

Load balanced Birkhoff–von Neumann switches, part II: multi-stage buffering

Cheng-Shang Chang^{*}, Duan-Shin Lee, Ching-Ming Lien

Institute of Communications Engineering, National Tsing Hua University, Hsinchu 300, Taiwan, ROC

Received 8 October 2001; accepted 8 October 2001

Abstract

The main objective of this sequel is to solve the out-of-sequence problem that occurs in the load balanced Birkhoff–von Neumann switch with one-stage buffering. We do this by adding a load-balancing buffer in front of the first stage and a resequencing-and-output buffer after the second stage. Moreover, packets are distributed at the first stage according to their *flows*, instead of their *arrival times* in part I. In this paper, we consider multicasting flows with two types of scheduling policies: the first come first served (FCFS) policy and the earliest deadline first (EDF) policy. The FCFS policy requires a jitter control mechanism in front of the second stage to ensure proper ordering of the traffic entering the second stage. For the EDF scheme, there is no need for jitter control. It uses the departure times of the corresponding FCFS output-buffered switch as deadlines and schedules packets according to their deadlines. For both policies, we show that the end-to-end delay through our multi-stage switch is bounded above by the sum of the delay from the corresponding FCFS output-buffered switch and a constant that only depends on the size of the switch and the number of multicasting flows supported by the switch. © 2002 Elsevier Science B.V. All rights reserved.

Keywords: Multi-stage switches; Load balancing; Scheduling; Multicasting; Performance bounds

1. Introduction

In part I [5], we proposed the load balanced Birkhoff–von Neumann switch with one-stage buffering (see Fig. 1). Such a switch consists of two stages of crossbar switching fabrics and one stage of buffering. The buffer at the input port of the second stage uses the virtual output queuing (VOQ) technique to solve the problem of head-of-line blocking. In such a switch, packets are of the same size. Also, time is slotted and synchronized so that exactly one packet can be transmitted within a time slot. In a time slot, both crossbar switches set up connection patterns corresponding to permutation matrices that are periodically generated from a one-cycle permutation matrix.

The reasoning behind such a switch architecture is as follows: since the connection patterns are periodic, packets from the same input port of the first stage are distributed in a round-robin fashion to the second stage according to their arrival times. Thus, the first stage performs load balancing for the incoming traffic. As the traffic coming into the

second stage is load balanced, it suffices to use simple periodic connection patterns to perform switching at the second stage. This is shown in Ref. [5] as a special case of the original Birkhoff [1]–von Neumann [13] decomposition used in Refs. [3,4]. There are several advantages of using such an architecture, including scalability, low hardware complexity, 100% throughput, low average delay in heavy load and bursty traffic, and efficient buffer usage. However, the main drawback of the load balanced Birkhoff–von Neumann switch with one-stage buffering is that packets might be out of sequence.

The main objective of this sequel is to solve the out-of-sequence problem that occurs in the load balanced Birkhoff–von Neumann switch with one-stage buffering. One quick fix is to add a resequencing-and-output buffer after the second stage. However, as packets are distributed according to their *arrival times* at the first stage, there is no guarantee on the size of the resequencing-and-output buffer to prevent packet losses. For this, one needs to distributed packets according to their *flows*, as indicated in the paper by Iyer and McKeown [9]. This is done by adding a flow splitter and a load-balancing buffer in front of the first stage (see Fig. 2). For an $N \times N$ switch, the load-balancing buffer at each input port of the first stage consists of N virtual output queues (VOQ) destined for the N output

^{*} Corresponding author.

E-mail addresses: cschang@ee.nthu.edu.tw (C.-S. Chang),
lds@cs.nthu.edu.tw (D.-S. Lee),
keiichi@gibbs.ee.nthu.edu.tw (C.-M. Lien).

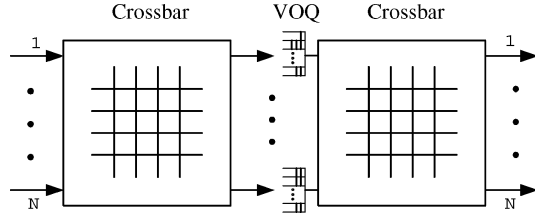


Fig. 1. The load balanced Birkhoff-von Neumann switch with one-stage buffering.

ports of that stage. Packets from the same *flow* are split in the round-robin fashion to the N virtual output queues and scheduled under the first come first served (FCFS) policy. By so doing, load balancing can be achieved for each flow as packets from the same flow are split almost evenly to the input ports of the second stage. More importantly, as pointed out in Ref. [9], the delay and the buffer size of the load-balancing buffer are bounded by constants that only depend on the size of the switch and the number of flows. The resequencing-and-output buffer after the second stage not only performs resequencing to keep packets in sequence, but also stores packets waiting for transmission from the output links.

In this paper, we consider a traffic model with multicasting flows. This is a more general model than the point-to-point traffic model in Ref. [9]. A multicasting flow is a stream of packets that has one common input and a set of common outputs. For the multicasting flows, fan-out splitting (see e.g. Ref. [8]) is performed at the central buffers (the VOQ in front of the second stage). The central buffers are assumed to be infinite so that no packets are lost in the switch. We consider two types of scheduling policies in the central buffers: the FCFS policy and the earliest deadline first (EDF) policy. For the FCFS policy, a jitter control mechanism is added in the VOQ in front of the second stage. Such a jitter control mechanism delays every packet to its maximum delay at the first stage so that the flows entering the second stage are simply time-shifted flows of the original ones. Our main result for the FCFS scheme with

jitter control is the following theorem. The proof of Theorem 1 will be given in Section 2.

Theorem 1. *Suppose that all the buffers are empty at time 0. Then the following results hold for FCFS scheme with jitter control.*

- (i) *The end-to-end delay for a packet through our switch with multi-stage buffering is bounded above by the sum of the delay through the corresponding FCFS output-buffered switch and $(N - 1)L_{\max} + NM_{\max}$, where L_{\max} (resp. M_{\max}) is the maximum number of flows at an input (resp. output) port.*
- (ii) *The load-balancing buffer at an input port of the first stage is bounded above by NL_{\max} .*
- (iii) *The delay through the load-balancing buffer at an input port of the first stage is bounded above by $(N - 1)L_{\max}$.*
- (iv) *The resequencing-and-output buffer at an output port of the second stage is bounded above by NM_{\max} .*

In the EDF scheme (see Fig. 3), every packet is assigned a deadline that is the departure time from the corresponding FCFS output-buffered switch. Packets are scheduled according to their deadlines in the central buffers. For the EDF scheme, there is no need to implement the jitter control mechanism in the FCFS scheme. As such, average packet delay can be greatly reduced. However, as there is no jitter control, one might need a larger resequencing buffer than that in the FCFS scheme with jitter control. Since the first stage is the same as that in the FCFS scheme, the delay and the buffer size are still bounded by $(N - 1)L_{\max}$ and NL_{\max} respectively. Moreover, we show the following theorem for the EDF scheme. Its proof will be given in Section 3.

Theorem 2. *Suppose that all the buffers are empty at time 0. Then the following results hold for the EDF scheme.*

- (i) *The end-to-end delay for a packet through our switch with multi-stage buffering is bounded above by the sum of*

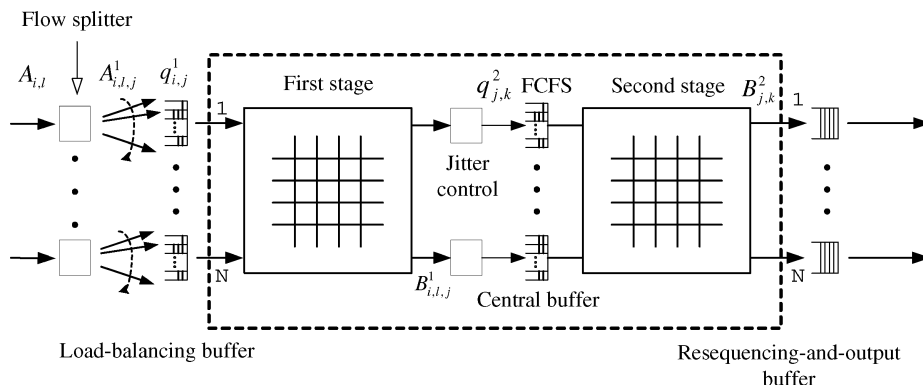


Fig. 2. The load balanced switch with multi-stage buffering under FCFS.

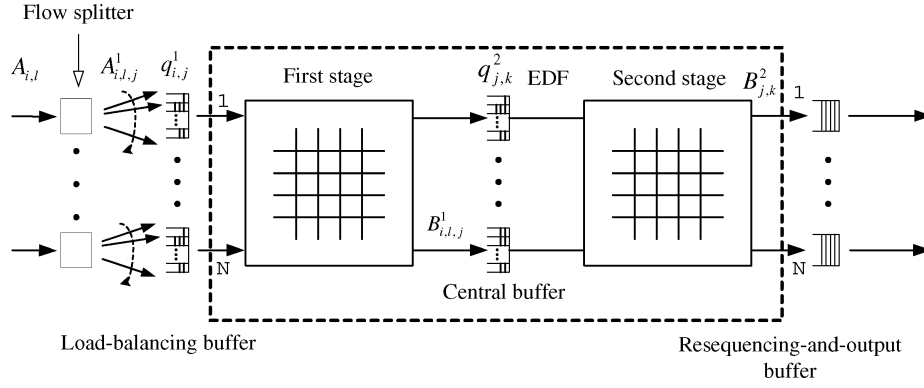


Fig. 3. The load balanced switch with multi-stage buffering under EDF.

the delay through the corresponding FCFS output-buffered switch and $(N-1)(L_{\max} + M_{\max})$.

(ii) The resequencing-and-output buffer at an output port of the second stage is bounded above $(N-1)(L_{\max} + M_{\max})$.

Computing the departure times from the corresponding FCFS output-buffered switch needs global information of all the inputs. A simple way is to use the packet arrival times as deadlines. Then the EDF scheme based on arrival times yields the same departure order except those packets that arrives at same time. Since there are at most M_{\max} packets that can arrive at the same time to an output port of the corresponding output-buffered switch, the end-to-end delay for a packet through the multi-stage switch using arrival times as deadlines is bounded above by the sum of the delay through the corresponding FCFS output-buffered switch and $(N-1)L_{\max} + NM_{\max}$. Also, the resequencing-and-output buffer at an output port of the second stage in this case is bounded above $(N-1)L_{\max} + NM_{\max}$.

It is of some interest to compare our schemes with the combined input–output queuing (CIOQ) switch (see e.g. Refs. [6,12]). The CIOQ switch provides the exact emulation of the corresponding output-buffered switch. For our schemes, we only have bounded differences between the departure times in our schemes and the departure times of the corresponding output-buffered switch. However, the CIOQ switch requires an internal rate speedup of two and a complicated scheduling algorithm that prohibits its practical use. In our schemes, there is no need for internal rate speedup (it is done by two switching fabrics) and the

complexity of the scheduling algorithm is low ($O(1)$ for the FCFS scheme).

2. The FCFS scheme

In this section, we prove Theorem 1. As discussed in Section 1, we consider $N \times N$ switches with multicasting flows under the FCFS scheduling policy. To be precise, let L_i be the number of multicasting flows through the i th input port. Denote by $A_{i,\ell}(t)$ the cumulative number of packet arrivals by time t from the ℓ th multicasting flow at the i th input port, $i = 1, \dots, N$, $\ell = 1, \dots, L_i$. Also, let $S_{i,\ell}$ be the set of outputs of that flow, $S^*(k) = \{(i, \ell) : k \in S_{i,\ell}\}$ be the set of multicasting flows through the k th output, and $M_k = |S^*(k)|$ be the number of multicasting flows through the k th output port. Define $L_{\max} = \max_{1 \leq i \leq N} L_i$ as the maximum number of multicasting flow through an input port and $M_{\max} = \max_{1 \leq k \leq N} M_k$ as the maximum number of multicasting flow through an output port.

2.1. Analysis for the output-buffered switch under FCFS

Now consider feeding these multicasting flows to an $N \times N$ output-buffered switch under the FCFS policy (see Fig. 4). Assume that there is an infinite buffer at each output port and that all the buffers are empty at time 0. Let $A_k^0(t)$ be the cumulative number of arrivals at the k th output buffer by time t , $q_k^0(t)$ be the number of packets at the k th output buffer at time t , and $B_k^0(t)$ be the cumulative number of departures at the k th output buffer by time t . Since an output-buffered switch is a work-conserving link with the constant capacity

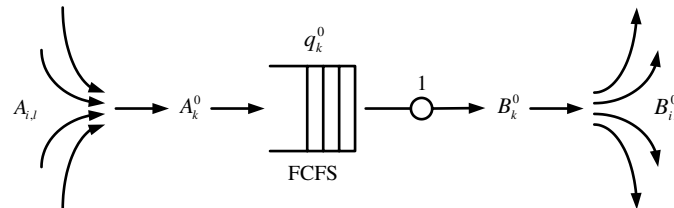
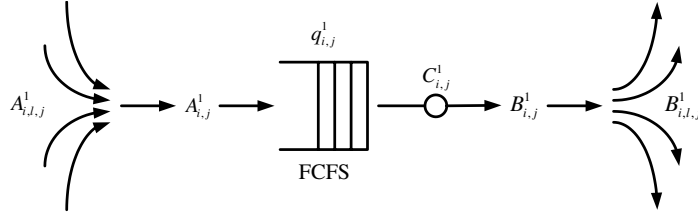


Fig. 4. The FCFS output-buffered switch with multicasting flows.

Fig. 5. The j th VOQ at the i th input port of the first stage.

1, one has the following well-known representation (see e.g. Ref. [2], Lemma 1.3.1)

$$q_k^o(t) = \max_{0 \leq s \leq t} [A_k^o(t) - A_k^o(s) - (t - s)], \quad (1)$$

where

$$A_k^o(t) = \sum_{(i,\ell) \in S^o(k)} A_{i,\ell}(t). \quad (2)$$

Moreover,

$$B_k^o(t) = A_k^o(t) - q_k^o(t) = \min_{0 \leq s \leq t} [A_k^o(s) + (t - s)]. \quad (3)$$

As the scheduling policy is FCFS and there are at most M_k arrivals per unit of time, packets that arrive at the k th buffer at time t will depart between $t + q_k^o(t) - M_k + 1$ and $t + q_k^o(t)$. Also, note from Eq. (3) that

$$B_k^o(t) - B_k^o(s) \leq t - s, \quad \text{for all } s \leq t, \quad (4)$$

as there is at most one packet coming out from an output port per time slot.

2.2. Analysis for the load-balancing buffer

As discussed in Section 1, the load-balancing buffer at each input port of the first stage consists of N virtual output queues (VOQ) destined for the N output ports of that stage. Packets from the same flow are split in the round-robin fashion to the N virtual output queues and scheduled under the FCFS policy. Without loss of generality, we assume that the first packet of a flow is always assigned to the first VOQ. To be precise, let $A_{i,\ell_j}^1(t)$ be the cumulative number of $A_{i,\ell}$ -flow packets that are split into the j th VOQ at the i th input port of the first stage by time t . Then

$$A_{i,\ell_j}^1(t) = \left\lceil \frac{A_{i,\ell}(t) - j + 1}{N} \right\rceil, \quad j = 1, \dots, N, \quad (5)$$

and

$$A_{i,\ell}(t) = \sum_{j=1}^N A_{i,\ell_j}^1(t), \quad (6)$$

where $\lceil x \rceil$ (resp. $\lfloor x \rfloor$) is the ceiling (resp. floor) function of x .

Now consider the j th VOQ at the i th input port of the first stage (see Fig. 5). Let $A_{i,j}^1(t)$ be the cumulative number of arrivals by time t to this queue, $C_{i,j}^1(t)$ be its cumulative number of time slots assigned to this queue by time t , $q_{i,j}^1(t)$ be the number of packets queued at time t , and $B_{i,j}^1(t)$ be the cumulative number of departures by time t

from this queue. Since such a queue can be viewed as a work conserving link with a time varying capacity under the FCFS policy, one has the following well-known representation (see e.g. Ref. [2])

$$q_{i,j}^1(t) = \max_{0 \leq s \leq t} [A_{i,j}^1(t) - A_{i,j}^1(s) - C_{i,j}^1(t) + C_{i,j}^1(s)], \quad (7)$$

where

$$A_{i,j}^1(t) = \sum_{\ell=1}^{L_i} A_{i,\ell_j}^1(t). \quad (8)$$

Moreover,

$$B_{i,j}^1(t) = \min_{0 \leq s \leq t} [A_{i,j}^1(s) + C_{i,j}^1(t) - C_{i,j}^1(s)]. \quad (9)$$

In Lemma 1, we show that both the queue length and the delay at the load-balancing buffer can be bounded by finite constants. A similar result was previously shown in Ref. [9].

Lemma 1.

(i) The maximum queue length of the j th VOQ at the i th input port of the first stage is bounded above by L_i , i.e.

$$q_{i,j}^1(t) \leq L_i, \quad \text{for all } t. \quad (10)$$

(ii) The maximum delay for a packet to depart the j th VOQ at the i th input port of the first stage is bounded above by $(N - 1)L_i$.

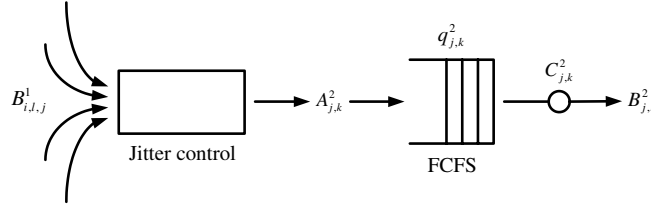
For the proof of Lemma 1, we need to introduce the following properties for the ceiling function and the floor function.

Proposition 1.

- (i) $\lceil a + b \rceil \leq \lceil a \rceil + \lceil b \rceil \leq \lceil a + b \rceil + 1$.
- (ii) $\lfloor a + b \rfloor \geq \lfloor a \rfloor + \lfloor b \rfloor$.
- (iii) $\lceil a \rceil \leq \lfloor a \rfloor + 1$.
- (iv) For an integer a , and $N > 0$,

$$\left\lceil \frac{a}{N} \right\rceil = \left\lfloor \frac{a + (N - 1)}{N} \right\rfloor.$$

Proof. Since the proofs of (i), (ii) and (iii) are trivial, we

Fig. 6. The k th VOQ at the j th input port of the second stage under FCFS.

only prove (iv). Without loss of generality, assume $\lceil \frac{a}{N} \rceil = P$ for some integer. Since a is an integer, we have

$$N(P-1) + 1 \leq a \leq NP.$$

Adding $(N-1)$ yields

$$NP \leq a + (N-1) \leq NP + (N-1).$$

This then implies

$$\left\lfloor \frac{NP}{N} \right\rfloor \leq \left\lfloor \frac{a + (N-1)}{N} \right\rfloor \leq \left\lfloor \frac{N(P+1) - 1}{N} \right\rfloor.$$

Thus,

$$\left\lfloor \frac{a + (N-1)}{N} \right\rfloor = P.$$

Proof (Proof of Lemma 1). (i) Note from Eqs. (8) and (5) that

$$\begin{aligned} A_{i,j}^1(t) - A_{i,j}^1(s) &= \sum_{\ell=1}^{L_i} (A_{i,\ell,j}^1(t) - A_{i,\ell,j}^1(s)) \\ &= \sum_{\ell=1}^{L_i} \left(\left\lfloor \frac{A_{i,\ell}(t) - j + 1}{N} \right\rfloor - \left\lfloor \frac{A_{i,\ell}(s) - j + 1}{N} \right\rfloor \right). \end{aligned}$$

Applying the first inequality in Proposition 1(i) yields

$$A_{i,j}^1(t) - A_{i,j}^1(s) \leq \sum_{\ell=1}^{L_i} \left\lfloor \frac{A_{i,\ell}(t) - A_{i,\ell}(s)}{N} \right\rfloor. \quad (11)$$

Since the connection patterns at the first stage are periodic with period N for some one-cycle permutation matrix, we have

$$C_{i,j}^1(t) - C_{i,j}^1(s) \geq \left\lfloor \frac{t-s}{N} \right\rfloor. \quad (12)$$

Observe that $\sum_{\ell=1}^{L_i} (A_{i,\ell}(t) - A_{i,\ell}(s))$ is the number of packet arrivals at the i th input port during the interval of length $t-s$. As there is at most one packet arrival at an input port per time slot, we have

$$\sum_{\ell=1}^{L_i} (A_{i,\ell}(t) - A_{i,\ell}(s)) \leq t-s. \quad (13)$$

In conjunction with Eq. (12),

$$\begin{aligned} C_{i,j}^1(t) - C_{i,j}^1(s) &\geq \left\lfloor \frac{\sum_{\ell=1}^{L_i} (A_{i,\ell}(t) - A_{i,\ell}(s))}{N} \right\rfloor \\ &\geq \sum_{\ell=1}^{L_i} \left\lfloor \frac{A_{i,\ell}(t) - A_{i,\ell}(s)}{N} \right\rfloor, \end{aligned} \quad (14)$$

where we use Proposition 1(ii) in the last inequality. From Eqs. (7), (11), and (14), and Proposition 1(iii),

$$\begin{aligned} q_{i,j}^1(t) &\leq \max_{0 \leq s \leq t} \left[\sum_{\ell=1}^{L_i} \left(\left\lfloor \frac{A_{i,\ell}(t) - A_{i,\ell}(s)}{N} \right\rfloor \right. \right. \\ &\quad \left. \left. - \left\lfloor \frac{A_{i,\ell}(t) - A_{i,\ell}(s)}{N} \right\rfloor \right) \right] \leq L_i. \end{aligned}$$

(ii) Since the scheduling policy at this queue is FCFS, it suffices to show that

$$B_{i,j}^1(t + (N-1)L_i) \geq A_{i,j}^1(t).$$

Note from Eq. (9) that

$$\begin{aligned} B_{i,j}^1(t + (N-1)L_i) - A_{i,j}^1(t) &= \min_{0 \leq s \leq t + (N-1)L_i} \left[A_{i,j}^1(s) - A_{i,j}^1(t) + C_{i,j}^1(t + (N-1)L_i) \right. \\ &\quad \left. - C_{i,j}^1(s) \right] \\ &= \min \left[\min_{0 \leq s \leq t} \left[C_{i,j}^1(t + (N-1)L_i) - C_{i,j}^1(s) - (A_{i,j}^1(t) - A_{i,j}^1(s)) \right] \right. \\ &\quad \left. - A_{i,j}^1(s) \right], \\ &= \min_{t+1 \leq s \leq t + (N-1)L_i} \left[C_{i,j}^1(t + (N-1)L_i) - C_{i,j}^1(s) - (A_{i,j}^1(t) - A_{i,j}^1(s)) \right]. \end{aligned}$$

All the terms in the second minimum are clearly nonnegative as both $A_{i,j}^1(t)$ and $C_{i,j}^1(t)$ are non-decreasing in t . On the other hand, for $0 \leq s \leq t$, we have from Eqs. (12), (13), and

(11) and Proposition 1(ii) and (iv) that

$$\begin{aligned}
& C_{ij}^1(t + (N-1)L_i) - C_{ij}^1(s) - (A_{ij}^1(t) - A_{ij}^1(s)) \\
& \geq \left\lfloor \frac{t + (N-1)L_i - s}{N} \right\rfloor - (A_{ij}^1(t) - A_{ij}^1(s)) \\
& \geq \left\lfloor \frac{\sum_{\ell=1}^{L_i} (A_{i,\ell}(t) - A_{i,\ell}(s)) + (N-1)L_i}{N} \right\rfloor \\
& \quad - \sum_{\ell=1}^{L_i} \left\lfloor \frac{A_{i,\ell}(t) - A_{i,\ell}(s)}{N} \right\rfloor \\
& = \left\lfloor \frac{\sum_{\ell=1}^{L_i} (A_{i,\ell}(t) - A_{i,\ell}(s) + (N-1))}{N} \right\rfloor \\
& \quad - \sum_{\ell=1}^{L_i} \left\lfloor \frac{A_{i,\ell}(t) - A_{i,\ell}(s)}{N} \right\rfloor \\
& \geq \sum_{\ell=1}^{L_i} \left(\left\lfloor \frac{A_{i,\ell}(t) - A_{i,\ell}(s) + (N-1)}{N} \right\rfloor \right. \\
& \quad \left. - \left\lfloor \frac{A_{i,\ell}(t) - A_{i,\ell}(s)}{N} \right\rfloor \right) = 0.
\end{aligned}$$

2.3. Analysis for the central buffer under FCFS

To ease the presentation, we let $d_{1,\max} = (N-1)L_{\max}$. Note from Lemma 1 that $d_{1,\max}$ is also the maximum delay at the first stage. As discussed in the introduction, a jitter control mechanism is added in the VOQ in front of the second stage. As the delay through the load-balancing buffer is bounded above by $d_{1,\max}$, packets with delay less than $d_{1,\max}$ are delayed to $d_{1,\max}$. By so doing, every packet has the same delay before entering the buffer at an input port of the second stage.

Now consider the k th VOQ at the j th input port of the second stage (see Fig. 6). Let $A_{j,k}^2(t)$ be the cumulative number of arrivals by time t to this queue, $C_{j,k}^2(t)$ be its cumulative number of time slots assigned to this queue by time t , $q_{j,k}^2(t)$ be the number of packets queued at time t , and $B_{j,k}^2(t)$ be the cumulative number of departures by time t from this queue. Since every packet has the same delay $d_{1,\max}$ through the first stage and fan-out splitting is done

at this stage, we have

$$A_{j,k}^2(t) = \sum_{(i,\ell) \in S^*(k)} A_{i,\ell,j}^1(t - d_{1,\max}) \quad (15)$$

(Here we use the convention that $A_{i,\ell,j}^1(\tau) = 0$ for $\tau < 0$). Also, as this queue is a work conserving link with a time varying capacity under the FCFS policy, one has

$$q_{j,k}^2(t) = \max_{0 \leq s \leq t} [A_{j,k}^2(t) - A_{j,k}^2(s) - (C_{j,k}^2(t) - C_{j,k}^2(s))], \quad (16)$$

and

$$B_{j,k}^2(t) = \min_{0 \leq s \leq t} [A_{j,k}^2(s) + C_{j,k}^2(t) - C_{j,k}^2(s)]. \quad (17)$$

Lemma 2. For the FCFS scheme,

- (i) $q_{j,k}^2(t) \leq [(q_k^o(t - d_{1,\max})/N) + M_k]$, for all t , and
- (ii) $B_{j,k}^2(t + q_k^o(t - d_{1,\max}) + (N-1)M_k) \geq A_{j,k}^2(t)$, for all t .

Lemma 2(i) provides an upper bound for the queue length in terms of the queue length of the corresponding output-buffered switch. As the scheduling policy is FCFS, Lemma 2(ii) implies that a packet that arrives at the queue at time t will depart not later than $t + q_k^o(t - d_{1,\max}) + (N-1)M_{\max}$.

Proof. (i) Note from Eqs. (15) and (5) and the inequality in Proposition 1(i) that

$$\begin{aligned}
A_{j,k}^2(t) - A_{j,k}^2(s) &= \sum_{(i,\ell) \in S^*(k)} (A_{i,\ell,j}^1(t - d_{1,\max}) - A_{i,\ell,j}^1(s - d_{1,\max})) \\
&\leq \sum_{(i,\ell) \in S^*(k)} \left\lfloor \frac{A_{i,\ell}(t - d_{1,\max}) - A_{i,\ell}(s - d_{1,\max})}{N} \right\rfloor \quad (18)
\end{aligned}$$

$$\leq \left\lfloor \frac{\sum_{(i,\ell) \in S^*(k)} (A_{i,\ell}(t - d_{1,\max}) - A_{i,\ell}(s - d_{1,\max}))}{N} \right\rfloor + M_k - 1. \quad (19)$$

Observe from Eqs. (1) and (2) that

$$\begin{aligned}
& \sum_{(i,\ell) \in S^*(k)} (A_{i,\ell}(t - d_{1,\max}) - A_{i,\ell}(s - d_{1,\max})) \\
&= A_k^o(t - d_{1,\max}) - A_k^o(s - d_{1,\max}) \leq q_k^o(t - d_{1,\max}) + (t - s). \quad (20)
\end{aligned}$$

Thus,

$$A_{j,k}^2(t) - A_{j,k}^2(s) \leq \left\lfloor \frac{q_k^o(t - d_{1,\max}) + (t - s)}{N} \right\rfloor + M_k - 1. \quad (21)$$

Since the connection patterns at the second stage are also periodic with period N for some one-cycle permutation matrix,

$$C_{j,k}^2(t) - C_{j,k}^2(s) \geq \left\lfloor \frac{t-s}{N} \right\rfloor. \quad (22)$$

From Eqs. (21), (22) and (16) and Proposition 1(iii), it then follows that

$$q_{j,k}^2(t) \leq \left\lceil \frac{q_k^0(t - d_{1,\max})}{N} \right\rceil + M_k.$$

for all t .

(ii) Let $d = q_k^0(t - d_{1,\max}) + (N-1)M_k$. Note from Eq. (17) that

$$\begin{aligned} B_{j,k}^2(t) - A_{j,k}^2(t) &= \min_{0 \leq s \leq t+d} [A_{j,k}^2(s) - A_{j,k}^2(t) + C_{ij}^2(t+d) - C_{ij}^2(s)] \\ &= \min \left[\min_{0 \leq s \leq t} [C_{j,k}^2(t+d) - C_{j,k}^2(s) - (A_{j,k}^2(t) - A_{j,k}^2(s))], \right. \\ &\quad \left. \min_{t+1 \leq s \leq t+d} [C_{j,k}^2(t+d) - C_{j,k}^2(s) - (A_{j,k}^2(t) - A_{j,k}^2(s))] \right] \end{aligned}$$

Clearly, all the terms in the second minimum are nonnegative as both $A_{j,k}^2(t)$ and $C_{j,k}^2(t)$ are nondecreasing in t . On the other hand, for $0 \leq s \leq t$, we have from Eq. (18) that

$$\begin{aligned} A_{j,k}^2(t) - A_{j,k}^2(s) &\leq \sum_{(i,\ell) \in S^*(k)} \left\lceil \frac{A_{i,\ell}(t - d_{1,\max}) - A_{i,\ell}(s - d_{1,\max})}{N} \right\rceil. \quad (23) \end{aligned}$$

Also, it follows from Eqs. (22) and (20), and Proposition 1(ii) that

$$\begin{aligned} C_{j,k}^2(t+d) - C_{j,k}^2(s) &\geq \left\lfloor \frac{t+d-s}{N} \right\rfloor \\ &= \left\lfloor \frac{q_k^0(t - d_{1,\max}) + (t-s) + (N-1)M_k}{N} \right\rfloor \\ &\geq \left\lfloor \frac{\sum_{(i,\ell) \in S^*(k)} (A_{i,\ell}(t - d_{1,\max}) - A_{i,\ell}(s - d_{1,\max})) + (N-1)M_k}{N} \right\rfloor \\ &= \left\lfloor \frac{\sum_{(i,\ell) \in S^*(k)} (A_{i,\ell}(t - d_{1,\max}) - A_{i,\ell}(s - d_{1,\max}) + (N-1))}{N} \right\rfloor \\ &\geq \sum_{(i,\ell) \in S^*(k)} \left\lfloor \frac{A_{j,k}^2(t - d_{1,\max}) - A_{j,k}^2(s - d_{1,\max}) + (N-1)}{N} \right\rfloor \quad (24) \end{aligned}$$

As in the proof of Lemma 1, we then have from Eqs. (23)

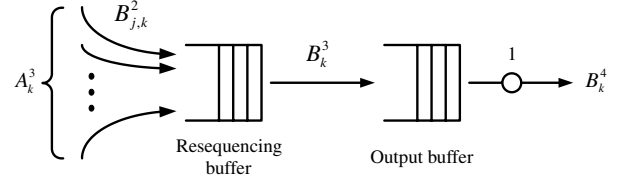


Fig. 7. The resequencing-and-output buffer.

and (24) and Proposition 1(iv) that

$$C_{j,k}^2(t+d) - C_{j,k}^2(s) - (A_{j,k}^2(t) - A_{j,k}^2(s)) \geq 0.$$

for $0 \leq s \leq t$.

2.4. Analysis for the resequencing-and-output buffer

In this section, we analyze the resequencing-and-output buffer. The resequencing-and-output buffer conceptually consists of two virtual buffers (see Fig. 7): (i) the resequencing buffer and (ii) the output buffer. The objective of the resequencing buffer is to reorder the packets so that packets of the same flow depart in the same order as they arrive. After resequencing, packets are stored in the output buffer waiting for transmission from the output link. Let $A_k^3(t)$ be the cumulative arrivals by time t to the k th resequencing buffer, and $B_k^3(t)$ be its cumulative departures. Note that $B_k^3(t)$ is also the cumulative arrivals by time t to the k th output buffer. Let $B_k^4(t)$ be the cumulative departures by time t from the k th output buffer. Clearly, from the input-output relation of a work-conserving link, we have

$$B_k^4(t) = \min_{0 \leq s \leq t} [B_k^3(s) + (t-s)]. \quad (25)$$

Lemma 3. *The following results hold for the FCFS scheme.*

- (i) $A_k^3(t) \leq B_k^0(t - d_{1,\max})$.
- (ii) $B_k^3(t + d_{1,\max} + NM_k) \geq B_k^0(t)$.
- (iii) $B_k^4(t) \geq B_k^0(t - d_{1,\max} - NM_k)$.
- (iv) *The number of packets queued at the k th resequencing-and-output buffer is bounded above by NM_k , i.e.*

$$A_k^3(t) - B_k^4(t) \leq NM_k, \text{ for all } t, \quad (26)$$

Proof. (i) Note that

$$A_k^3(t) = \sum_{j=1}^N B_{j,k}^2(t).$$

From Eq. (17), it follows that

$$\begin{aligned} A_k^3(t) &= \sum_{j=1}^N \min_{0 \leq s \leq t} [A_{j,k}^2(s) + C_{j,k}^2(t) - C_{j,k}^2(s)] \\ &\leq \min_{0 \leq s \leq t} \left[\sum_{j=1}^N A_{j,k}^2(s) + C_{j,k}^2(t) - C_{j,k}^2(s) \right]. \end{aligned}$$

As the connection patterns at the second stage are periodic

with period N for some one-cycle permutation matrix,

$$\sum_{j=1}^N C_{j,k}^2(t) = t.$$

Thus,

$$A_k^3(t) \leq \min_{0 \leq s \leq t} \left[\sum_{j=1}^N A_{j,k}^2(s) + t - s \right]. \quad (27)$$

Note from Eqs. (15), (6), and (2) that for all s

$$\begin{aligned} \sum_{j=1}^N A_{j,k}^2(s) &= \sum_{j=1}^N \sum_{(i,\ell) \in S^*(k)} A_{i,\ell}^1(s - d_{1,\max}) \\ &= \sum_{(i,\ell) \in S^*(k)} A_{i,\ell}(s - d_{1,\max}) = A_k^o(s - d_{1,\max}). \end{aligned}$$

Using this and Eq. (3) in Eq. (27) yields

$$\begin{aligned} A_k^3(t) &\leq \min_{0 \leq s \leq t} [A_k^o(s - d_{1,\max}) + t - s] \\ &= \min_{d_{1,\max} \leq s \leq t} [A_k^o(s - d_{1,\max}) + t - s] \\ &= \min_{0 \leq \tau \leq t - d_{1,\max}} [A_k^o(\tau) + t - d_{1,\max} - \tau] = B_k^o(t - d_{1,\max}) \end{aligned}$$

(ii) Consider a packet that is destined for the k th output port. Without loss of generality, suppose that this packet arrives at time t and it is routed through the k th VOQ at j th input port of the second stage. From 3 (ii), the packet leaves the first stage not later than $t + d_{1,\max}$. After the jitter control, it arrives at the j th input port of the second stage exactly at $t + d_{1,\max}$. Since FCFS is used in the k th VOQ at the j th input port of the second stage, we have from Lemma 2(ii) that the packet leaves the second stage not later than $t + q_k^o(t) + (N-1)M_k + d_{1,\max}$. As the bound is independent of j , we conclude that every packet that arrives at the first stage at time t will depart from the second stage not later than $t + q_k^o(t) + (N-1)M_k + d_{1,\max}$. Note from Eq. (1) that $t + q_k^o(t)$ is non-decreasing in t . Thus, any other packets that are destined for the k th output port and arrives before t leave the second stage not later than $t + q_k^o(t) + (N-1)M_k + d_{1,\max}$. This implies the packet (that arrives at the first stage at time t and destined for the k th output port) leaves the resequencing buffer not later than $t + q_k^o(t) + (N-1)M_k + d_{1,\max}$.

On the other hand, since the packet that arrives at time t departs from the corresponding output-buffered switch between $t + q_k^o(t) - M_k + 1$ and $t + q_k^o(t)$, the departure time for a packet to leave the resequencing buffer is not later than the sum of that from the corresponding output-buffer switch and $d_{1,\max} + NM_k$. This shows that for all t

$$B_k^3(t + d_{1,\max} + NM_k) \geq B_k^o(t).$$

(iii) From Eq. (25) and (ii) of this lemma, it follows that

$$B_k^4(t) \geq \min_{0 \leq s \leq t} [B_k^o(s - d_{1,\max} - NM_k) + t - s].$$

Since $B_k^o(t) - B_k^o(s) \leq t - s$ in Eq. (4), the above minimum occurs at $s = t$. Thus,

$$B_k^4(t) \geq B_k^o(t - d_{1,\max} - NM_k).$$

(iv) From (i) and (iii) of this lemma and Eq. (4), it follows that

$$\begin{aligned} A_k^3(t) - B_k^4(t) &\leq B_k^o(t - d_{1,\max}) - B_k^o(t - d_{1,\max} - NM_k) \\ &\leq NM_k. \end{aligned}$$

Proof (Proof of Theorem 1).

- (i) This is a direct consequence of Lemma 3(iii).
- (ii) Since there are N VOQ at each input port, the result then follows from Lemma 1(i).
- (iii) This is also shown in Lemma 1(ii).
- (iv) It is shown in Lemma 3(iv).

3. The EDF scheme

In this section, we prove Theorem 2. Consider the same traffic model as in Section 2. Also, let $B_{i,\ell}^o(t)$ be the cumulative departures of the $A_{i,\ell}$ -flow by time t from the corresponding FCFS output-buffered switch (see Fig. 4). Since the EDF scheme uses packet departure times as deadlines, $B_{i,\ell}^o(t)$ is also the number of packets from the $A_{i,\ell}$ -flow that have deadlines not greater than t . Note that

$$B_k^o(t) = \sum_{(i,\ell) \in S^*(k)} B_{i,\ell}^o(t). \quad (28)$$

We first establish the following inequalities that will be used in the proof of Theorem 2.

Proposition 2. For the FCFS output-buffered switch,

$$\sum_{(i,\ell) \in S(k)} B_{i,\ell}^o(t) \leq \min_{0 \leq s \leq t} \left[\sum_{(i,\ell) \in S(k)} A_{i,\ell}(s) + t - s \right], \quad (29)$$

for all t and for all $S(k)$ that is a subset of $S^*(k)$.

Proof. If one assigns priority to the set of traffic from $S(k)$ in the output-buffered switch, then the right hand side of Eq. (29) is the cumulative departures by time t for this set of traffic (cf. Eq. (3)). This should be larger than or equal to the left hand of Eq. (29), which is the cumulative departures by time t for this set of traffic under FCFS.

We also note that results in Section 2.2 can be directly applied as the operations for the load-balancing buffer in

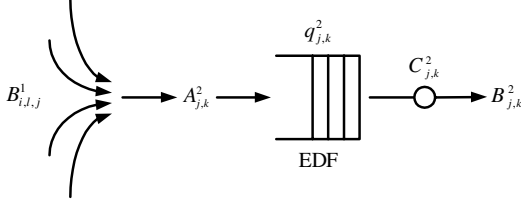


Fig. 8. The k th VOQ at the j th input port of the second stage under EDF.

front of the first stage in the EDF scheme is exactly the same as that in Section 2.2.

3.1. Analysis for the central buffer under EDF

Now consider the k th VOQ at the j th input port of the second stage (see Fig. 8). As in Section 2.3, let $A_{j,k}^2(t)$ be the cumulative number of arrivals by time t to this queue and $C_{j,k}^2(t)$ be its cumulative number of time slots assigned to this queue by time t .

Also, let $B_{i,\ell,j}^1$ be the departure process of $A_{i,\ell,j}^1$ from the j th VOQ at the i th input port of the first stage. Since there is no jitter control in the EDF scheme, the departure process from the first stage is simply the arrival process at the second stage. Thus,

$$A_{j,k}^2(t) = \sum_{(i,\ell) \in S^*(k)} B_{i,\ell,j}^1(t). \quad (30)$$

Lemma 4. Every packet leaves the k th VOQ at the j th input port of the second stage not later than the sum of its deadline and $(N-1)(L_{\max} + M_k)$.

Note that if we add the deadline of every packet by the same constant $(N-1)(L_{\max} + M_k)$, then Lemma 4 is equivalent to that every packet leaves the k th VOQ at the j th input port of the second stage not later than its deadline.

Proof. Let $d_1 = (N-1)L_{\max}$ and $d_2 = (N-1)M_k$. Let $D_{i,\ell,j}(t)$ be the number of packets of the $A_{i,\ell,j}^1$ -flow that have deadlines not greater than t . From [2], Theorem 5.6.1, it suffices to show that

$$\begin{aligned} & \sum_{(i,\ell) \in S(k)} D_{i,\ell,j}(t - d_1 - d_2) \\ & \leq \min_{0 \leq s \leq t} \left[\sum_{(i,\ell) \in S(k)} B_{i,\ell,j}^1(s) + C_{j,k}^2(t) - C_{j,k}^2(s) \right], \end{aligned} \quad (31)$$

for all t and for all $S(k)$ that is a subset of $S^*(k)$.

As the delay through the j th VOQ at an input port of the first stage is bounded above by $d_1 = (N-1)L_{\max}$ in Lemma

1(ii), we then have

$$\begin{aligned} & \min_{0 \leq s \leq t} \left[\sum_{(i,\ell) \in S(k)} B_{i,\ell,j}^1(s) + C_{j,k}^2(t) - C_{j,k}^2(s) \right] \\ & \geq \min_{0 \leq s \leq t} \left[\sum_{(i,\ell) \in S(k)} A_{i,\ell,j}^1(s - d_1) + C_{j,k}^2(t) - C_{j,k}^2(s) \right] \\ & = \min_{0 \leq \tau \leq t - d_1} \left[\sum_{(i,\ell) \in S(k)} A_{i,\ell,j}^1(\tau) + C_{j,k}^2(t) - C_{j,k}^2(\tau + d_1) \right]. \end{aligned}$$

Using Eqs. (5) and (22), one has

$$\begin{aligned} & \min_{0 \leq s \leq t} \left[\sum_{(i,\ell) \in S(k)} B_{i,\ell,j}^1(s) + C_{j,k}^2(t) - C_{j,k}^2(s) \right] \\ & \geq \min_{0 \leq \tau \leq t - d_1} \left[\sum_{(i,\ell) \in S(k)} \left\lceil \frac{A_{i,\ell}(\tau) - j + 1}{N} \right\rceil + \left\lfloor \frac{t - d_1 - \tau}{N} \right\rfloor \right]. \end{aligned} \quad (32)$$

On the other hand, since we use the departure times from the corresponding output-buffered switch as deadlines and packets are split in the round robin fashion, we have (cf. Eq. (5))

$$D_{i,\ell,j}(t) = \left\lceil \frac{B_{i,\ell}^o(t) - j + 1}{N} \right\rceil. \quad (33)$$

In conjunction with Eq. (32), the inequalities in Eq. (31) hold if we can show that for all $0 \leq \tau \leq t - d_1$,

$$\begin{aligned} & \sum_{(i,\ell) \in S(k)} \left\lceil \frac{B_{i,\ell}^o(t - d_1 - d_2) - j + 1}{N} \right\rceil \\ & \leq \sum_{(i,\ell) \in S(k)} \left\lceil \frac{A_{i,\ell}(\tau) - j + 1}{N} \right\rceil + \left\lfloor \frac{t - d_1 - \tau}{N} \right\rfloor. \end{aligned} \quad (34)$$

Observe from Proposition 2 that

$$\sum_{(i,\ell) \in S(k)} B_{i,\ell}^o(t - d_1 - d_2) \leq \sum_{(i,\ell) \in S(k)} A_{i,\ell}(\tau) + t - d_1 - d_2 - \tau, \quad (35)$$

for all τ , t and all subsets $S(k)$. In particular, choosing $\tau = t - d_1 - d_2$ and a subset that contains only one flow yields

$$A_{i,\ell}(t - d_1 - d_2) \geq B_{i,\ell}^o(t - d_1 - d_2),$$

for all $(i,\ell) \in S^*(k)$. Thus, for $t - d_1 - d_2 + 1 \leq \tau \leq t - d_1$, one has

$$A_{i,\ell}(\tau) \geq A_{i,\ell}(t - d_1 - d_2) \geq B_{i,\ell}^o(t - d_1 - d_2),$$

and the inequality in Eq. (34) is satisfied trivially. For

$$0 \leq \tau \leq t - d_1 - d_2,$$

we have from the inequality Eq. (35) that

$$t - d_1 - \tau \geq \sum_{(i,\ell) \in S(k)} (B_{i,\ell}^o(t - d_1 - d_2) - A_{i,\ell}(\tau)) + d_2.$$

Recall that

$$d_2 = (N - 1)M_k.$$

Also, we have $|S(k)| \leq M_k$ as $S(k)$ is a subset of $S^*(k)$. Thus,

$$\begin{aligned} & \left\lfloor \frac{t - d_1 - \tau}{N} \right\rfloor \\ & \geq \left\lfloor \frac{\sum_{(i,\ell) \in S(k)} (B_{i,\ell}^o(t - d_1 - d_2) - A_{i,\ell}(\tau)) + (N - 1)M_k}{N} \right\rfloor \\ & \geq \left\lfloor \frac{\sum_{(i,\ell) \in S(k)} (B_{i,\ell}^o(t - d_1 - d_2) - A_{i,\ell}(\tau)) + (N - 1)|S(k)|}{N} \right\rfloor \\ & = \left\lfloor \frac{\sum_{(i,\ell) \in S(k)} (B_{i,\ell}^o(t - d_1 - d_2) - A_{i,\ell}(\tau)) + (N - 1)}{N} \right\rfloor. \end{aligned}$$

Applying Proposition 1(ii) yields

$$\begin{aligned} & \left\lfloor \frac{t - d_1 - \tau}{N} \right\rfloor \\ & \geq \sum_{(i,\ell) \in S(k)} \left\lfloor \frac{B_{i,\ell}^o(t - d_1 - d_2) - A_{i,\ell}(\tau) + (N - 1)}{N} \right\rfloor. \end{aligned} \quad (36)$$

On the other hand, we also have from Proposition 1(i) that

$$\begin{aligned} & \sum_{(i,\ell) \in S(k)} \left\lfloor \frac{B_{i,\ell}^o(t - d_1 - d_2) - j + 1}{N} \right\rfloor \\ & - \sum_{(i,\ell) \in S(k)} \left\lfloor \frac{A_{i,\ell}(\tau) - j + 1}{N} \right\rfloor \\ & \leq \sum_{(i,\ell) \in S(k)} \left\lfloor \frac{B_{i,\ell}^o(t - d_1 - d_2) - A_{i,\ell}(\tau)}{N} \right\rfloor. \end{aligned} \quad (37)$$

The inequalities in Eq. (34) then follows from Eqs. (36) and (37) and Proposition 1(iv).

3.2. Analysis for the resequencing-and-output buffer under EDF

In this section, we analyze the resequencing-and-output buffer under the EDF scheme. As in Section 2.4, let $A_k^3(t)$ be the cumulative arrivals by time t to the k th resequencing buffer, $B_k^3(t)$ be its cumulative departures, and $B_k^4(t)$ be the

cumulative departures by time t from the k th output buffer (see Fig. 7).

Lemma 5. *The following results hold for the EDF scheme.*

- (i) $A_k^3(t) \leq B_k^o(t)$.
- (ii) $B_k^3(t + (N - 1)(L_{\max} + M_k)) \geq B_k^o(t)$.
- (iii) $B_k^4(t) \geq B_k^o(t - (N - 1)(L_{\max} + M_k))$.
- (iv) *The number of packets queued at the k th resequencing-and-output buffer is bounded above by $(N - 1)(L_{\max} + M_k)$, i.e.*

$$A_k^3(t) - B_k^4(t) \leq (N - 1)(L_{\max} + M_k), \quad \text{for all } t. \quad (38)$$

Proof.

- (i) Note from Eq. (30), $B_{i,\ell_j}^1(s) \leq A_{i,\ell_j}^1(s)$, Eqs. (6) and (2) that for all s

$$\begin{aligned} \sum_{j=1}^N A_{j,k}^2(s) &= \sum_{j=1}^N \sum_{(i,\ell) \in S^*(k)} B_{i,\ell_j}^1(s) \leq \sum_{j=1}^N \sum_{(i,\ell) \in S^*(k)} A_{i,\ell_j}^1(s) \\ &= \sum_{(i,\ell) \in S^*(k)} A_{i,\ell}(s) = A_k^o(s). \end{aligned}$$

As the EDF scheme is work conserving, Eq. (27) in the proof of Lemma 3 is still applicable. It then follows from Eq. (3) that

$$A_k^3(t) \leq \min_{0 \leq s \leq t} [A_k^o(s) + t - s] = B_k^o(t).$$

- (ii) Since we use the departure times from the corresponding output-buffered switch as deadlines, the result then follows from Lemma 4.

- (iii) This follows the same argument in Lemma 3(iii).

- (iv) From (i) and (iii) of this lemma and Eq. (4), it follows that

$$\begin{aligned} A_k^3(t) - B_k^4(t) &\leq B_k^o(t) - B_k^o(t - (N - 1)(L_{\max} + M_k)) \\ &\leq (N - 1)(L_{\max} + M_k). \end{aligned}$$

Proof (Proof of Theorem 2). (i) This is a direct consequence of Lemma 5(iii). (ii) It is shown in Lemma 5(iv).

4. Conclusions

In this paper, we showed that resequencing can be done efficiently by distributing packets according to their flows in the load balanced Birkhoff–von Neumann switches. This is done by adding a load-balancing buffer in front of the first stage and a resequencing-and-output buffer after the second stage. We considered two scheduling policies: the FCFS policy and the EDF policy. For both schemes, we showed that the end-to-end delay

is bounded by the sum of the delay from the corresponding FCFS output-buffered switch and a constant that only depends on the number of flows and the size of the switch. The FCFS policy requires jitter control in front of the second stage so that packets entering the second stage have the same delay. Its scheduling complexity is simpler than that of the EDF scheme. The EDF policy uses the departure times of the corresponding output-buffered switch as deadlines and schedules packets according to their deadlines. As such, there is no need for jitter control. Since computing departure times of the corresponding output-buffered switch requires the global information of all the inputs, we also suggested using the arrival times as deadlines for the EDF scheme. This reduces the scheduling complexity at the cost of slightly increasing the delay bound and the size of the resequencing-and-output buffer.

There are several problems that require further study.

(i) We only looked at the worst case of the end-to-end delay. It is of the same importance to study the average delay. In particular, the jitter control mechanism in the FCFS scheme that delays every packet to its maximum (worst case) delay is expected to perform poorly in the light load (the main reason that we have jitter control in the FCFS scheme is for the sake of the proof for the worst case bound). Our preliminary simulation study shows that the average delay with jitter control is mainly dominated by the delay in the jitter control even when the load is moderate.

(ii) Though we only compared our schemes with the FCFS output-buffered switches, we expect that our approaches could also be used for other work conserving scheduling policies, such as the service curve earliest deadline first (SCED) policy in Ref. [7] and the generalized processor sharing (GPS) policy in Ref. [11]. For both SCED and GPS, it is well-known that quality of service (QoS) guarantees can be achieved in output-buffered switches. If there is a bound for the difference between our multi-stage switch and the corresponding output-buffered switch under these two policies, then QoS guarantees can also be achieved in our multi-stage switch. Research along this line will be reported separately.

(iii) In this paper, packets are assumed to be of the same size. To allow variable length packets, one needs to introduce packetizers (see e.g. Ref. [2]) for segmentation and reassembly. The bounds in Ref. [2] for packetizers might be used to extend the results in this paper to the case with variable length packets.

(iv) In the recent paper by Keslassy and McKeown [10], a ‘full frame first’ algorithm is proposed so that packets in the two-stage switches can be departed in sequence. However, this is at the cost of communications between VOQs. It would be interesting to see if there are other scheduling algorithms that achieve the same objective with lower communication cost.

Acknowledgements

This research is supported in part by the National Science Council, Taiwan, ROC, under Contracts NSC-89-2213-E007-127, NSC-89-2219-E007-015 and the program for promoting academic excellence of universities 89-E-FA04-1-4.

References

- [1] G. Birkhoff, Tres observaciones sobre el algebra lineal, Univ. Nac. Tucumán Rev. Ser. A 5 (1946) 147–151.
- [2] C.S. Chang, Performance Guarantees in Communication Networks, Springer, London, 2000.
- [3] C.S. Chang, W.J. Chen, H.Y. Huang, On service guarantees for input buffered crossbar switches: a capacity decomposition approach by Birkhoff and von Neumann, IEEE IWQoS'99, pp. 79–86, London, UK, 1999 (US patent pending).
- [4] C.S. Chang, W.J. Chen, H.Y. Huang, Birkhoff–von Neumann input buffered crossbar switches, IEEE INFOCOM2000, pp. 1614–1623, Tel Aviv, Israel, 2000.
- [5] C.S. Chang, D.S. Lee, Y.S. Jou, Load balanced Birkhoff–von Neumann switches, part I: one-stage buffering, Computer Communications, special issue on Current Issues in Terabit Switching, 2001, in press.
- [6] S.-T. Chuang, A. Goel, N. McKeown, B. Prabhkar, Matching output queuing with a combined input output queued switch, IEEE INFOCOM'99, pp. 1169–1178, New York, 1999.
- [7] R.L. Cruz, Quality of service guarantees in virtual circuit switched networks, IEEE J. Selected Areas Commun. August (1995) special issue on Advances in the Fundamentals of Networking.
- [8] J.Y. Hui, T. Renner, Queueing strategies for multicast packet switching, IEEE GLOBECOM'90 3 (1990) 1431–1437.
- [9] S. Iyer, N. McKeown, Making parallel packet switch practical, Proceedings of IEEE INFOCOM 2001, Anchorage, Alaska, USA.
- [10] I. Keslassy, N. McKeown, Maintaining packet order in two-stage switches, preprint, 2001.
- [11] A.K. Parekh, R.G. Gallager, A generalized processor sharing approach to flow control in integrated service networks: the single-node case, IEEE/ACM Trans. Network. 1 (1993) 344–357.
- [12] I. Stoica, H. Zhang, Exact emulation of an output queuing switch by a combined input output queuing switch, IEEE IWQoS'98, pp. 218–224, Napa, California, 1998.
- [13] J. von Neumann, A certain zero-sum two-person game equivalent to the optimal assignment problem, Contributions to the Theory of Games, vol. 2, Princeton University Press, Princeton, NJ, 1953 pp. 5–12.



Cheng-Shang Chang received the B.S. degree from the National Taiwan University, Taipei, Taiwan, in 1983, and the M.S. and Ph.D. degrees from Columbia University, New York, NY, in 1986 and 1989, respectively, all in Electrical Engineering. From 1989 to 1993 he was employed as a Research Staff Member at the IBM Thomas J. Watson Research Center, Yorktown Heights, N.Y. In 1993, he joined the Department of Electrical Engineering at National Tsing Hua University, Taiwan, R.O.C., where he is a Professor. His current research interests are concerned with queueing theory, stochastic scheduling and performance evaluation of telecommunication networks and parallel processing systems. Dr. Chang received the IBM Outstanding Innovation Award in 1992, and the Outstanding Research Award from the National Science Council, Taiwan, in 1999 and 2001. He is the author of the book “Performance Guarantees in Communication Networks”, and he served as an editor for Operations Research from 1992 to 1999.



Duan-Shin Lee was born in Taiwan in 1961. He received the B.S. degree from National Tsing Hwa University, Taiwan, in 1983, and the M.S. and Ph.D. degrees from Columbia University, New York, in 1987 and 1990, all in electrical engineering. He worked as a research staff member at the C&C Research Laboratory of NEC USA, Inc. in Princeton, New Jersey from 1990 to 1998. Since 1998, he has been an associate professor in the Department of Computer Science of National Tsing Hwa University in Hsinchu, Taiwan. His research interests are switch and router design, personal communication systems, performance analysis computer networks and queueing theory. Dr. Lee is a senior member of IEEE.



Ching-Ming Lien received the B.S. and M.S. degrees from National Tsing Hua University, Taiwan, both in electrical engineering, in 1999 and 2001, respectively. He was the administrator of the Bulletin Board System of the Department of Electrical Engineering and worked as a free lancer of primary news media in Taiwan when he was an undergraduate student in the National Tsing Hua University. His graduate research is focused on modeling and performance analysis of high-speed networks, with emphasis on the architecture of switches.