), 32–39 (2004). doi:[10.1109/MCOM.2004.1321385](https://doi.org/10.1109/MCOM.2004.1321385)
- [7] Yang, Y., Nho, J., Mahalik, N.P., Kim, K., Ahn, B.: QoS provisioning in the EPON systems with traffic-class burst-polling based delta DBA. *IEICE Trans. Commun.* **89**, 419–426 (2006)
- [8] Wu, S., Ding, Q., Chung, K.C.: Improving the network performance using prediction based longest queue first (PLQF) scheduling algorithm. In: *ATM (ICATM 2001) and 4th IEEE International Conference on High Speed Intelligent Internet Symposium*, Seoul, South Korea, pp. 344–348, ISBN: 0-7803-7093-7. (2001)
- [9] Luo, Y., Ansari, N.: Bandwidth allocation for multiservice access on EPONs. *IEEE Commun. Mag.* **43**(2), 16–21 (2005). doi:[10.1109/MCOM.2005.1391498](https://doi.org/10.1109/MCOM.2005.1391498)
- [10] Luo, Y., Ansari, N.: Limited sharing with traffic prediction for dynamic bandwidth allocation and QoS provisioning over Ethernet passive optical networks. *OSA J. Opt. Netw.* **4**(9), 561–572 (2005)
- [11] Hwang, I.S., Shyu, Z.D., Ke, L.Y., Chang, C.C.: A novel early DBA mechanism with prediction-based fair excessive bandwidth reallocation scheme in EPON. In: *IEEE Proceedings of the Sixth International Conference on Networking*, Sainte-Luce, Martinique, France, pp. 75–80, ISBN:0-7695-2805-8 (2008)
- [12] Yin, S., Luo, Y., Ansari, N., Wang, T.: Non-linear predictor-based dynamic bandwidth allocation over TDM-PONs: stability analysis and controller design. *IEEE International Conference on Communications*, Beijing, pp. 5186–5190, ISBN: 978-1-4244-2075-9 (2008)
- [13] Haykin, S., Li, L.: Nonlinear adaptive prediction of nonstationary signals. *IEEE Trans. Signal Process.* **43**(2), 526–535 (1995)
- [14] Baltersee, L., Chambers, J.A.: Nonlinear adaptive prediction of speech using a pipelined recurrent neural network. *IEEE Trans. Signal Process.* **46**(8), 2207–2216 (1998)
- [15] Chang, C.J., Chen, B.W., Liu, T.Y., Ren, F.C.: Fuzzy/neural congestion control for integrated voice and data DS-CDMA/FRMA cellular networks. *IEEE J. Sel. Areas Commun.* **18**(2), 183–293 (2000)
- [16] Floyd, S., Jacobson, V.: Random early detection gateways for congestion avoidance. *IEEE/ACM Trans. Netw.* **1**(4), 397–413 ISSN:1063-6692 (1993)
- [17] Connor, J.T., Martin, R.D., Atlas, L.E.: Recurrent neural networks and robust time series prediction. *IEEE Trans. Neural Netw.* **5**(2), 240–254 (1994)
- [18] Willinger, W., Taqqu, M.S., Sherman, R., Wilson, D.V.: Self-similarity through high-variability: statistical analysis of Ethernet LAN traffic at the source level. *IEEE/ACM Trans. Netw.* **5**(1), 71–86 ISSN: 1063-6692 (1997)

## Author Biographies



**Jan-Wen Peng** received the BSEE from Tamkang University, Taiwan, ROC, and MSEE degrees from University of Florida, USA, in 1986 and 1991, respectively. Since 1992, he joined the Telecommunication Laboratories, Directorate General of Telecommunications, Ministry of Communications, Taiwan, as an assistant engineer in Subscriber Loop Research Department. Currently he is an associate research engineer of

Broadband Transport and Access Technology Laboratory, Chunghwa Telecom Laboratories. His specialty is in the area of digital communications, broadband access technology, and optical access network. He is now a candidate for the Ph.D. degree in the area of Electrical Engineering at National Chiao Tung University, Taiwan. He is also a member of the IEEE.



**Chung-Ju Chang** was born in Taiwan, ROC, in August, 1950. He received the B.E. and M.E. degrees in Electronics Engineering from National Chiao Tung University, Hsinchu, Taiwan, in 1972 and 1976, respectively, and the Ph.D. degree in Electrical Engineering from National Taiwan University, Taiwan, in 1985. From 1976 to 1988, he was with Telecommunication Laboratories, Directorate General of Telecommunications,

Ministry of Communications, Taiwan, as a Design Engineer, Supervisor, Project Manager, and then Division Director. In 1988, he joined the Faculty of the Department of Communication Engineering, College of Electrical Engineering and Computer Science, National Chiao Tung University, as an Associate Professor. He has been a Professor since 1993. He was Director of the Institute of Communication Engineering from August 1993 to July 1995, Chairman of Department of Communication Engineering from August 1999 to July 2001, and Dean of the Research and Development Office from August 2002 to July 2004. Also, he was an Advisor for the Ministry of Education to promote the education of communication science and technologies for colleges and universities in Taiwan during 1995–1999. Moreover, he serves as Editor for *IEEE Communications Magazine* and Associate Editor for *IEEE Transactions Vehicular Technology*. His research interests include performance evaluation, radio resources management for wireless communication networks, and traffic control for broadband networks. Dr. Chang is a member of the Chinese Institute of Engineers (CIE) and *IEEE Fellow*.



**Po-Lung Tien** received the B.S. degree in applied mathematics, the M.S. degree in computer and information science, and the Ph.D. degree in computer and information engineering from the National Chiao Tung University, Hsinchu, Taiwan, ROC, in 1992, 1995, and 2000, respectively. In 2000, he joined National Chiao Tung University, where he is currently a Research Assistant Professor of the Department of Computer

Science and Information Engineering. His current research interests include optical networking, wireless local networking, multimedia communications, performance modeling and analysis, and applications of soft computing.