# Can We Schedule Traffic More Efficiently in Optical Packet Switches?

Bin Wu, *Student Member, IEEE*, Xin Wang, and Kwan L. Yeung, *Senior Member, IEEE*

*Abstract*—We consider traffic scheduling in non-blocking electronic-buffered optical packet switches (OPS) with bounded packet delay. Due to the reconfiguration overhead of the switch fabric, the two commonly used optimization objectives, minimizing packet delay and minimizing switch speedup, conflict with each other. Intelligent scheduling algorithms have been designed to provide tradeoff between these two objectives. In this paper, we propose a more efficient approach to schedule OPS traffic, resulting in significantly reduced speedup and/or packet delay. However, our approach is based on a very interesting conjecture, which has not been strictly proved so far. We would like to put forward this conjecture as an open question, and call for a proof or disproof.

*Index Terms*—Conjecture, optical packet switch (OPS), performance guaranteed switching, scheduling.

Fig. 1. A scalable high speed optical packet switch.

## I. INTRODUCTION

RECENT progress on optical switching technologies [1-4] has enabled the implementation of electronic-buffered optical packet switches (OPS) as shown in Fig. 1. The core of this architecture is the optical switch fabric, which can efficiently provide huge switching capacity as demanded by the backbone routers in the Internet. Since optical connections (i.e. optical fibers) are used to interconnect the input/output line-cards with the central switch fabric, the input/output line-cards can be distributed into several racks, which may locate at hundreds of meters away from each other. As a result, power consumption in each rack can be reduced, and the switch becomes more scalable.

On the other hand, the optical switch fabric usually needs some guard time to change its inter-connection pattern from one to another, and to synchronize the signals arriving at the input ports [5]. This guard time is called *reconfiguration overhead*. During this period, no packet can be transmitted across the switch fabric. Accordingly, packet transmission rate in the switch fabric must be faster than the external line-rate (i.e. a *speedup* is required in the switch fabric) in order to achieve *performance guaranteed switching* (i.e. 100% throughput with bounded packet delay) [6-12]. It is shown in [6, 8] that minimizing speedup and minimizing delay are two conflicting goals, where a higher speedup gives a shorter packet delay and vice versa.
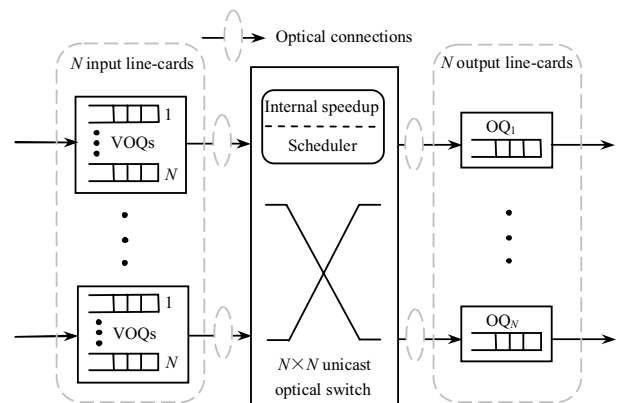
Based on switch architectures similar to Fig. 1, several algorithms have been recently proposed [6-12] to schedule OPS traffic with guaranteed switching performance. Among them, MIN [6], $\alpha^i$-SCALE [9] and QLEF [10] aim primarily at minimizing the packet delay, whereas reducing speedup is a secondary objective. DOUBLE [6] is the first algorithm that allows the tradeoff between speedup and delay. Let $N$ denote the switch size and $N_S$ be the (maximum) number of switch configurations required for scheduling. DOUBLE needs no more than $N_S=2N$ configurations to schedule any *legitimate* traffic matrix, with a speedup of $S_{schedule}=2$ (detailed in Section II). However, this algorithm does not consider the amount of reconfiguration overhead $\delta$ in its scheduling decision, and thus it is not optimized for switches with different $\delta$. Besides, $N_S=2N$ only represents a single point in the solution space [6, 8], and the characteristics for other $N_S$ values are not studied in [6]. To address those issues, ADAPTIVE [8] is proposed. It is shown [8] that DOUBLE can be regarded as a special case of ADAPTIVE at $N_S=2N$.

In this paper, we explore the possibility of beating the performance of DOUBLE and ADAPTIVE. We show that this can be achieved based on a very interesting conjecture. We put forward this conjecture as an open question, and hope that a proof or disproof can be found soon.

The rest of the paper is organized as follows. In Section II, we review the generic scheduling procedure [6-12] and the existing scheduling algorithms DOUBLE and ADAPTIVE. In

Bin Wu, Xin Wang and Kwan L. Yeung are with the Department of Electrical and Electronic Engineering, The University of Hong Kong, Pokfulam, Hong Kong (Tel: 852-2857-8493; Fax: 852-2559-8738; e-mail: {binwu, xinwang, kyeung}@eee.hku.hk).
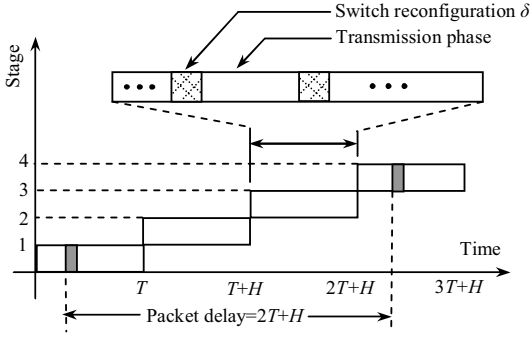
Fig. 2. Optical packet switch scheduling stages.

Section III, we propose an approach to improve the scheduling efficiency. Section IV gives some further discussion, and we conclude the paper in Section V.

## II. TRAFFIC SCHEDULING AND SPEEDUP-DELAY TRADEOFF

### A. Scheduling Procedure

The generic four-stage scheduling procedure as shown in Fig. 2 is followed. In Stage 1, incoming packets are periodically accumulated in the input buffers over $T$ time slots to construct an $N \times N$ traffic matrix $C(T) = \{c_{ij}\}$. Each entry $c_{ij}$ denotes the number of packets received at input $i$ and destined to output $j$. $C(T)$ is *legitimate* if each of its line sums (i.e. row sum or column sum) is no larger than $T$. Throughout the paper, we only consider legitimate $C(T)$. The scheduling algorithm takes $H$ time slots in Stage 2 to generate $N_S$ configurations $P_n = \{p^{(n)}_{ij}\}$, $n \in \{1, \ldots, N_S\}$, each weighted by $\phi_n$ to cover $C(T)$. "Cover" means that $\sum_{n=1}^{N_S} \phi_n p^{(n)}_{ij} \geq c_{ij}$ for any $i, j \in \{1, \ldots, N\}$. $P_n$ is an $N \times N$ permutation matrix with at most a single "1" in each line (row or column). $p^{(n)}_{ij} = 1$ indicates that a packet can be sent from input $i$ to output $j$ in one slot; $p^{(n)}_{ij} = 0$ otherwise. In Stage 3, the switch fabric is reconfigured according to the $N_S$ configurations obtained in Stage 2. An internal switch fabric speedup $S$ is applied, resulting in compressed/shortened time slots, to ensure that this stage occupies only $T$ (regular) slots. The fabric holds each configuration $P_n$ for $\phi_n$ *compressed slots* for packet transmission. Finally in Stage 4 packets are sent onto the output lines from output buffers (in $T$ slots).

From the tagged packet in Fig. 2, we can see that the bounded delay of any packet is $2T+H$ slots. Because $\delta N_S$ slots are used to reconfigure the switch for $N_S$ times in Stage 3, only $T-\delta N_S$ slots are left for transmitting $C(T)$. Since there are at most $T$ packets waiting at each input port for transmission, a speedup factor $S_{\text{reconfigure}} = T/(T-\delta N_S)$ is necessary to compensate for the idle time caused by reconfigurations. At the same time, the scheduling algorithm may produce some empty slots (i.e. underutilize the bandwidth provided by the configurations [6-12]). As a result, more than $T$ compressed slots are usually needed in Stage 3 to transmit $C(T)$. Therefore another speedup factor

$$S_{\text{schedule}} = \frac{1}{T}\sum_{n=1}^{N_S} \phi_n \qquad (1)$$

is required to compensate for the inefficient scheduling. In fact, $S_{\text{schedule}}$ denotes the efficiency of the scheduling algorithm adopted. A smaller $S_{\text{schedule}}$ indicates a more efficient scheduler (with less empty slots in the schedule). The overall internal speedup $S$ is then given by

$$S = S_{\text{reconfigure}} \times S_{\text{schedule}} = \frac{T}{T-\delta N_S} S_{\text{schedule}} = \frac{1}{T-\delta N_S}\sum_{n=1}^{N_S} \phi_n . \quad (2)$$

### B. Speedup-Delay Tradeoff and Scheduling Algorithms

For a given $C(T)$, we divide it by $T/(N_S-N)$ [1] to get a quotient matrix $Q = \{q_{ij}\}$ and a residue matrix $R = \{r_{ij}\}$:

$$C(T) = \frac{T}{N_S - N} \times Q + R . \qquad (3)$$

Since each line of $C(T)$ sums to at most $T$, the maximum line sum of $Q$ is at most $N_S-N$. So we can apply edge-coloring [13] to the bipartite multigraph of $Q$, and get $N_S-N$ configurations to cover $Q$ [6-8]. On the other hand, all the entries in $R$ are not larger than $T/(N_S-N)$. So, $R$ can be covered by *any N* non-overlapping configurations (i.e. any two of them do not cover the same entry), with each weighted by $T/(N_S-N)$. All in all, $C(T)$ can be covered by $(N_S-N)+N=N_S$ configurations, each weighted by $\phi_n = T/(N_S-N)$. From (1), $S_{\text{schedule}}$ can be found.

$$S_{\text{schedule}} = \frac{1}{T}\sum_{n=1}^{N_S} \phi_n = \frac{1}{T} \times \frac{T}{N_S-N} \times N_S = 1 + \frac{N}{N_S-N} . \qquad (4)$$

The above formula (4) is referred to as *speedup function* in [8]. In essence, it depicts the tradeoff relationship between the speedup ($S_{\text{schedule}}$ due to the inefficient scheduling) and the delay (in terms of $N_S$). Recall that the bounded delay of any packet is $2T+H$ slots and $T>\delta N_S$. Therefore the minimum achievable delay is given by $2\delta N_S+H$.

DOUBLE [6] requires $N_S=2N$ configurations to cover $C(T)$. This is obtained by replacing $N_S$ in (3) by $2N$ to get $C(T)=[T/N]\times Q+R$. $N$ configurations are required to cover $Q$ and $R$ respectively, and each configuration is equally weighted by $\phi_n=T/N$. From (4), DOUBLE achieves $S_{\text{schedule}}=2$. Fig. 3 gives an example of DOUBLE execution.

Unlike DOUBLE, ADAPTIVE [8] substitutes (4) into (2), and minimizes the overall speedup $S$ by solving

$$\frac{\partial S}{\partial N_S} = 0 . \qquad (5)$$

Therefore the schedule generated by ADAPTIVE is optimized with respect to the value of $\delta$.

## III. IMPROVING SCHEDULING EFFICIENCY

In this section, we aim at achieving a better scheduling efficiency than DOUBLE and ADAPTIVE. Since DOUBLE is a special case of ADAPTIVE at $N_S=2N$, for simplicity, we only focus on DOUBLE below. We further assume that the switch

---

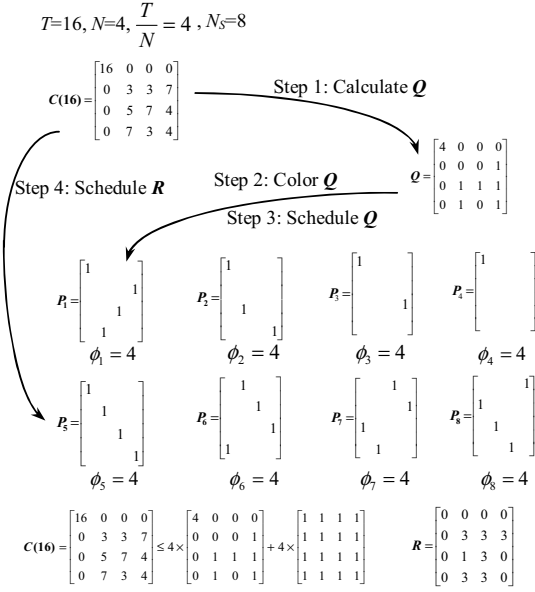[1] If $T/(N_S-N)$ is not an integer, use $\lceil T/(N_S-N) \rceil$ as the substitute [8].

Fig. 3. An example of DOUBLE execution. The all-1 matrix used to cover *R* equals to the sum of the *N* non-overlapping configurations (*P₅-P₈*).



Fig. 4. Move the circled LERs from *R* to *Q*.

size *N* is an even number.

### A. Observation and Motivation

In DOUBLE, the traffic matrix $C(T)$ is decomposed as $C(T)=[T/N]\times Q+R$. For any $r_{ij}\in R$, if $r_{ij}>T/(2N)$, we call it an *LER* (large entry in *R*). Otherwise it is an *SER* (small entry in *R*). We have the following Lemma 1 (proved in Appendix A).

*Lemma 1:* In DOUBLE, if a particular line (row *i* or column *j*) of *R* contains *k* LERs, then in *Q* we have

$$\sum_{j=1}^{N}q_{ij}\le N-\left\lceil\frac{k}{2}\right\rceil \text{ for row } i, \text{ or } \sum_{i=1}^{N}q_{ij}\le N-\left\lceil\frac{k}{2}\right\rceil \text{ for column } j.$$

For example, the second row of *R* in Fig. 3 contains *k*=3 LERs (>T/(2N)=2). Then the entries in the second row of *Q* sum to at most $N-\lceil k/2 \rceil=2$.

Based on Lemma 1, we can *move* some packets from *R* to *Q*, while keeping the maximum line sum of *Q* not more than *N*. Note that each configuration in DOUBLE is equally weighted by $\phi_n=T/N$. Without loss of generality, if row *i* of *R* contains *k* LERs, we can move *at most* $\lceil k/2 \rceil$ of these LERs to *Q* by setting them to 0s in *R*, and at the same time increasing the corresponding entries in *Q* by one. Fig. 4 shows an example based on the *Q* and *R* in Fig. 3. We use *Q′* and *R′* to denote the modified *Q* and *R*. We can see that *Q′* can still be covered by *N* configurations with a weight $\phi_n=T/N$ each. So we can pack more packets in the *N* configurations used to cover the quotient matrix (than DOUBLE). On the other hand, since some LERs are moved to *Q*, it may not be necessary to use another *N* equally weighted configurations (with $\phi_n=T/N$) to cover *R′*.

The above observation motivates us to explore a more efficient scheduling algorithm than DOUBLE. In fact, there may be at most *N* LERs in each line of the original *R*. From Lemma 1, "half" of these LERs in each line of *R* can be moved to *Q*, while keeping the maximum line sum of *Q′* not more than *N*. Therefore, it is reasonable to expect that each line of *R′* contains at most *N/2* LERs after packet moving. As a result, we may be able to find *N/2* non-overlapping configurations, each weighted by $\phi_n=T/N$, to cover all the remaining LERs in *R′*. At the same time, we may find another set of *N/2* non-overlapping configurations, each with a *reduced* weight of $\phi_n=T/(2N)$, to cover the remaining SERs. If this can be done, then $S_{schedule}$ of DOUBLE can be reduced to

$$S_{schedule}=\frac{1}{T}\sum_{n=1}^{2N}\phi_n=\frac{1}{T}\times\left(\frac{T}{N}\times N+\frac{T}{N}\times\frac{N}{2}+\frac{T}{2N}\times\frac{N}{2}\right)=1.75 \quad (6)$$

### B. Issues

To achieve the goal mentioned above, we need to solve the following issues.

- Determine the set of LERs to be moved to *Q*, such that *each* line sum of *Q′* does not exceed *N*, and *R′* contains at most *N/2* LERs in *each* line.
- Among the *N* non-overlapping configurations used to cover *R′*, *N/2* of them should cover all the remaining LERs in *R′*, and the other *N/2* configurations should cover all the SERs not yet covered by the first *N/2* configurations.

Generally, it is not easy to determine the set of LERs that should be moved to *Q*. This can be seen from the example in Fig. 5. Assume all the non-zero entries (3s) in *R* are LERs. The number next to each line is the number of LERs that can be moved from this line to *Q*, which is obtained from $\lceil k/2 \rceil$ based on Lemma 1. If the four circled entries are moved, then we cannot further move any other LERs without violating the quota of the corresponding line. At this point, the last row still contains more than *N/2* LERs. For larger switch size *N*, it will be more difficult to figure out a proper set of LERs to move.

### C. Methodology

We first define two important notions, *PC*s (*Predefined Configurations*) and *DHH matrix*. For an *N×N* matrix, we can use *N* predefined non-overlapping configurations (or *PC*s) to cover all of its entries. As an example, eight non-overlapping *PC*s (i.e. *PC₁∼PC₈*), as defined in Fig. 6, can be used to cover all the entries of an 8×8 matrix. Note that the number at each entry of this matrix denotes the particular *PC* that covers this



Fig. 5. A carefully designed mechanism is necessary in order to move LERs properly.



Fig. 6. A possible set of predefined configurations for an 8×8 matrix.

entry. For easy reading, the entries covered by $PC_3$ are circled in Fig. 6.

Given an arbitrary 0/1 matrix, we can use two lines to partition it into four $(N/2)\times(N/2)$ zones/sub-matrices $A$, $B$, $C$ and $D$ as shown in Fig. 9 (and Fig. 10) in Appendix B. For each row/column, if the number of 1s in zone $A$ or $C$ is no less than that in zone $B$ or $D$, *or is less by at most one*, then this 0/1 matrix is called a *DHH matrix*. In other words, the two diagonal zones ($A$ and $C$) of a DHH matrix contain at least "half" number of 1*s for each row and column*. (Please refer to Appendix B for a more rigorous definition.)

Our approach in improving the scheduling efficiency (i.e. minimizing $S_{\text{schedule}}$) is based on the concepts of $PC$ and DHH matrix. We first convert the residue matrix $R=\{r_{ij}\}$ into a 0/1 indicator matrix $\Omega=\{\omega_{ij}\}$ such that $\omega_{ij}=1$ if $r_{ij}$ is an LER and $\omega_{ij}=0$ otherwise. The *DHH conjecture* given in Appendix B says that we can always find two permutation matrices $U$ and $V$, such that $\Omega'=U\Omega V$ is a DHH matrix. For the 8×8 matrix, assume that a DHH matrix $\Omega'$ is obtained by $\Omega'=U\Omega V$. Then, $PC_5\sim PC_8$ defined in Fig. 6 (which span over the sub-matrices $A$ and $C$ as illustrated in Fig. 9) can cover "more than half" of the 1s for *each row and column* of $\Omega'$. The remaining 1s covered by $PC_1\sim PC_4$ "correspond" to LERs that should be moved to $Q$. Since $\Omega'$ is obtained from $\Omega$ after some row/column permutations (i.e. $\Omega'=U\Omega V$), the 1s in $\Omega'$ do not directly match the original LER entries. Therefore, we need to invoke an inverse transform to get our desired configurations.

Fig. 7 gives an example, where the execution steps are indexed by the numbers in the dashed circles. In Step 1, we construct the indicator matrix $\Omega$ from $R$. In Step 2, we find two

$$T=32,\ N=8,\ \frac{T}{N}=4,\ N_S=16$$

$$C(32)=\begin{bmatrix}7&11&7&5&1&0&1&0\\7&7&11&1&0&0&2&1\\7&3&11&5&1&0&1&1\\9&0&0&5&8&4&3&3\\0&0&2&0&11&3&3&11\\0&4&0&0&3&7&15&3\\1&5&0&9&3&3&3&7\\1&0&1&1&3&15&3&3\end{bmatrix}=4Q+R=4\times\begin{bmatrix}1&2&1&1&0&0&0&0\\1&1&2&0&0&0&0&0\\1&0&2&1&0&0&0&0\\2&0&0&1&2&1&0&0\\0&0&0&0&2&0&0&2\\0&1&0&0&0&1&3&0\\0&1&0&2&0&0&0&1\\0&0&0&0&0&3&0&0\end{bmatrix}+\begin{bmatrix}3&3&3&1&1&0&1&0\\3&3&3&1&0&0&2&1\\3&3&3&1&1&0&1&1\\1&0&0&1&0&0&3&3\\0&0&2&0&3&3&3&3\\0&0&0&0&3&3&3&3\\1&1&0&1&3&3&3&3\\1&0&1&1&3&3&3&3\end{bmatrix}$$

$\xrightarrow{\ (1)\ }\ \Omega=\begin{bmatrix}1&1&1&0&0&0&0&0\\1&1&1&0&0&0&0&0\\1&1&1&0&0&0&0&0\\0&0&0&0&0&0&1&1\\0&0&0&0&1&1&1&1\\0&0&0&0&1&1&1&1\\0&0&0&0&1&1&1&1\\0&0&0&0&1&1&1&1\end{bmatrix}$

$(2)\downarrow$

$$\Omega'=U\Omega V=\begin{bmatrix}0&0&0&0&0&0&0&1\\0&0&0&0&1&0&0&0\\0&0&1&0&0&0&0&0\\0&0&0&1&0&0&0&0\\0&1&0&0&0&0&0&0\\0&0&0&0&0&1&0&0\\0&0&0&0&0&0&1&0\\1&0&0&0&0&0&0&0\end{bmatrix}\begin{bmatrix}1&1&1&0&0&0&0&0\\1&1&1&0&0&0&0&0\\1&1&1&0&0&0&0&0\\0&0&0&0&0&0&1&1\\0&0&0&0&1&1&1&1\\0&0&0&0&1&1&1&1\\0&0&0&0&1&1&1&1\\0&0&0&0&1&1&1&1\end{bmatrix}\begin{bmatrix}0&0&0&1&0&0&0&0\\0&1&0&0&0&0&0&0\\0&0&1&0&0&0&0&0\\0&0&0&0&0&0&0&1\\1&0&0&0&0&0&0&0\\0&0&0&0&0&1&0&0\\0&0&0&0&0&0&1&0\\0&0&1&0&0&0&0&0\end{bmatrix}=\left[\begin{array}{cccc|cccc}1&0&0&1&0&1&1&0\\1&0&0&1&0&1&1&0\\0&1&1&0&1&0&0&0\\0&0&0&1&0&0&1&0\\\hline0&1&1&0&1&0&0&0\\1&0&0&1&0&1&1&0\\1&0&0&1&0&1&1&0\\0&1&1&0&1&0&0&0\end{array}\right]$$

$(3)\downarrow$

$$\Delta_n=U^{-1}(PC_n)V^{-1}$$

$$A_1=\begin{bmatrix}0&0&0&0&1&0&0&0\\0&0&0&0&0&0&0&1\\0&0&0&0&0&1&0&0\\1&0&0&0&0&0&0&0\\0&0&0&0&0&1&0&0\\0&0&1&0&0&0&0&0\\0&1&0&0&0&0&0&0\\0&0&0&1&0&0&0&0\end{bmatrix}\ A_2=\begin{bmatrix}0&0&0&0&0&0&0&1\\0&0&1&0&0&0&0&0\\0&0&0&0&0&1&0&0\\0&0&0&0&1&0&0&0\\0&0&1&0&0&0&0&0\\0&1&0&0&0&0&0&0\\0&0&0&0&1&0&0&0\\1&0&0&0&0&0&0&0\end{bmatrix}\ A_3=\begin{bmatrix}0&0&1&0&0&0&0&0\\0&1&0&0&0&0&0&0\\0&0&0&1&0&0&0&0\\0&0&0&0&0&0&1&0\\1&0&0&0&0&0&0&0\\0&0&0&0&1&0&0&0\\0&0&0&0&0&0&0&1\\0&0&0&0&0&1&0&0\end{bmatrix}\ A_4=\begin{bmatrix}0&1&0&0&0&0&0&0\\0&0&0&0&1&0&0&0\\1&0&0&0&0&0&0&0\\0&0&0&1&0&0&0&0\\0&0&0&0&0&0&0&1\\0&0&1&0&0&0&0&0\\0&0&0&0&0&0&1&0\end{bmatrix}$$

$Q=\begin{bmatrix}1&2&1&1&0&0&0&0\\1&1&2&0&0&0&0&0\\1&0&2&1&0&0&0&0\\2&0&0&1&2&1&0&0\\0&0&0&0&2&0&0&2\\0&1&0&0&0&1&3&0\\0&1&0&2&0&0&0&1\\0&0&0&0&0&3&0&0\end{bmatrix}$ $(4)\rightarrow$ $R=\begin{bmatrix}3&\textcircled{3}&\textcircled{3}&1&1&0&1&0\\\textcircled{3}&\textcircled{3}&3&1&0&0&2&1\\\textcircled{3}&3&3&1&1&0&1&1\\1&0&0&1&0&0&\textcircled{3}&3\\0&0&2&0&3&\textcircled{3}&\textcircled{3}&3\\0&0&0&0&\textcircled{3}&3&3&\textcircled{3}\\1&1&0&1&\textcircled{3}&3&3&\textcircled{3}\\1&0&1&1&3&\textcircled{3}&\textcircled{3}&3\end{bmatrix}$

$$A_5=\begin{bmatrix}0&0&0&1&0&0&0&0\\1&0&0&0&0&0&0&0\\0&0&1&0&0&0&0&0\\0&0&0&0&0&0&0&1\\0&1&0&0&0&0&0&0\\0&0&0&0&0&1&0&0\\0&0&0&0&0&0&1&0\\0&0&0&1&0&0&0&0\end{bmatrix}\ A_6=\begin{bmatrix}1&0&0&0&0&0&0&0\\0&0&0&0&0&0&0&1\\0&1&0&0&0&0&0&0\\0&0&1&0&0&0&0&0\\0&0&0&1&0&0&0&0\\0&0&0&0&0&1&0&0\\0&0&0&1&0&0&0&0\\0&0&0&0&0&0&1&0\end{bmatrix}\ A_7=\begin{bmatrix}0&0&0&0&1&0&0&0\\0&0&0&0&0&0&1&0\\0&1&0&0&0&0&0&0\\0&0&0&0&0&0&1&0\\0&0&1&0&0&0&0&0\\1&0&0&0&0&0&0&0\\0&0&0&0&0&0&0&0\end{bmatrix}\ A_8=\begin{bmatrix}0&0&0&0&0&1&0&0\\0&0&0&0&1&0&0&0\\0&0&0&0&0&0&0&1\\0&0&1&0&0&0&0&0\\0&0&0&1&0&0&0&0\\0&0&0&0&0&1&0&0\\1&0&0&0&0&0&0&0\\0&0&1&0&0&0&0&0\end{bmatrix}$$

$Q'=\begin{bmatrix}1&3&2&1&0&0&0&0\\1&2&3&0&0&0&0&0\\2&0&2&1&0&0&0&0\\2&0&0&1&2&1&1&2\\0&0&0&2&1&1&2\\0&1&0&0&1&1&3&1\\0&1&0&2&1&0&0&2\\0&0&0&0&0&4&1&0\end{bmatrix}$ $R'=\begin{bmatrix}3&0&0&1&1&0&1&0\\3&0&0&1&0&0&2&1\\0&3&3&1&1&0&1&1\\1&0&0&1&0&0&0&3\\0&0&2&0&3&0&0&3\\0&0&0&0&0&3&3&0\\1&1&0&1&0&3&3&0\\1&0&1&1&3&0&0&3\end{bmatrix}$

$(5)\downarrow$

$$P_1=\ldots\quad P_2=\ldots\quad P_3=\ldots\quad P_4=\ldots\quad P_5=\ldots\quad P_6=\ldots\quad P_7=\ldots\quad P_8=\ldots$$

$(6)\downarrow$

$$C(32)\le 4\times Q'+R'\le 4\times\sum_{n=1}^{8}P_n+\left(2\times\sum_{n=1}^{4}A_n+4\times\sum_{n=5}^{8}A_n\right)\qquad S_{\text{schedule}}=\frac{1}{32}\times(4\times8+2\times4+4\times4)=1.75$$

Fig. 7. An illustrative example of our proposed approach.

184

permutations $U$ and $V$ to permute $\boldsymbol{\Omega}$ into a DHH matrix $\boldsymbol{\Omega'}$. Then, Step 3 imposes $U$ and $V$ on the predefined configurations $PC_1 \sim PC_8$ and uses $\Delta_n = U^{-1}(PC_n)V^{-1}$ (note $U^{-1}=U$ and $V^{-1}=V$ in this particular example) to generate eight configurations $\Delta_1 \sim \Delta_8$. For $\Delta_1 \sim \Delta_4$, if they cover some LERs in the original $R$, then these LERs (circled entries in Step 4) are to be moved to $Q$. Based on Lemma 1, we can set these LERs to 0s in $R$, and increase the corresponding entries in $Q$ by one. $Q'$ and $R'$ are thus obtained in Step 4. In Step 5, to cover $Q'$ we simply use the same edge-coloring algorithm [13] as in DOUBLE to determine configurations $P_1 \sim P_8$, each weighted by $T/N=4$. On the other hand, $R'$ can be covered by $\Delta_1 \sim \Delta_8$, with a weight of 4 for $\Delta_5 \sim \Delta_8$, and a reduced weight of $T/(2N)=2$ for $\Delta_1 \sim \Delta_4$. As a result, $S_{\text{schedule}}$ can be reduced to 1.75, as shown in Step 6. In fact, this 12.5% reduction in $S_{\text{schedule}}$ is independent of switch size $N$, as stipulated by (6). Note that there is a one-to-one mapping from entry to entry between $\boldsymbol{\Omega}$ and $\boldsymbol{\Omega'}$, which is determined by the linear transform $\boldsymbol{\Omega'}=U\boldsymbol{\Omega}V$. Therefore, the resulting $\Delta_1 \sim \Delta_8$ are non-overlapping, and they can cover every entry of $R'$.

## IV. Discussion

If we define LERs as $r_{ij}>T/[2\times(N_S-N)]$ (instead of $r_{ij}>T/(2N)$ for DOUBLE) and decompose $C(T)$ as discussed in Part B of Section II, then the speedup function (4) used in ADAPTIVE [8] can also be reduced to

$$S_{\text{schedule}} = \frac{1}{T}\sum_{n=1}^{N_S} \phi_n = 1 + \frac{3}{4} \times \frac{N}{N_S - N} . \qquad (7)$$

For $N_S<2N$, we can achieve a greater gain (than 12.5%) over the original ADAPTIVE algorithm. Since speedup and packet delay can trade one for another, this also means that packet delay can be smaller than that given in [8] if the same speedup is applied to the OPS switch.

From Fig. 11 in Appendix B, we know that $U$ and $V$ are used to record the row/column permutations involved, and they can be constructed from a *square unit matrix* $E$ (i.e. an $N\times N$ matrix with $N$ 1s at its diagonal entries and all other entries are 0s). Therefore, it is not necessary to get $U^{-1}$ and $V^{-1}$ by algebraic calculations. Instead, we can also start from a square unit matrix $E$, and permute its lines *in a reverse order* to generate $U^{-1}$ and $V^{-1}$.

Finally, it is important to note that our proposed approach is based on the DHH conjecture. We have tried to prove it for a very long time but without luck. We also cannot find a single counterexample by checking extensive samples using computer programs. Many mathematicians including the authors of [14] have reviewed our conjecture. As for now, the problem is still open.

## V. Conclusion

Due to the reconfiguration overhead, speedup and packet delay are two main issues for traffic scheduling in optical packet switches (OPS). In this paper, we proposed a new approach to improve the scheduling efficiency in OPS with guaranteed switching performance. Compared with the existing scheduling algorithms, our approach can significantly reduce speedup and/or packet delay. However, the proposed approach is based on the DHH conjecture given in this paper. We call for a proof or disproof for this conjecture.

### APPENDIX A   CORRECTNESS PROOF OF LEMMA 1

*Lemma 1:* In DOUBLE, if a particular line (row $i$ or column $j$) of $R$ contains $k$ LERs, then in $Q$ we have

$$\sum_{j=1}^{N} q_{ij} \le N - \left\lceil \frac{k}{2} \right\rceil \text{ for row } i, \text{ or } \sum_{i=1}^{N} q_{ij} \le N - \left\lceil \frac{k}{2} \right\rceil \text{ for column } j.$$

*Proof:* After $C(T)$ is divided by $T/N$, we have

$$C(T) = \frac{T}{N}Q + R \quad or \quad c_{ij} = \frac{T}{N}q_{ij} + r_{ij} , \qquad (8)$$

Without loss of generality, we assume that row $i$ of $R$ contains $k$ LERs. Because

$$\sum_{j=1}^{N} c_{ij} = \frac{T}{N}\sum_{j=1}^{N} q_{ij} + \sum_{j=1}^{N} r_{ij} \le T , \qquad (9)$$

we have

$$\sum_{j=1}^{N} q_{ij} \le \frac{T - \sum_{j=1}^{N} r_{ij}}{\frac{T}{N}} < \frac{T - \frac{T}{2N} \times k}{\frac{T}{N}} = N - \frac{k}{2} . \qquad (10)$$

Since $\sum_{j=1}^{N} q_{ij}$ is an integer, we then have

$$\sum_{j=1}^{N} q_{ij} \le N - \left\lceil \frac{k}{2} \right\rceil . \qquad (11)$$

### APPENDIX B   DHH CONJECTURE

In this appendix, we assume that the size $N$ of the matrices/vectors is an even number.

*Definition 1 (halve):* Given an arbitrary 0/1 vector (row or column whose entries are either 0 or 1), use a line to separate it into two equal parts as shown in Fig. 8. Let $x$ and $y$ denote the number of 1s in each part. If $|x-y|\le 1$, we say that the 1s in the vector are halved by the line.

Fig. 8 gives some examples, where the 1s are halved in (a) and (c), but not in (b) and (d).

*Definition 2 (DHH matrix):* Given an arbitrary $N\times N$ 0/1 matrix, use two lines to partition it into four $(N/2)\times(N/2)$ zones/sub-matrices $A$, $B$, $C$ and $D$ as in Fig. 9. For *each* row/column of the matrix, if the number of 1s in zone $A$ or $C$ is more than that in zone $B$ or $D$, or the 1s in this row/column are halved by one of the two lines, then this matrix is called a DHH matrix (it means that the diagonal half-size sub-matrices $A$ and $C$ contain at least "half" number of 1s *for each row and column*).

The matrix in Fig. 10a is a DHH matrix, whereas the matrix in Fig. 10b is not, because its second column has two more 1s in zone $D$ than that in zone $A$.

*DHH conjecture:* Given an *arbitrary $N\times N$ 0/1 matrix $\boldsymbol{\Omega}$*, we

$$[1\ 0\ 1\ 0\,|\,0\ 1\ 1\ 0]\qquad [0\ 0\ 0\ 0\,|\,0\ 1\ 1\ 0]$$

(a)                  (b)

$$\begin{array}{c}\begin{bmatrix}1\\1\\0\\0\\\hline1\\0\\1\\1\end{bmatrix}\end{array}\qquad\begin{bmatrix}1\\1\\0\\0\\\hline1\\1\\1\\1\end{bmatrix}$$

(c)        (d)

Fig. 8. Illustration of definition "halve".

$$\left[\begin{array}{c|c} A & B \\ \hline D & C \end{array}\right]$$

Fig. 9. Sub-matrices.

$$\left[\begin{array}{cc|cc}1&0&0&0\\1&1&1&1\\\hline0&1&0&1\\0&1&0&0\end{array}\right]\qquad\left[\begin{array}{cc|cc}1&0&0&0\\1&0&1&1\\\hline0&1&0&1\\0&1&0&0\end{array}\right]$$

(a)              (b)

Fig. 10. (a) is a DHH matrix, and (b) is not.



A "×" is used to indicate a row or a column that violates the requirement defined in the conjecture.

The final result is a DHH matrix.

$U$ and $V$ are used to record the row/column permutations. They can be constructed as follows: Initialize $U$ and $V$ as square unit matrices $E$. If we permute two rows in $\Omega$, then we also permute the corresponding rows in $U$; if we permute two columns in $\Omega$, then we permute the corresponding columns in $V$ too. After $\Omega$ is turned into a DHH matrix, the corresponding $U$ and $V$ are also obtained.
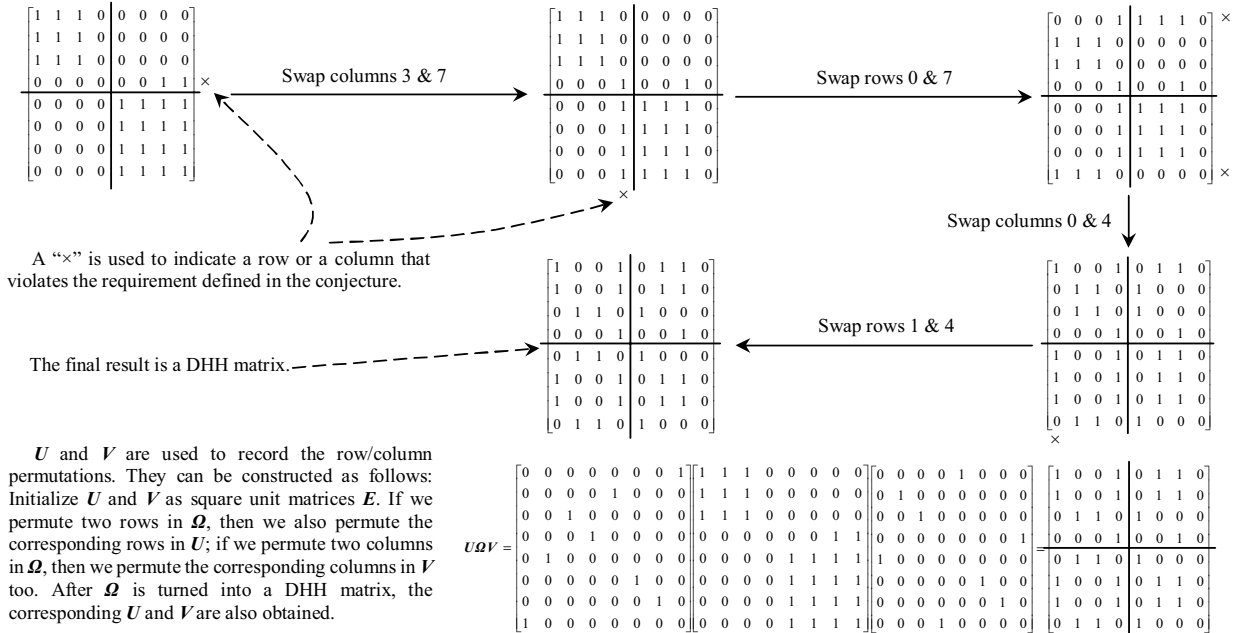
Fig. 11. Permute a matrix $\Omega$ into a DHH matrix (the rows/columns are numbered from 0 to 7).

can permute its rows or columns[2] for a limited number of times, such that $\Omega$ can be turned into a DHH matrix. In other words, there exist two permutation matrices $U$ and $V$, such that $U\Omega V$ is a DHH matrix.

For example, if we swap the first row and the last row in Fig. 10b, then the resulting matrix is a DHH matrix. Fig. 11 gives a more complex example for $\Omega$ in Fig. 7.

### REFERENCES

[1] J.E Fouquet et. al, "A compact, scalable cross-connect switch using total internal reflection due to thermally-generated bubbles", *IEEE LEOS Annual Meeting*, pp. 169-170, Dec. 1998.

[2] L. Y. Lin, "Micromachined free-space matrix switches with submilli-second switching time for large-scale optical crossconnect", *OFC'98 Tech. Digest*, pp. 147-148, Feb. 1998.

[3] O. B. Spahn, C. Sullivan, J. Burkhart, C. Tigges, and E. Garcia, "GaAs-based microelectromechanical waveguide switch", *Proc. 2000 IEEE/LEOS Intl. Conf. on Optical MEMS*, pp. 41-42, Aug. 2000.

[4] A. J. Agranat, "Electroholographic wavelength selective crossconnect", *1999 Digest of the LEOS Summer Topical Meetings*, pp. 61-62, Jul. 1999.

[5] K. Kar, D. Stiliadis, T. V. Lakshman, and L. Tassiulas, "Scheduling algorithms for optical packet fabrics", *IEEE Journal on Selected Areas in Communications*, vol. 21, issue 7, pp. 1143-1155, Sept. 2003.

[6] B. Towles and W. J. Dally, "Guaranteed scheduling for switches with configuration overhead", *IEEE/ACM Trans. Networking*, vol. 11, no. 5, pp. 835-847, Oct. 2003.

[7] Xin Li and Hamdi, M., "On scheduling optical packet switches with reconfiguration delay", *IEEE Journal on Selected Areas in Communications*, vol. 21, issue 7, pp. 1156-1164, Sept. 2003.

[8] Bin Wu and Kwan L. Yeung, "Minimizing internal speedup for performance guaranteed optical packet switches", *IEEE GLOBECOM '04*, vol. 3, pp. 1742-1746, Dec. 2004.

[9] Bin Wu and Kwan L. Yeung, "Scheduling optical packet switches with minimum number of configurations", *IEEE ICC '05*, vol. 3, pp. 1830-1835, May 2005.

[10] Bin Wu and Kwan L. Yeung, "Traffic scheduling in non-blocking optical packet switches with minimum delay", *IEEE GLOBECOM '05*, vol. 4, pp. 2041-2045, Dec. 2005.

[11] Bin Wu, Kwan L. Yeung and V. O. K. Li, "Two-layer parallel switching: A practical and survivable design for performance guaranteed optical packet switches", *IEEE GLOBECOM '05*, vol. 4, pp. 1905-1909, Dec. 2005.

[12] Bin Wu and Kwan L. Yeung, "On optimization of optical packet switches with reconfiguration overhead", *IEEE HPSR '05*, pp. 217-221, May 2005.

[13] R. Cole and J. Hopcroft, "On edge coloring bipartite graphs", *SIAM Journal on Computing*, vol. 11, pp. 540-546, Aug. 1982.

[14] R. A. Brualdi and H. J. Ryser, "Combinatorial matrix theory", Cambridge University Press, 1991.

[2] i.e., swap its rows or columns. Note that a row can only be swapped with another row, and a column can only be swapped with another column.