

RZ 3603 (# 99613) 04/28/05  
Computer Science 26 pages

# Research Report

## Performance Evaluation of the Data Vortex Photonic Switch

Ilias Iliadis<sup>1</sup>, Nikolaos Chrysos<sup>2</sup>, Cyriel Minkenberg<sup>1</sup>

<sup>1</sup>IBM Research GmbH  
Zurich Research Laboratory  
8803 Rüschlikon  
Switzerland

<sup>2</sup>Institute of Computer Science  
Foundation for Research and Technology - Hellas (FORTH) ICS-FORTH  
P.O. Box 1385  
Vassilika Vouton  
Heraklion, Crete  
GR-711-10 Greece

### LIMITED DISTRIBUTION NOTICE

This report has been submitted for publication outside of IBM and will probably be copyrighted if accepted for publication. It has been issued as a Research Report for early dissemination of its contents. In view of the transfer of copyright to the outside publisher, its distribution outside of IBM prior to publication should be limited to peer communications and specific requests. After outside publication, requests should be filled only by reprints or legally obtained copies of the article (e.g., payment of royalties). Some reports are available at <http://domino.watson.ibm.com/library/Cyberdig.nsf/home>.

**IBM** Research  
Almaden · Austin · Beijing · Delhi · Haifa · T.J. Watson · Tokyo · Zurich

# Performance Evaluation of the Data Vortex Photonic Switch

Ilias Iliadis, Nikolaos Chrysos, Cyriel Minckenberg

## Abstract

The Data Vortex photonic packet-switching architecture features an all-optical transparent data path, highly distributed control, low latency, and a high degree of scalability. These characteristics make it attractive as a routing fabric in future photonic packet switches. We analyze the performance of the Data Vortex architecture as a function of its height and angle dimensions,  $H$  and  $A$ . The investigation is based on two performance measures: the average delay and the maximum throughput of the switch. We present an analytical model assuming uniform traffic and derive closed-form expressions for these measures. Our results obtained demonstrate that as  $H$  increases, the saturation throughput decreases and approaches  $\frac{2}{9} = 0.22$  when  $A$  is small and  $H$  is large. Furthermore, for fixed switch size, the saturation throughput is maximized when  $A$  is minimal. We also present simulation results for the maximum throughput under uniform and nonuniform traffic, as well as for the mean number of hops and the mean input-queue packet delay as a function of input load, and address the issue of resequencing delay. The results obtained advocate that to support more ports, it is preferable to increase the height dimension and to keep the angle dimension as small as possible.

**Index Terms**—Photonic switching systems, modeling, optical communication.

## I. INTRODUCTION

Massively parallel high-performance computing systems require large-scale, high-bandwidth, low-latency packet-switched interconnection networks. Existing electronic solutions are limited by the problems of power and pin-out density. Optical technologies appear well positioned to overcome these issues. However, challenges such as contention resolution, buffering, scalability, and latency must be addressed to obtain a successful solution. The Data Vortex switch is a photonic packet-switch architecture tailored for optical computer interconnection networks [1, 2] that addresses these issues. Data Vortex comprises an all-optical fabric datapath and electronic units for the control operation as well as for the input and output interfaces to the fabric. The fabric is self-routing, i.e. the routing and contention control units are physically distributed along the nodes. Thus simple and small electronic control circuits can be used, resulting in inexpensive nodes. Moreover, owing to its deflection routing policy, the fabric is *bufferless*, enabling distributed contention resolution and a transparent, low-latency, end-to-end optical datapath. This architecture can support an arbitrarily large number of computing nodes using only simple  $2 \times 2$  switching nodes as a building block.

The performance of the Data Vortex architecture was studied by means of simulation in [3–5]. Here we develop an analytical model and provide a theoretical explanation of the performance behavior of this architecture. In particular, this model allows us to derive closed-form expressions for the maximum throughput

I. Iliadis and C. Minckenberg are with IBM Research, Zurich Research Laboratory, Säumerstrasse 4, CH-8803 Rüschlikon, Switzerland.

N. Chrysos is with the Institute of Computer Science, Foundation for Research and Technology - Hellas (FORTH) ICS-FORTH, P.O. Box 1385, Vassilika Vouton, Heraklion, Crete, GR-711-10 Greece.

and the mean packet latency as well as gain more insight into the properties of several other performance-related measures. In addition, we perform simulations that take into account the effects of angle decoding, unbalanced traffic, input queuing, and out-of-order delivery, which have not been considered in previous work.

We describe the architecture of the Data Vortex, including its topology and the routing and deflection policies in Section II. Our main contribution is the development of an analytical model for studying the behavior of this architecture. In Section III we derive closed-form expressions for the deflection probabilities and the mean delay through the switch for any  $A$  and  $H$  under uniform uncorrelated arrivals. In Section IV we derive the saturation throughput of the switch and investigate how to select the parameters in order to maximize it. Finally, the performance of the Data Vortex architecture is discussed in Section V using analytical and simulation results. Simulation results for various performance measures such as the throughput under unbalanced traffic, the hop count distribution, the mean packet latency and the resequencing delay are presented.

## II. DATA VORTEX ARCHITECTURE

### A. Topology

We briefly review the Data Vortex architecture presented in [1, 2]. The Data Vortex switch fabric contains a set of nodes arranged in  $C$  concentric cylinders, where each cylinder level  $c$  (hereafter referred to as  $CL_c$  with  $0 \leq c \leq C - 1$ ) is considered as a routing level. Packets enter the switch in  $CL_0$  (outermost), and depart from  $CL_{C-1}$  (innermost). The topology and the size of this switch are characterized by the angle dimension  $A$  and the height dimension  $H$ . The maximum number of input and output ports  $N$  supported by the switch is  $N = A \cdot H$ , where  $H$  specifies the number of nodes along the height of each of the concentric cylinders and  $A$  the number of nodes along their circumference. In total, each cylinder contains  $N$  nodes, and there are  $C = \log_2(H) + 1$  cylinders. The total number of nodes in the switch is  $A \cdot H \cdot C$ . The  $N$  input interfaces of the fabric are connected to the outermost  $CL_0$  and the  $N$  output interfaces to the innermost  $CL_{C-1}$ . The remaining cylinders constitute a bufferless routing and deflection network. Although the routing fabric is unbuffered, the input and output interfaces of the switch contain electronic buffer memories. Packets may have to wait inside these memories until they can enter the switch fabric (input buffering), or until their final destination is ready to accept them (output buffering).

Every node in the routing network is a  $2 \times 2$  switching element having two input and two output links, connecting the node to four different neighbor nodes. Each node has an incoming link from a node in the same cylinder (WEST), and one from a node in the outer adjacent cylinder (NORTH). It also has an outgoing link to a node in the same cylinder (EAST) and one to a node in the inner adjacent cylinder (SOUTH). Each switch input interface connects to the NORTH input of a distinct node in the outermost  $CL_0$ , i.e., there is a one-to-one mapping between the outermost cylinder nodes and the input interfaces. Similarly, the SOUTH output of a node in the innermost  $CL_{C-1}$  directs packets to a specific output interface. Contention between

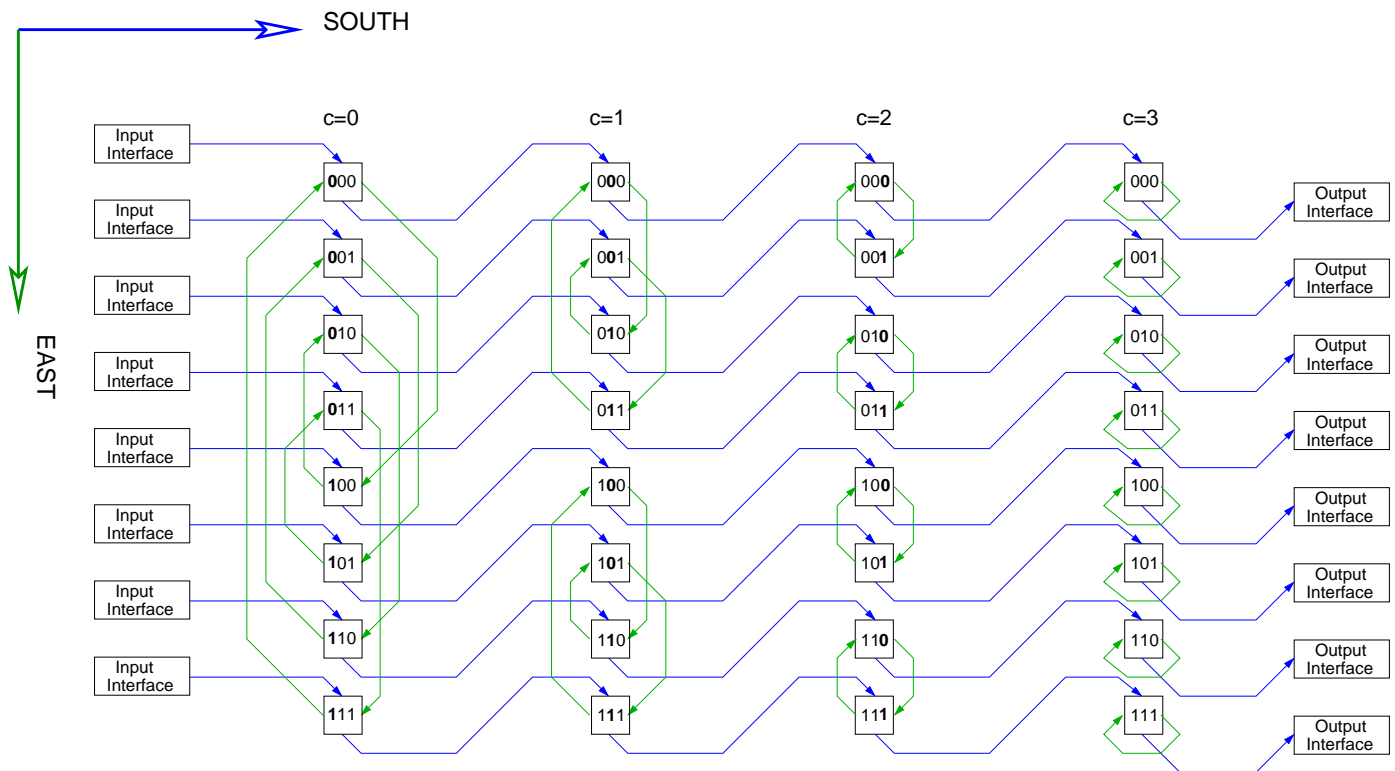


Fig. 1. The Data Vortex topology with  $H=8$  ( $C = 4$ ) and  $A=1$ .

the ports at every node is resolved by a deflection technique as described in Section II-B which results in the routing method as described in Section II-C.

The switch-internal header of a packet contains a height-destination address and an angle-destination address. All nodes in the same cylinder, except the innermost one, decode a specific bit, which is different for each cylinder, out of a packet's height-destination address. As a result, when packets reach  $CL_{C-1}$  after passing through  $\log_2(H)$  cylinder levels, they are located at the desired height. The innermost cylinder nodes decode the destination angle, thus ensuring that packets are sent to the proper height *and* angle destination.

The exact interconnection patterns for any given  $H$  and  $A$  are defined in [5, Sec. 2.2]. Fig. 1 illustrates the configuration with  $H = 8$  and  $A = 1$ , where the binary numbers in the boxes represent the height address. The interconnection pattern exhibits the following important characteristics. First, the height never changes when moving SOUTH. Second, at any  $CL_c$ , the positions of the  $c$  most significant bits representing the node height remain unchanged for any EAST connection. This implies that every cylinder fixes one additional height bit (the  $c+1$ -th bit shown in bold), until the packet has arrived at its destination height. No height change is possible in the innermost  $CL_{C-1}$  because packets reaching this cylinder are already located at the desired height. This cylinder is intended exclusively to move packets to the correct destination angle. Third, every connection (EAST as well as SOUTH) results in a change of angle by one modulo  $N$ . This feature is shown in Fig. 2, but it is not apparent in Fig. 1 because  $A = 1$ .

## B. Deflection

Let us consider a pair of nodes, say  $n_1$  and  $n_2$ , that might send a packet to a common destination node, say  $n$ . Suppose  $n_1$  is connected to the NORTH input of node  $n$  (outer cylinder), and  $n_2$  is connected to its WEST input (same cylinder), see Fig. 2. To resolve a potential conflict, there is a control connection from  $n_2$  to  $n_1$  over which an appropriate control signal, either of type *permit* or *no-permit* is sent. Node  $n_2$  has priority over  $n_1$ , i.e.  $n_2$  can always send a packet EAST to node  $n$ . When it does so, it also sends a *no-permit* control signal to  $n_1$  to halt any potential transmission from  $n_1$  to  $n$ . Upon reception of this signal, the halted node  $n_1$  cannot send a packet SOUTH to  $n$ ; it can only send it EAST and also issue a corresponding *no-permit* signal to its control connection. This mechanism is referred to as *deflection*. It also implies that the WEST input always preempts the NORTH input of a node, i.e., packets traveling along a given cylinder have priority over packets from the outer adjacent cylinder. This deflection method ensures that at any time, at most one of a node's inputs receives a new packet, and accordingly, only one of its outputs sends one. This property, named *single-packet-routing*, significantly simplifies the routing and the flow-control circuits required in the Data Vortex architecture [2]. As we will see in Section II-C, deflections may also occur under *permit* signals for the purpose of routing.

The deflection mechanism also applies at the switch fabric interfaces, as each input interface module has a SOUTH output to the NORTH input of some node  $n$  in  $CL_0$ . The node WEST of  $n$  can deflect a packet from the input interface  $n$ , preventing it from entering the switch fabric. At the output side, the output interface that can be reached through an innermost node  $n$  can exert backpressure towards  $n$  to deflect a packet at the output interface. The deflected packet will then circle around in  $CL_{C-1}$ , following EAST paths, waiting for a new chance to enter the output interface. This backpressure can be used, for instance, if the output interface buffer is full.

## C. Routing

Time in Data Vortex is divided into equal-length time slots; the switch operates in a synchronous fashion on fixed-length packets. Henceforth, we select the time slot length as a unit of time. We assume that both EAST and SOUTH movements occur at slot boundaries. We identify a node by its three coordinates  $(a, c, h)$ , where  $a \in \{0, \dots, A - 1\}$ ,  $h \in \{0, \dots, H - 1\}$  and  $c \in \{0, \dots, C - 1\}$ . The routing algorithm executed by the nodes is independent of the packet's incoming direction. At the beginning of a time slot, a node, say  $(a, c, h)$ , accepts a packet (if any), determines the routing to be applied, and forwards the packet in the resulting direction. Suppose that the destination of a packet arriving at node  $(a_1, c_1, h_1)$ , with  $c < C - 1$ , is  $(a, h)$ . The node is eligible to send the packet SOUTH if its  $c_1 + 1$ -th most significant height bit is equal to the corresponding bit in  $h$ . Otherwise, the packet is deflected and sent EAST, resulting in an angle increase as discussed in Section II-A. Note that the interconnection pattern within a cylinder guarantees that the  $c_1 + 1$ -th most significant height bit of the EAST neighbor will now be equal to the corresponding bit in  $h$ , because the

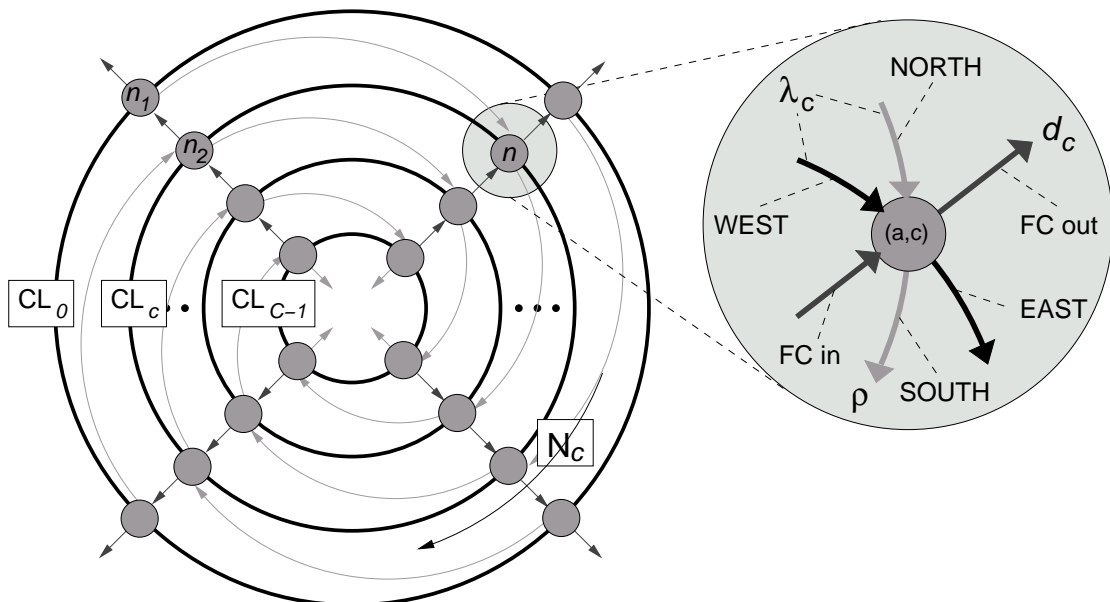


Fig. 2. A schematic to illustrate the throughput derivation procedure.

topology guarantees that the relevant height bit alternates with every angle increase in every cylinder except  $C - 1$ . Hence, in general, any packet has an opportunity to arrive at an eligible node and move SOUTH every other angle increase, i.e., after at most one hop within a cylinder level. An eligible node will in fact send the packet SOUTH if the flow control from the competing node allows it, i.e., in the presence of a *permit* signal, otherwise it will send it EAST owing to deflection as discussed in Section II-B.

A packet in the innermost  $CL_{C-1}$  is guaranteed to be at a node  $(a_1, C - 1, h)$  located at the right height. Then the routing will compare  $a_1$  with  $a$ : if these match, and the flow-control signal allows it, it will send the packet SOUTH towards the output interface, otherwise it will send it EAST. Note that a packet reaches its destination angle after at most  $A - 1$  EAST hops in the innermost  $CL_{C-1}$ .

### III. SYSTEM ANALYSIS

Here we derive the mean delay of a packet through the system. We assume that traffic is uniform, i.e. the destination of an arbitrary packet can be any of the  $N$  output ports with equal probability  $1/N$ . We also assume Bernoulli i.i.d. arrivals, with  $\rho$  denoting the probability that a packet arrives at a given slot of an input interface. The maximum or saturation throughput for given  $A$  and  $C$  is denoted by  $\rho_{A,C}^*$ .

#### A. Deflection probabilities

First, we derive expressions for the deflection probabilities and study their behavior. Under a stable operation (i.e. for  $\rho < \rho_{A,C}^*$ ),  $\rho$  is also the rate at which packets enter from NORTH or depart to SOUTH at any node in any  $CL_c$ ,  $c \in \{0, \dots, C - 1\}$ , see Fig. 2. Packets may have to follow EAST hops, either for routing purposes or because they are deflected by a packet in the inner adjacent cylinder. Whenever a packet goes EAST, it sends a *no-permit* signal, thus possibly deflecting another packet in the outer adjacent cylinder.

For any time slot, we denote by  $\lambda_c$  the probability that a node in  $\text{CL}_c$  has an active input, and by  $d_c$  the probability that it sends a *no-permit* signal due to sending a packet EAST. Therefore,  $d_c$  is also equal to the probability that, in a given time slot, a node in  $\text{CL}_{c-1}$  is prohibited from sending SOUTH, regardless of whether it has a packet to send. We denote by  $N_c$  the average number of EAST hops a packet traverses within  $\text{CL}_c$ . Thus, on average, a packet follows the EAST path  $N_c$  times and the SOUTH path only once. As we assume uniform traffic, it follows that an active node in  $\text{CL}_c$  sends its packet EAST with probability  $\frac{N_c}{N_c+1}$  and SOUTH with probability  $\frac{1}{N_c+1}$ . Now, for each  $\text{CL}_c$ , we have  $\rho = \lambda_c \cdot \frac{1}{N_c+1}$ ,  $d_c = \lambda_c \cdot \frac{N_c}{N_c+1}$ , which imply that

$$d_c = \rho \cdot N_c, \quad \forall c \in \{0, \dots, C-1\}. \quad (1)$$

We now calculate the mean number of EAST hops  $N_{c-1}$  in the outer adjacent cylinder based on  $d_c$ . Owing to the uniform traffic assumption and the routing mechanism described in Section II-C, a packet enters  $\text{CL}_{c-1}$  at a node eligible to send it directly SOUTH with probability  $\frac{1}{2}$ , whereas with probability  $\frac{1}{2}$  it has to follow an EAST hop to reach such an eligible node. If a packet is prevented from moving SOUTH (deflected), it will have to follow two EAST hops in  $\text{CL}_{c-1}$  before it again reaches an eligible node. Assuming that each time a packet arrives at an eligible node the deflection probability is the same (equal to  $d_c$ ) independently of the number of hops traversed, the average number of hops a packet will traverse in  $\text{CL}_{c-1}$  is given by

$$N_{c-1} = \frac{1}{2} \sum_{k=0}^{\infty} 2k(1-d_c)d_c^k + \frac{1}{2} \sum_{k=0}^{\infty} (2k+1)(1-d_c)d_c^k = \frac{1}{2}(1-d_c) \left( 4 \sum_{k=0}^{\infty} kd_c^k + \sum_{k=0}^{\infty} d_c^k \right). \quad (2)$$

Note that  $d_c = \lambda_c \frac{N_c}{N_c+1} < \lambda_c \leq 1$ , i.e.  $d_c < 1$ . Consequently, we have that  $\sum_{k=0}^{\infty} d_c^k = \frac{1}{1-d_c}$  and  $\sum_{k=0}^{\infty} kd_c^k = \frac{d_c}{(1-d_c)^2}$ , hence (2) reduces to

$$N_{c-1} = \frac{3d_c + 1}{2(1-d_c)}, \quad \forall c \in \{1, \dots, C-1\}. \quad (3)$$

Finally, combining (3) with (1), we obtain the following recursive relation:

$$d_{c-1} = \rho \cdot \frac{3d_c + 1}{2(1-d_c)}, \quad \forall c \in \{1, \dots, C-1\}. \quad (4)$$

According to the discussion in Section II-C and assuming that no backpressure is exerted by the output interfaces, the number of hops that packets will move EAST in the innermost  $\text{CL}_{C-1}$  to reach their destination angle is uniformly distributed in the interval  $[0, A-1]$ . Hence,  $N_{C-1} = \frac{A-1}{2}$ , which by virtue of (1) yields

$$d_{C-1} = \rho \cdot \frac{A-1}{2}. \quad (5)$$

Closed-form expressions for the deflection probabilities  $d_c^{(C)}$  as a function of  $C$ ,  $c$ ,  $A$ , and  $\rho$  are obtained through the deflection base-functions  $g_c$  with the following proposition.

*Proposition 1:* It holds that

$$d_c^{(C)}(\rho) = \begin{cases} g_{C-1-c}(\rho), & \rho \leq \rho_{A,C}^*, \\ g_{C-1-c}(\rho_{A,C}^*), & \rho \geq \rho_{A,C}^*, \end{cases} \quad \forall c \in \{0, \dots, C-1\}, \quad (6)$$

where

$$g_c(\rho) = \begin{cases} \frac{a_1 z_1^c + a_2 z_2^c}{b_1 z_1^c + b_2 z_2^c}, & \rho \neq \frac{2}{9}, \\ \frac{1}{3} \left[ 1 - \frac{2(4-A)}{(4-A)c+6} \right], & \rho = \frac{2}{9}, \end{cases} \quad \forall c \in \mathbb{N}_0, \quad \rho \in \mathbb{R}, \quad (7)$$

$$z_{1,2} = \frac{2 + 3\rho \pm \sqrt{\Delta}}{16\rho}, \quad \text{with } \Delta = 9\rho^2 - 20\rho + 4 = 9\left(\rho - \frac{2}{9}\right)(\rho - 2), \quad (8)$$

and

$$a_i = \frac{4(A-1)\rho z_i - (A-2)}{8(z_i - z_{3-i})}, \quad b_i = \frac{8z_i - (A+2)}{8(z_i - z_{3-i})} = \frac{2z_{3-i}}{2z_{3-i} - 1} a_i, \quad \text{for } i = 1, 2. \quad (9)$$

*Proof:* See Appendix A. ■

From (6) we deduce that the properties of the deflection probabilities  $d_c^{(C)}$  depend on those of the deflection base-functions  $g_c$  as stated by the following lemmas. Let  $\mathcal{R}_c$  be the continuous interval, either open  $[0, u_c)$  or closed  $[0, u_c]$  (with  $u_c \leq 1$ ), in which  $g_c(\rho) < 1, \forall \rho \in \mathcal{R}_c$ .

LEMMA 1: The  $g$ -functions are increasing in  $\rho$ , i.e.  $g_c(\rho_1) \leq g_c(\rho_2), \forall A \in \mathbb{N}, c \in \mathbb{N}_0, \rho_1 < \rho_2 \in \mathcal{R}_c$ , with the equality holding iff  $A = 1$  and  $c = 0$ .

*Proof:* See Appendix A. ■

LEMMA 2: It holds that

$$\left\{ \begin{array}{l} g_{c-1}(\rho) < g_c(\rho), \quad \text{for } A = 1, 2, \\ g_{c-1}(\rho) \geq g_c(\rho) \Leftrightarrow \rho \leq \hat{\rho}_A, \quad \text{for } A \geq 3, \end{array} \right\} \quad \forall c \in \mathbb{N}, \quad \rho \in \mathcal{R}_c, \quad (10)$$

with

$$\hat{\rho}_A \triangleq \frac{2(A-2)}{(A-1)(A+2)}, \quad \text{for } A \geq 3. \quad (11)$$

*Proof:* See Appendix A. ■

LEMMA 3: The  $g$ -functions are increasing in  $A$ , i.e.  $g_c(\rho; A) < g_c(\rho; A+1), \forall A \in \mathbb{N}, c \in \mathbb{N}_0, \rho \in \mathbb{R}$ .

*Proof:* By using induction on  $A$ . It is similar to that of Lemma 1 and it is therefore omitted. ■

LEMMA 4: For all  $c, c \in \mathbb{N}_0$ , it holds that

$$g_{c+1}(\rho; A=1) = g_c(\rho; A=2), \quad \forall \rho \in [0, 1]; \quad g_{c+1}(\rho; A=2) < g_c(\rho; A=4) < 1, \quad \forall \rho \in \mathcal{R}_c^{(A=4)}. \quad (12)$$

*Proof:* See Appendix A. ■

LEMMA 5: Let us define

$$g_\infty(\rho) \triangleq \lim_{c \rightarrow \infty} g_c(\rho), \quad \text{for } \rho \in \mathcal{R}_\infty, \quad \text{where } \mathcal{R}_\infty \triangleq \lim_{c \rightarrow \infty} \mathcal{R}_c. \quad (13)$$

Then it holds that

$$g_\infty(\rho) = \begin{cases} y_2(\rho), & \text{for } \rho \in [0, \frac{2}{9}], \quad 1 \leq A \leq 4 \\ \left\{ \begin{array}{l} y_2(\rho), \quad \text{for } \rho \in [0, \hat{\rho}_A), \\ \frac{A-2}{A+2}, \quad \text{for } \rho = \hat{\rho}_A, \end{array} \right\}, & A \geq 5, \end{cases} \quad (14)$$

where

$$y_2(\rho) \triangleq \frac{2 - 3\rho - \sqrt{\Delta}}{4} = \frac{a_1(\rho)}{b_1(\rho)} = \frac{2z_2 - 1}{2z_2} = 1 - 4\rho z_1, \quad \text{for } \rho \in \left[0, \frac{2}{9}\right], \quad (15)$$



$$\text{and } \mathcal{R}_\infty = [0, u_\infty], \quad \text{with } u_\infty = \lim_{c \rightarrow \infty} u_c = \begin{cases} \frac{2}{9}, & \text{for } 1 \leq A \leq 4, \\ \hat{\rho}_A, & \text{for } A \geq 4. \end{cases} \quad (16)$$

*Proof:* See Appendix A. ■

The behavior of the deflection base-functions  $g_c(\rho)$  for  $A \leq 6$  is depicted in Fig. 3. As  $c$  increases, they approach the  $y_2(\rho)$  function for  $\rho \leq u_\infty$  and exhibit a sharp increase at  $\rho = u_\infty^+$  for  $A \leq 4$  or at  $\rho = u_\infty^-$  for  $A \geq 5$ . This is due to the fact that  $\lim_{c \rightarrow \infty} g'_c(u_\infty) = \infty$  because either  $g'_\infty(u_\infty) = y'_2(\frac{2}{9}) = \infty$  for  $A \leq 4$ , or  $g_\infty(\rho)$  is discontinuous at  $\rho = u_\infty$ , given that  $g_\infty(u_\infty) = \frac{A-2}{A+2} > \frac{1}{A-1} = y_2(u_\infty) = \lim_{\rho \rightarrow u_\infty} g_\infty(\rho)$  for  $A \geq 5$ .

The properties of the deflection probabilities are obtained from Lemmas 1, 2 and 3, by using (6) and (5), and are given by the following corollaries.

*Corollary 1:* The deflection probabilities  $d_c$  are increasing functions in  $\rho$ , i.e.  $d_c^{(C)}(\rho_1) \leq d_c^{(C)}(\rho_2)$ ,  $\forall A \in \mathbb{N}$ ,  $C \in \mathbb{N}$ ,  $c \in \{0, 1, \dots, C-1\}$ ,  $\rho_1 < \rho_2 \leq \rho_{A,C}^*$ , with the equality holding iff  $A = 1$  and  $c = C-1$ .

*Corollary 2:* It holds that

$$\left\{ \begin{array}{ll} d_0^{(C)}(\rho) > d_1^{(C)}(\rho) > \dots > d_{C-1}^{(C)}(\rho) = 0, & \text{for } A = 1, \\ d_0^{(C)}(\rho) > d_1^{(C)}(\rho) > \dots > d_{C-1}^{(C)}(\rho) = \frac{\rho}{2}, & \text{for } A = 2, \\ d_0^{(C)}(\rho) \underset{\geq}{\leq} d_1^{(C)}(\rho) \underset{\geq}{\leq} \dots \underset{\geq}{\leq} d_{C-1}^{(C)}(\rho) = \frac{A-1}{2}\rho \Leftrightarrow \rho \underset{\geq}{\leq} \hat{\rho}_A, & \text{for } A \geq 3, \end{array} \right\} \quad \forall C \in \mathbb{N}, \quad 0 < \rho \leq 1. \quad (17)$$

*Corollary 3:* The deflection probabilities  $d_c$  are increasing functions in  $A$ , i.e.  $d_c^{(C)}(\rho; A) < d_c^{(C)}(\rho; A+1)$ ,  $\forall A \in \mathbb{N}$ ,  $C \in \mathbb{N}$ ,  $c \in \{0, 1, \dots, C-1\}$ ,  $\rho \in [0, 1]$ .

These properties are illustrated in Fig. 4, which shows the deflection probabilities  $d_c^{(5)}$  corresponding to five cylinders ( $c = 0, 1, 2, 3, 4$ ) for  $A \leq 6$ .

**Remark 1.** At saturation, the deflection probability at  $CL_0$  is, as expected, larger than those at the remaining cylinders (this is proven later in Corollary 4). Note however, that for  $A \geq 5$ , this deflection probability is the smallest, compared with those at the remaining cylinders, practically in the entire range of (stable) loads. Only for loads close to the saturation throughput it does increase sharply and become the largest.

## B. Delay analysis

The packet delay  $D$  comprises the delay  $D_i$  in the input buffers and the delay  $D_f$  through the fabric. The in-fabric delay  $D_f$  a packet experiences by propagating along the links through the fabric is equal to the number of hops it traverses to reach its output interface, starting from the input interface. This is an important measure in the Data Vortex, because no signal reshaping or regeneration is performed inside the fabric. As the number of hops traversed by a packet inside the switch increases, the signal-to-noise ratio of the packet drops, leading to an increase in the bit-error rate, which, depending on the quality of the optical components, will fall below an acceptable level after a certain number of hops. Note that  $D_f$  is at least  $C$  time slots. Also, as the input interface performs cut-through, the input delay  $D_i$  can be as small as zero. Hence,  $D \geq C$ .

Next, we derive the mean delay  $\bar{D}$  of a packet through the system, where  $\bar{D} = \bar{D}_i + \bar{D}_f$ . We consider an input interface and assume that the type of the control signal received at an arbitrary time slot is independent

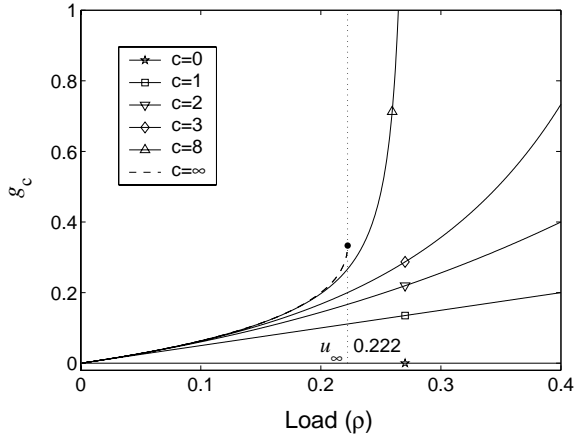
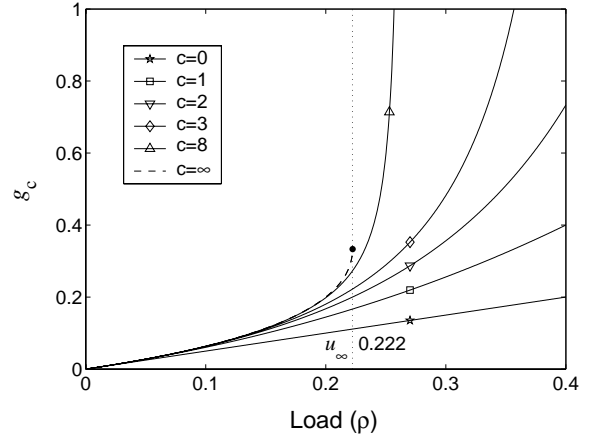
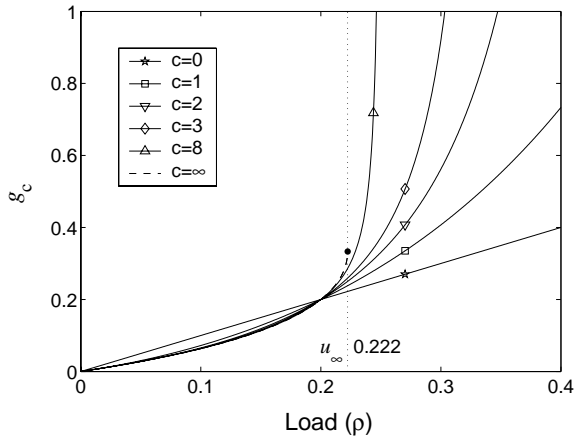
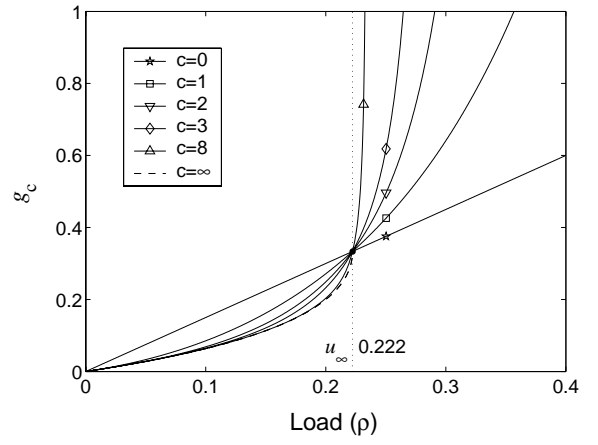
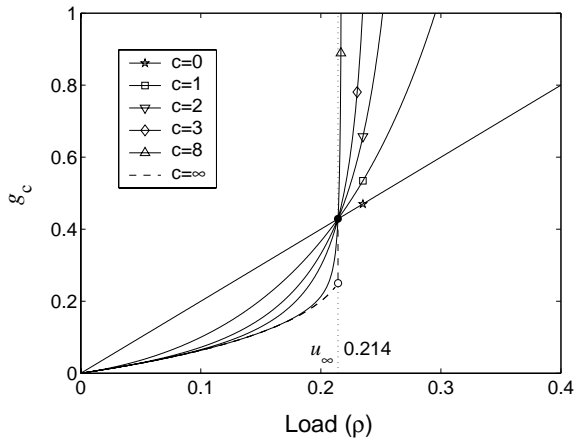
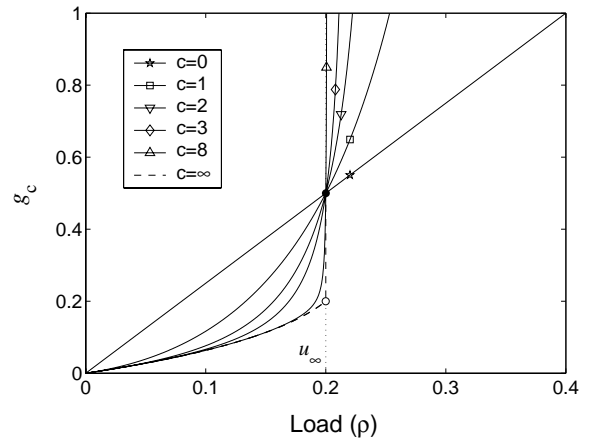
(a)  $A = 1$ .(b)  $A = 2$ .(c)  $A = 3$ .(d)  $A = 4$ .(e)  $A = 5$ .(f)  $A = 6$ .

Fig. 3. The deflection base-functions  $g_c(\rho)$  for  $c = 0, 1, 2, 3, 8, \infty$ , and  $\rho \in \mathcal{R}_c$ , for  $A = 1, 2, 3, 4, 5, 6$ .

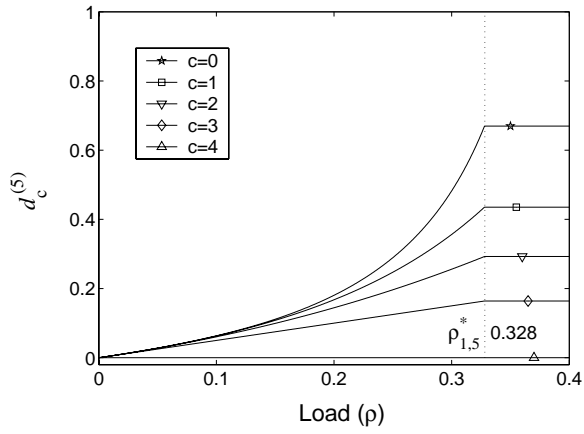
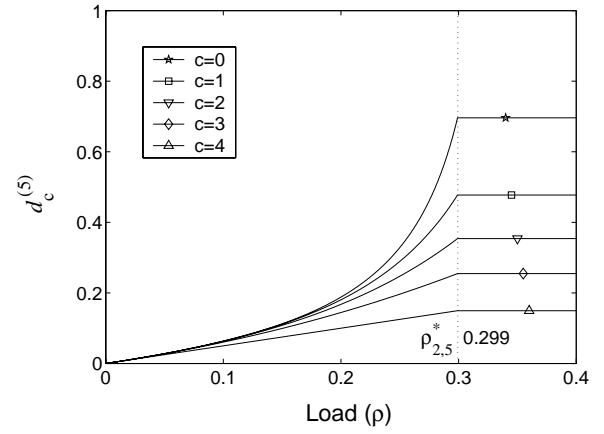
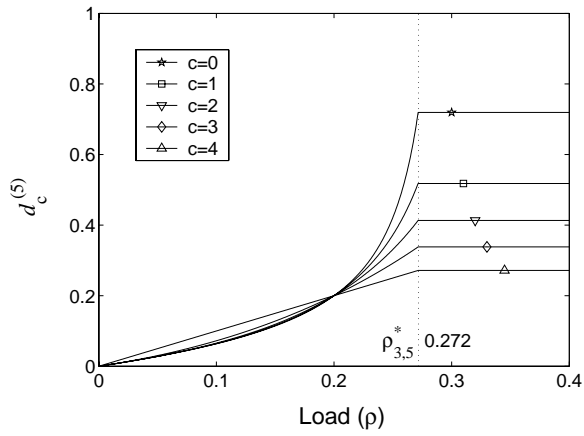
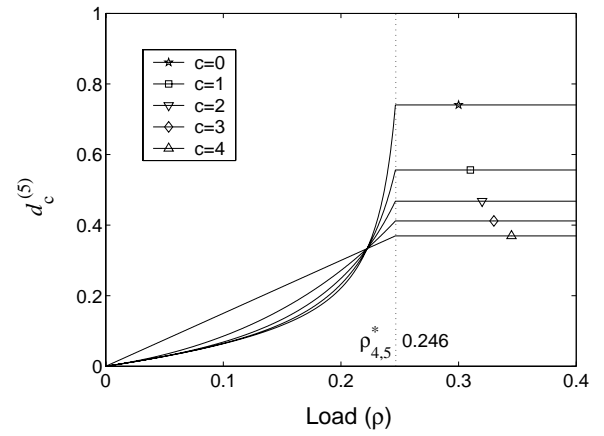
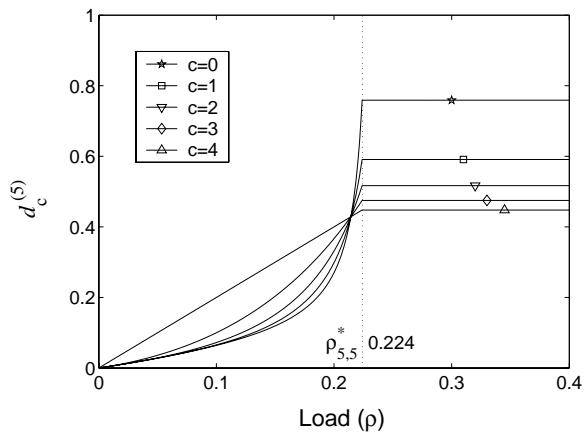
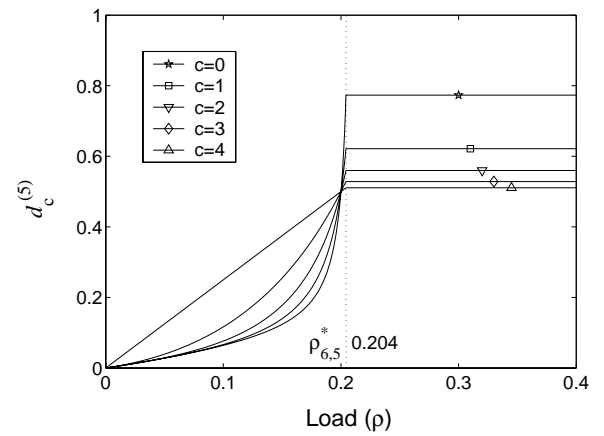
(a)  $A = 1$ .(b)  $A = 2$ .(c)  $A = 3$ .(d)  $A = 4$ .(e)  $A = 5$ .(f)  $A = 6$ .

Fig. 4. The deflection probabilities  $d_c^{(5)}(\rho)$  for  $C = 5$ ,  $c = 0, 1, 2, 3, 4$ , and  $A = 1, 2, 3, 4, 5, 6$ .

of the type of the control signals received at the preceding time slots. This implies that a packet at the head-of-line of the input interface receives at any time slot a *permit* signal with probability  $1 - d_0$ , and that its service time  $X$  is geometrically distributed with  $P(X = k) = (1 - d_0)d_0^{k-1}$ , for  $k = 1, 2, \dots$ . Thus, the mean input delay  $\overline{D}_i$  is obtained from a discrete-time Geo/G/1 queueing system as follows [6, Eq.(1.13)]:

$$\overline{D}_i(\rho) = \frac{1 - \rho}{1 - d_0(\rho) - \rho} - 1 = \frac{d_0(\rho)}{1 - d_0(\rho) - \rho}. \quad (18)$$

The mean in-fabric delay  $\overline{D}_f$  is given by  $\overline{D}_f = \sum_{c=0}^{C-1} (N_c + 1)$  which by virtue of (1) yields

$$\overline{D}_f(\rho) = C + \frac{1}{\rho} \sum_{c=0}^{C-1} d_c(\rho). \quad (19)$$

**Remark 2.** From (18) and Proposition 1 it follows that  $\overline{D}_i(0) = 0$ . From (19) and (6) and given that  $g'_0(0) = \frac{A-1}{2}$  and  $g'_c(0) = \frac{1}{2}$  for  $c \in \mathbb{N}$ , it follows that  $\overline{D}_f(0) = C + \sum_{c=0}^{C-1} d'_c(0) = C + \frac{C-1}{2} + \frac{A-1}{2}$ . The first term represents the traversing of cylinders, the second term the deflections due to height address routing and the third term the angle routing at the innermost  $CL_{C-1}$ . Consequently, for uniform random traffic at low input load, the in-fabric delay predominates the total delay given by  $\overline{D}(0) = C + \frac{C-1}{2} + \frac{A-1}{2} = \frac{3C+A}{2} - 1$ .

#### IV. SATURATION ANALYSIS

We now examine the impact of  $A$  and  $C$  on the maximum or saturation throughput  $\rho_{A,C}^*$  of the switch. We also derive the values of  $A$  and  $C$  corresponding to a given switch size  $N$  that maximize the saturation throughput. Equation (18) implies that the system is stable iff  $1 - d_0^{(C)}(\rho) - \rho > 0$ . Consequently, at saturation it holds that

$$1 - \rho_{A,C}^* - d_0^{(C)}(\rho_{A,C}^*) = 0. \quad (20)$$

**Remark 3.** From (6) and (20) it follows that  $\rho_{A,C}^*$  is a root of the function  $f_C(\rho)$  defined as follows:

$$f_C(\rho) \triangleq 1 - \rho - g_{C-1}(\rho), \quad \forall C \in \mathbb{N}. \quad (21)$$

In particular, owing to Lemma 1, the function  $f_C(\rho)$  is strictly decreasing in  $\rho$ , implying that it has at least one root in  $(0, 1]$ . Furthermore,  $\rho_{A,C}^*$  is obtained as the smallest positive root of  $f_C(\rho)$ .

In general, the saturation throughput  $\rho_{A,C}^*$  depends on both  $A$  and  $C$ . It is of interest to determine the pairs  $(A, C)$  for which the saturation throughput of a switch of a given size  $N$  is maximized. Note that  $N = AH$ , with  $C = 1 + \log_2(H)$  or  $H = 2^{C-1}$ . We therefore consider the following maximization problem

$$\max_{(A,C)} \rho_{A,C}^* \quad \text{s.t. } A \times 2^{C-1} = N \quad (22)$$

The maximum throughputs corresponding to various switch sizes  $N$  ( $N = 2^k$ ) and angles  $A$  are listed in Table I. The results lead us to conjecture that, for a fixed angle, the maximum throughput decreases as the switch size increases. Is this conjecture, however, also valid for large  $N$ ? If it indeed is, does the maximum

TABLE I  
MAXIMUM THROUGHPUT  $\rho_A^*$  UNDER CONSTANT  $N$ .

$A \setminus N$	2	4	8	16	32	64	128	256	512
1	0.666	0.472	0.377	0.328	0.299	0.281	0.268	0.260	0.253
2	0.666	0.472	0.377	0.328	0.299	0.281	0.268	0.260	0.253
4	—	0.400	0.309	0.273	0.256	0.246	0.240	0.236	0.233
8	—	—	0.222	0.187	0.177	0.174	0.172	0.171	0.171
16	—	—	—	0.117	0.106	0.104	0.103	0.103	0.103
32	—	—	—	—	0.060	0.057	0.056	0.056	0.056
64	—	—	—	—	—	0.030	0.029	0.029	0.029
128	—	—	—	—	—	—	0.015	0.015	0.015
256	—	—	—	—	—	—	—	0.007	0.007
512	—	—	—	—	—	—	—	—	0.003

throughput reduce to zero or does it asymptotically approach another value? Deriving the maximum throughput for large values of  $N$  (or  $C$ ) becomes, however, a formidable task, because it amounts to finding the roots of polynomials whose degree as well as coefficients are extremely large: the degree grows linearly in  $C$ ; the coefficients are given as  $C$ -th powers of integers. Consequently, answers can only be found by means of analysis. It turns out that the conjecture indeed holds and that the maximum throughput does not reduce to zero. In particular, for  $A = 1$ , as  $C$  increases, the maximum throughput decreases, approaching 0.222, as suggested by the following general theorems.

*Theorem 1:* It holds that

$$\rho_\infty^*(A) < \cdots < \rho_{A,C+1}^* < \rho_{A,C}^* < \cdots < \rho_{A,1}^* = \frac{2}{A+1}, \quad \forall A \in \mathbb{N}, \quad (23)$$

where

$$\rho_\infty^*(A) = u_\infty = \begin{cases} \frac{2}{9} = 0.222\dots, & \text{for } 1 \leq A \leq 4, \\ \frac{2(A-2)}{(A+2)(A-1)} = \hat{\rho}_A, & \text{for } A \geq 4. \end{cases} \quad (24)$$

*Proof:* See Appendix B. ■

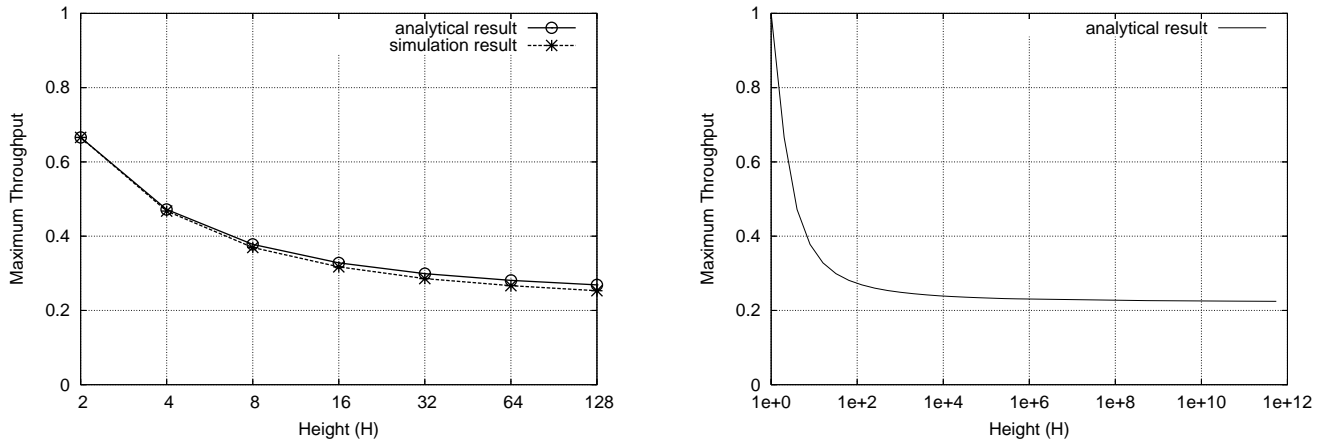
*Theorem 2:* It holds that  $\lim_{C \rightarrow \infty} \rho_{A,C}^* = \rho_\infty^*(A)$ ,  $\forall A \in \mathbb{N}$ .

*Proof:* See Appendix B. ■

*Theorem 3:* The sequence  $\{\rho_{A,C}^*\}$  decreases and approaches 0 as  $A$  increases, i.e.  $\rho_{1,C}^* > \cdots > \rho_{A,C}^* > \cdots > \rho_{A+1,C}^* > \cdots$ , with  $\lim_{A \rightarrow \infty} \rho_{A,C}^* = 0$ ,  $\forall C \in \mathbb{N}$ .

*Proof:* Immediate from Remark 1, Lemma 3, and (23). ■

*Corollary 4:* At saturation, the largest deflection probability is the one at the outermost  $\text{CL}_0$ .



(a) Analytical vs. simulation results.

(b) Maximum throughput as a function of height  $H$ .Fig. 5. Analytical and simulation results for the maximum throughput under uniform traffic for  $A = 1$ .

*Proof:* Immediate consequence of (23), (24), (10) and (6). ■

Furthermore, the results reveal that the maximum throughput increases as the angle decreases and is maximized when  $A = 1, 2$ . The following theorem states that this holds for any typical switch size  $N$ , i.e.

**Theorem 4:** It holds that  $\rho_{(N,1)}^* < \dots < \rho_{(A, \log_2(2N/A))}^* < \dots < \rho_{(4, \log_2(N)-1)}^* < \rho_{(2, \log_2(N))}^* = \rho_{(1, \log_2(N)+1)}^*$ ,  $\forall N = 2^k$ , ( $k \in \mathbb{N}$ ).

*Proof:* See Appendix B. ■

Fig. 5a compares the results for the switch throughput obtained using the analytic procedure presented in this section, with the results obtained by simulation (see Section V), for  $A = 1$  and  $H \in [2, 128]$ . The two plots are in very good agreement. Fig. 5b plots the analytically obtained switch throughput for even higher switch heights: the maximum throughput decreases with increasing  $H$ , and approaches 0.222.

## V. NUMERICAL RESULTS

We developed a model using the OMNeT++ discrete event simulation environment [7] to simulate the Data Vortex switch. The two primary measures of interest are the mean packet delay and the maximum switch throughput. The results were obtained using the Akaroa environment [8] with a statistical confidence of 90% and a precision of 1%. We consider Bernoulli i.i.d. uniform as well as nonuniform traffic. Furthermore, as the fabric does not guarantee in-order delivery, packets may be delivered out-of-order to the output interfaces. This issue is addressed in Section V-D.

### A. Throughput vs. Switch Size

#### A.1 Uniform Traffic

First, we measure the switch throughput under uniform traffic as a function of the switch dimensions  $A$  and  $H$ . Fig. 6a shows the switch throughput versus  $H$  for various  $A$ . Among the configurations simulated,

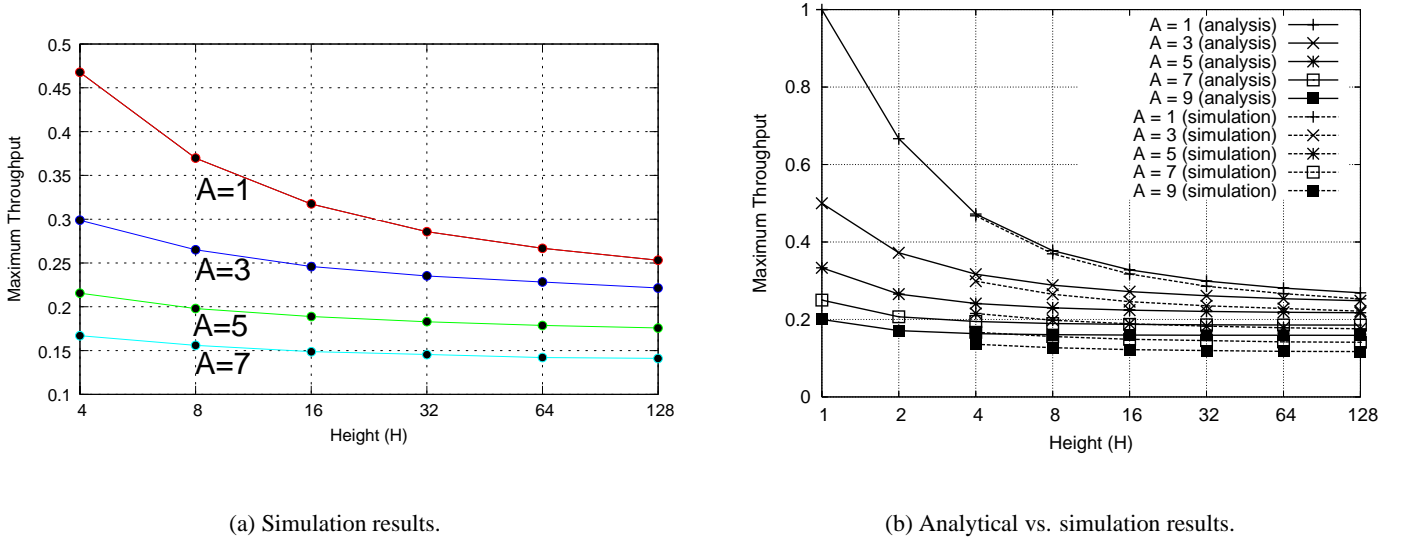


Fig. 6. Maximum throughput as a function of angle  $A$  and height  $H$ .

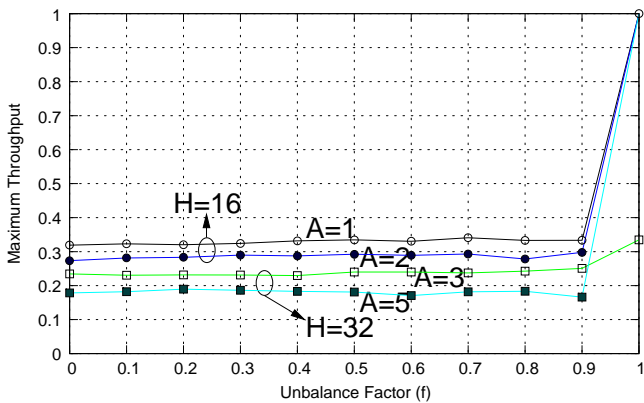
the highest throughput (0.46) is achieved when  $A = 1$  and  $H = 4$ , i.e., by a  $4 \times 4$  switch, whereas the lowest throughput is close to 0.12. These values correspond to the throughput results obtained in Section IV. The figure confirms that the throughput of the Data Vortex switch decreases with increasing  $H$  (or  $C$ ), as stipulated by Theorem 1. As the height increases, the number of cylinders increases, which introduces more levels of deflections, thus reducing the switch throughput. The throughput also decreases with increasing  $A$ , as stipulated by Theorem 3. This is due to angle decoding in the innermost  $CL_{C-1}$ . As the average number of deflections there increases and these deflections prevent other packets from moving SOUTH, they tend to propagate outward, finally reaching the outermost  $CL_0$ , thus reducing the injection rate and the throughput.

Fig. 6b compares the simulation and analytic results for  $A = 1, 3, 5, 7$ , and  $9$ . We see that the prediction of the analysis becomes less accurate as  $A$  increases; the maximum deviation is about 0.042 ( $A = 9, H = 128$ ). This is due to the independence assumption used in the derivation of (2) that ignores existing correlations, which apparently become more pronounced with increasing  $A$ .

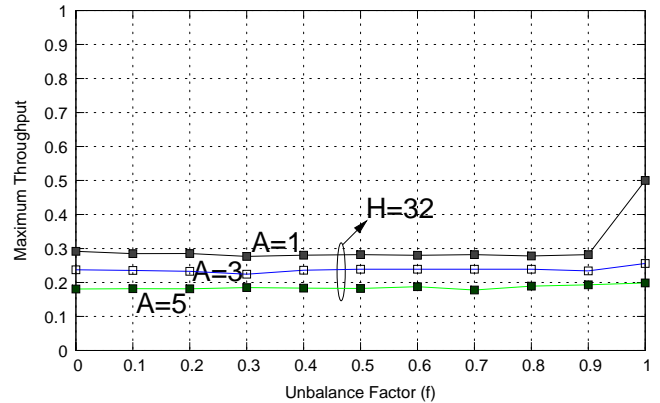
## A.2 Unbalanced Traffic

Next, we apply nonuniform—also called unbalanced—traffic to the switch and measure the switch throughput. The degree of unbalance is specified by the unbalance factor  $f$  [9]: input  $i$  sends to a predetermined favorite output  $\text{fav}(i)$ , with probability  $f + \frac{1-f}{N}$ , and to each one of the other outputs with probability  $\frac{1-f}{N}$ . The pattern is admissible for any value of  $f$  if each input has a different favorite output, i.e., if  $i \neq j \Leftrightarrow \text{fav}(i) \neq \text{fav}(j)$ . Fig. 7 illustrates the maximum switch throughput versus  $f$  for various switch sizes. For  $f = 0$ , the traffic is uniform, whereas for  $f = 1$ , it is completely unbalanced, i.e., each input sends traffic only to its favorite output.

We used two different mapping functions to determine  $\text{fav}(i)$ . The first one,  $\text{Id}$ , is the identity function:



(a) Id mapping function.



(b) Rev mapping function.

Fig. 7. Maximum throughput under unbalanced traffic.

$\text{Id}(a, h) = (a, h)$ . Fig. 7a presents the results for unbalanced traffic using the Id mapping function. Two aspects are worth pointing at: first, the throughput is largely immune to  $f$ , in the sense that it does not degrade compared with the uniform case. This behavior is not commonly found in other switch architectures, in which the throughput usually drops under unbalanced traffic [9]. Second, we see, unexpectedly, that with  $f = 1$  the normalized switch throughput is sometimes well below 100%. This is due to an internal path conflict that can occur even when there is no output contention, as explained below.

In general, 100% throughput is achieved if and only if no packet deflections occur in the switch, which means that packets follow only SOUTH paths. Therefore, each input sends packets to a *single* (favorite) output, implying that  $f = 1$ . Let  $\text{Map}(a, h) = (a_m, h_m)$  be the mapping function that determines the favorite output of a packet that starts from angle  $a$  and height  $h$ . When it reaches the innermost  $\text{CL}_{C-1}$  following only SOUTH paths, first, its height has not changed, i.e.  $h_m = h$ , and second it has traversed  $C - 1$  hops (cylinders) and therefore its angle  $a'$  is equal to  $(a + C - 1) \bmod A$ . Moreover, this angle is precisely its destination angle so that the next hop is SOUTH, i.e.  $a_m = a'$ . Recalling that  $C - 1 = \log_2(H)$ , we get  $a_m = (a + \log_2(H)) \bmod A$ . Consequently, 100% throughput is achieved if and only if the traffic is completely unbalanced, i.e.  $f = 1$ , and  $\text{Map}(a, h) = ((a + \log_2(H)) \bmod A, h)$ . Note also that  $\text{Id}(a, h) = \text{Map}(a, h)$  iff  $\log_2(H) = kA$  for some  $k \in \mathbb{N}_0$ . We now verify that the latter condition is satisfied by the  $(H = 2^k, A = 1)$ ,  $(H = 16, A = 2)$  and  $(H = 32, A = 5)$  configurations. However, it is not satisfied by the  $(H = 32, A = 3)$  configuration, so in this case the switch throughput with  $f = 1$  and the identity mapping is well below 100%.

Fig. 7b uses a different mapping function, namely, the Rev function:  $\text{Rev}(a, h) = (a, (h + 1) \bmod H)$ . In this configuration, deflections due to height translation occur at almost all cylinder levels, and the switch throughput does not reach 100% for any  $f$ . However, if we exclude the fully unbalanced case, the results here are virtually identical to those obtained with the Id mapping function.



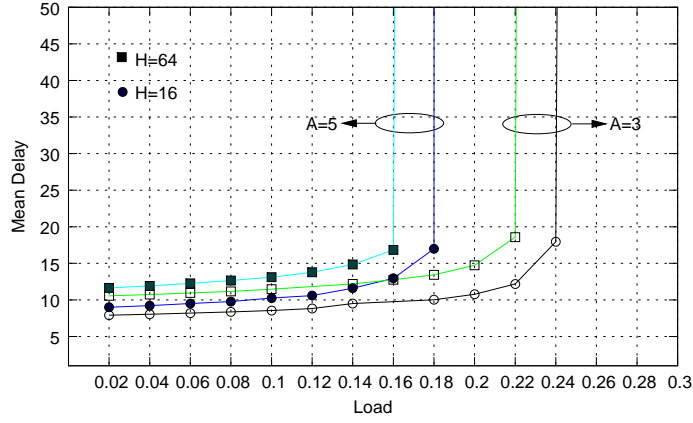
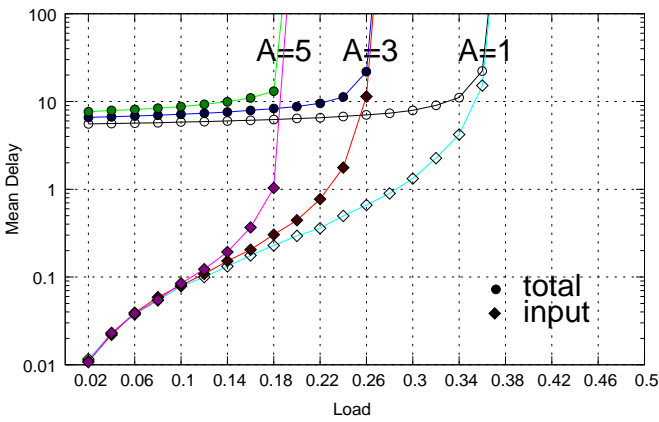
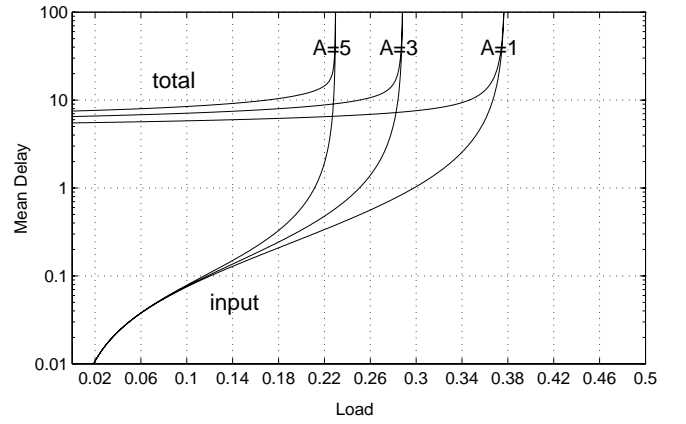


Fig. 8. Average packet delay in number of time slots. Bernoulli i.i.d. uniform traffic.



(a) Simulation.



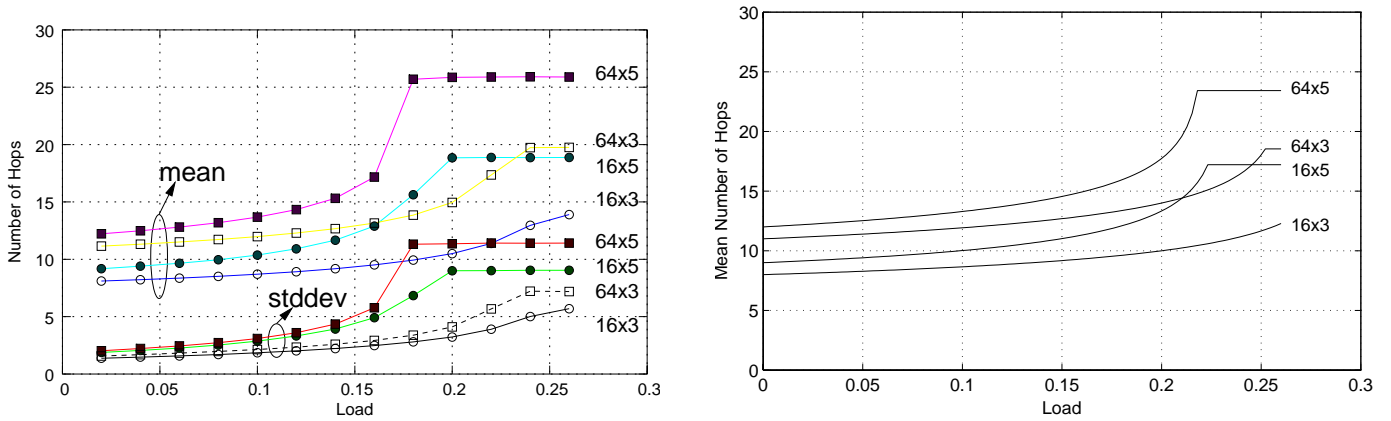
(b) Analysis.

Fig. 9. Mean total and input delay in number of time slots under Bernoulli i.i.d. uniform traffic ( $H = 8$ ).

## B. Packet Delay vs. Load

Next, we plot the average packet delay for different input loads under Bernoulli i.i.d. uniform traffic. Fig. 8 contains four plots corresponding to the average packet delay for four different switch size configurations:  $(H, A) = (16, 3), (64, 3), (16, 5),$  and  $(64, 5)$ . It shows that at low input load (0.02–0.04) the mean packet delay is very close to the value  $D^{\min} = C + \frac{A-1+\log_2(H)}{2}$  obtained by Remark 2. With increasing input load, the average delay increases slowly at first, then very sharply as the switch approaches its saturation point. According to Remark 1, this is due to the sharp increase of deflections at  $CL_0$ . Note that the  $(H = 64, A = 3)$  switch can support more ports ( $N = 192$ ) than the  $(H = 16, A = 5)$  switch ( $N = 80$ ) and, in addition, has better performance. This again confirms that in order to increase the size of the Data Vortex switch, it is generally better to increase  $H$  and keep  $A$  as small as possible, as stipulated by Theorem 4.

Fig. 9 plots the mean total packet delay (including the delays in the input buffers and in the fabric) and the mean packet delay in the input buffers separately. It shows that at low input load the input delay is almost zero; as the input load approaches saturation, the input delay predominates the aggregate packet delay. In



(a) Simulation results, mean and standard deviation.

(b) Analytic results, mean.

Fig. 10. Mean and standard deviation of the number of hops under Bernoulli i.i.d. uniform traffic.

fact at saturation, the input delay is infinite whereas the in-fabric delay is finite. We again note the excellent agreement between the simulation and analytic results for small  $A$ , i.e.  $A = 1$ .

In another set of experiments, we used virtual output queues (VOQ) at the input interfaces instead of FIFOs and round-robin scheduling to select one of the eligible VOQs. However, we found no significant throughput or delay improvements. This was to be expected as the deflections reaching the input interfaces are indiscriminate with respect to packet destination.

### C. Hop-Count Distribution

In this section we study the distribution of the number of hops that packets traverse in the fabric, i.e., the in-fabric packet delay, under uniform traffic. Fig. 10a plots the mean and standard deviation of the per-packet hop count for different switch configurations as a function of the traffic load at the inputs. The overall trend is clear: larger  $A$  and higher input load increase the number of hops being traversed. However, a closer look reveals that the average number of hops increases slowly before the switch reaches saturation; it increases sharply at the saturation point, and then stays almost constant even as the input load increases further. This result can also be verified in [3, Fig. 4] and is in agreement with the analytic results plotted in Fig. 10b according to (19). The standard deviation is relatively small, indicating a relatively narrow distribution.

Fig. 11 plots the probability density of the hop-count distribution for  $A = 1, 3$ , and  $5$  for input loads of  $0.1$  and  $1.0$ . We see that the probability of traversing a large number of hops is quite low, even at high load. However, at high load, the distributions have a long tail, and packets can traverse as many as  $150$  hops before departing the fabric. These points are not shown in the figure, because their probability is extremely low (smaller than  $10^{-6}$ ). In Fig. 12 we plot the complementary cumulative distribution and the probability density for different  $A$ . The input load used in each configuration is just a few percent below the saturation point of the switch. Specifically, the input load values are  $0.25$ ,  $0.16$ , and  $0.11$  for  $A = 1$ ,  $A = 3$ , and  $A = 5$ ,

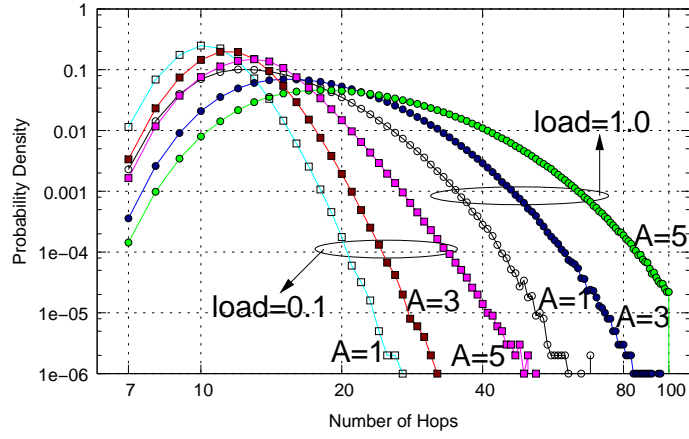


Fig. 11. Probability density of hop-count distribution; y- and x-axis in logarithmic scale; Bernoulli i.i.d. uniform traffic;  $H = 64$ .

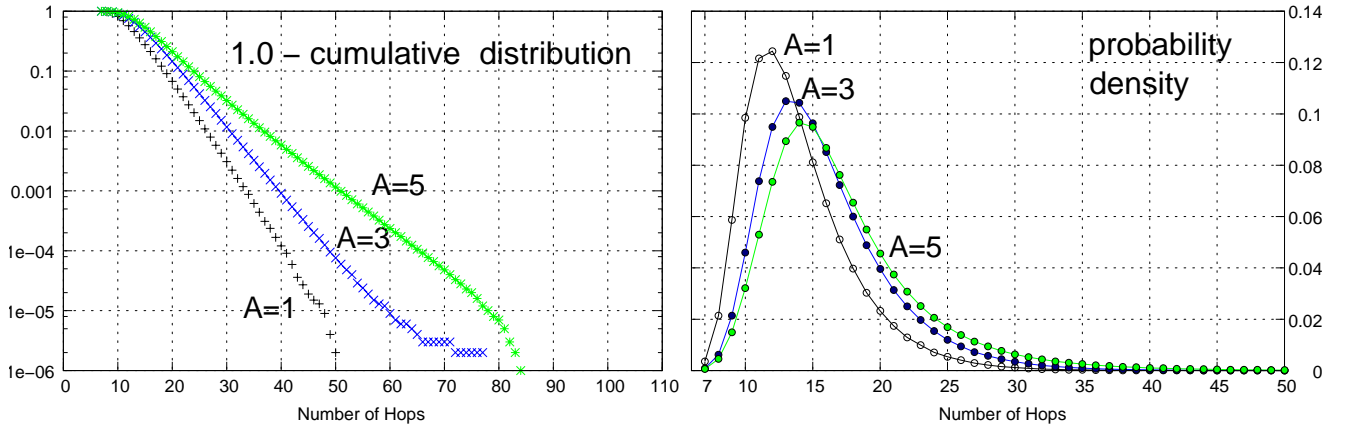


Fig. 12. Cumulative distribution and probability density of number of hops; Bernoulli i.i.d. uniform traffic;  $H = 64$ .

respectively. We see that more than 99% of the packets have an in-fabric delay of fewer than 40 cells.

Next we use the unbalanced traffic pattern of Section V-A.2 and again plot measurements for the distribution of the in-fabric packet delay. Fig. 13 presents the probability density for  $H = 64$ ,  $A = 3$  and 5, under the unbalanced scenario using the Id mapping function and  $f = 0.5$ . Both plots contain some form of periodicity: every  $A$  hops, the probability density exhibits a peak. This can be explained by the presence of the heavy flows under the unbalanced scenario. These flows, which have the form  $(a, h) \rightarrow (a, h)$ , must follow exactly  $kA$ ,  $k \in \mathbb{N}$ , hops before they reach their destination, which explains the higher probability of these points on the  $x$ -axis of Fig. 13.

#### D. Resequencing Delay

The basic Data Vortex routing mechanism can deliver packets out of order to the output interfaces. This requires the addition of *resequencing* buffers and logic at every output interface to temporarily store out-of-order packets until the missing ones arrive. Every output interface has  $N$  resequencing queues, one for every switch input. Every queue has a separate buffer space  $B_q$ . A packet is stored in the resequencing queue corresponding to its input in a position determined by its sequence number. A resequencing queue is eligible

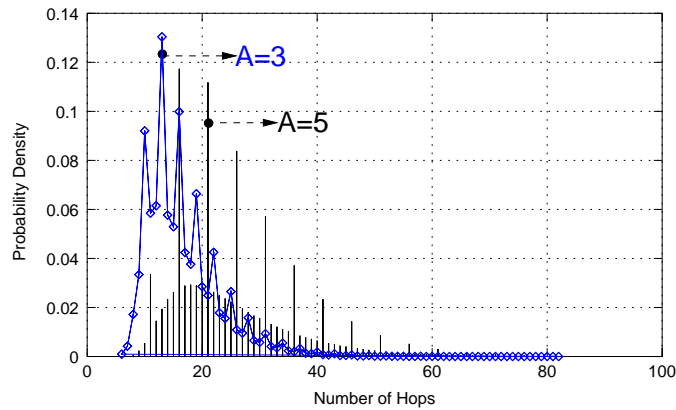


Fig. 13. Hop-count distribution under unbalanced traffic;  $H = 64$ .

if it contains the packet with the sequence number (of the corresponding input) expected next by the output. A round-robin scheduler is used at every output interface to select one of the eligible resequencing queues.

A scheme in which the output interface exerts backpressure towards  $CL_{C-1}$  by using the deflection mechanism could lead to a deadlock. A deadlock would occur when the output buffer is full with out-of-order cells; none of them can depart from the buffer, and no cells can enter the buffer because there is no free space left. We therefore consider an alternative scheme based on a per-output-credit flow control. An end-to-end credit-based flow-control mechanism guarantees that the resequencing buffers do not overflow: at the input interface, a packet is injected into the switching fabric only if buffer space (credits) for the corresponding output is available. Whenever a packet from input  $i$  is transmitted from the output interface, a credit is sent back to input  $i$ . For simplicity, we assume in our model that credit communication is performed through a point-to-point network connecting each output interface directly to all input interfaces, so that credits do not experience any contention; the credit propagation delay is set to be equal to the minimum delay that a packet can experience when traveling from an input interface to an output interface.<sup>1</sup> As packets destined to different outputs could block each other if they were placed in a common queue at the input interface, we used multiple input queues, one for each output (VOQ), and a round-robin scheduler to select one of those eligible.

Fig. 14 shows the resequencing delay as a function of the unbalance factor  $f$  for various switch configurations. Each resequencing queue has 16 packets' worth of buffer space ( $B_q = 16$ ), which is slightly larger than the round-trip time in packets for the  $H = 64$  configurations. In our model, sufficient credits are available at the input side for every output to achieve full link utilization with an uncongested persistent flow. The results were extracted after 10,000 packets had been sent from each input. For uniform traffic ( $f = 0$ ), the resequencing delay is almost zero, i.e. there is almost no out-of-order delivery. We attribute this to the round-robin VOQ service, which, with uniform Bernoulli traffic at 100% load, implies that each flow is served approximately once every  $N$  time slots. This means that packets of the same flow receive a temporal spacing of  $N$  time slots. Because the probability that the in-fabric latency is in general greater than  $N$  time slots is

<sup>1</sup>This is the ideal credit scenario in terms of performance. All other solutions, such as sending packets through the Data Vortex, piggybacked or through separate messages, will perform worse.

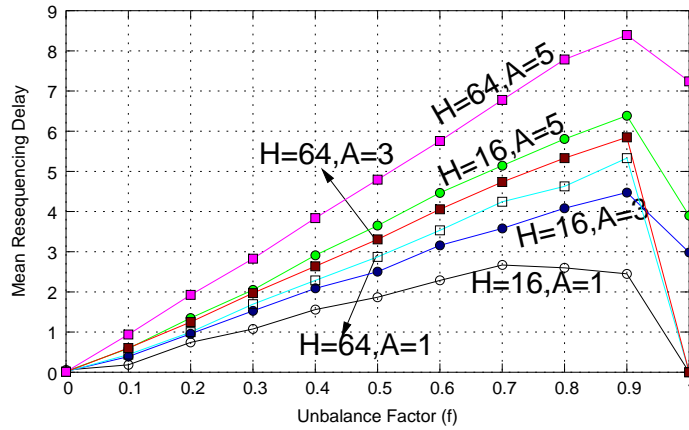


Fig. 14. Resequencing delay in number of time slots; unbalanced traffic with Id mapping function; input load 100%.

very low (see Fig. 10a and 12), the preceding packet of a given flow will almost always already have arrived at the output interface before the next packet is injected. However, this does not hold for nonuniform traffic, as the figure illustrates. The resequencing delay grows as  $f$  increases, owing to the heavy flows that have more packets concurrently pending inside the switching fabric, thus creating the necessary conditions for out-of-order arrivals at the output interfaces. With  $f = 1$  the configurations  $(H = 16, A = 1)$ ,  $(H = 64, A = 1)$  and  $(H = 64, A = 3)$  incur no deflections (see Section V-A.2, for the Id function), and hence the resequencing delay drops to zero. Finally, note that the resequencing delay is small compared with the total delay given in Section V-B.

## VI. CONCLUSIONS

The Data Vortex architecture is a promising choice for all-optical switching fabrics, as it is scalable to thousands of ports. We developed an analytical model to study its performance under uniform traffic and Bernoulli arrivals. We derived closed-form expressions for the mean packet delay through the switching fabric and for its maximum throughput. Our analysis has shown that the switch performance scales better with the height parameter  $H$  than with the angle parameter  $A$ : As  $A$  increases, the maximum throughput decreases and approaches zero, whereas as  $H$  increases, the maximum throughput decreases, approaching a positive lower bound. In particular, for  $A = 1$ , a switch of any size can sustain a throughput of 22%. We also demonstrated that for any typical switch size, the saturation throughput is maximized for  $A = 1$ . Our simulation results confirm the analytical findings. Simulation results were also obtained for nonuniform traffic, revealing that the maximum throughput is largely insensitive to nonuniformity. The latency results demonstrate that the in-fabric delay distribution is quite narrow; moreover, we have observed that with uniform traffic virtually no misordering occurs, although the Data Vortex routing mechanism does not guarantee in-order delivery. However, misordering increases with nonuniformity. Finally, we find no throughput or latency benefits in using virtual output queues at the input interfaces. There is, however, an advantage to VOQs in the context of resequencing.

## ACKNOWLEDGMENTS

The authors gratefully acknowledge the help and input of Benjamin Small and Prof. Keren Bergman of the Lightwave Research Laboratory, Columbia University, New York.

APPENDIX A  
DEFLECTION PROBABILITIES

**Proof of Proposition 1.**

We begin by noting that, according to (5), the deflection probability of the innermost  $\text{CL}_{C-1}$  is independent of the number of cylinders  $C$ . This leads us to the introduction of the following variables:

$$g_c(\rho) \triangleq d_{C-1-c}^{(C)}(\rho), \quad \text{implying that } d_c^{(C)}(\rho) = g_{C-1-c}(\rho), \quad \forall c \in \{0, 1, \dots, C-1\}, \quad \rho \leq \rho_{A,C}^*. \quad (25)$$

From (25), (4) and (5) we obtain the following recursive relation:

$$g_c(\rho) = \frac{\rho(3g_{c-1}(\rho) + 1)}{2(1 - g_{c-1}(\rho))}, \quad \forall c \in \{1, 2, \dots, C\}, \quad \text{with } g_0(\rho) = M\rho, \quad \forall \rho \leq \rho_{A,C}^*, \quad (26)$$

where 
$$M \triangleq \frac{A-1}{2}. \quad (27)$$

**Remark 4.** Note that  $\forall \rho \in \mathcal{R}_c$  such that  $g_c(\rho) < 1$ , by virtue of (26) it also holds that  $g_{c-1}(\rho) < 1$ , which implies that  $\rho \in \mathcal{R}_{c-1}$ . Furthermore,  $g_{c-1}(u_c) < 1$  as  $g_c(u_c) = 1$  for  $c \geq 2$ . Consequently,  $u_c < u_{c-1}$  and  $\mathcal{R}_c \subset \mathcal{R}_{c-1}$  for  $c \geq 2$ . This also holds for  $c = 1$  provided that  $A \geq 2$ , whereas for  $c = A = 1$  it holds that  $u_0 = u_1 = 1$  and  $\mathcal{R}_0 = \mathcal{R}_1 = [0, 1]$ .

**Remark 5.** The above formulation suggests that the values of  $g_c(\rho)$  for given  $A$  (or  $M$ ),  $C$  and  $\rho$  are independent of  $C$ . Therefore, we proceed to derive the entire sequence  $\{g_c\}$  by simply assuming that  $C = \infty$ . We also consider the extended domain  $\rho \in [0, 1]$  instead of  $\rho \in [0, \rho_{A,C}^*]$ .

From (26) it follows that  $g_c(\rho)$  ( $c = 0, 1, \dots$ ) can be written as the ratio of two polynomials  $p_c(\rho)$  and  $q_c(\rho)$  of degree  $c + \mathbf{1}_{\{M \neq 0\}}$  and  $c$ , respectively, i.e.

$$g_c(\rho) = \frac{p_c(\rho)}{q_c(\rho)}, \quad \text{for } c = 0, 1, \dots, \quad (28)$$

which satisfy the following set of recurrence equations:

$$\left\{ \begin{array}{l} p_c(\rho) = \rho(3p_{c-1}(\rho) + q_{c-1}(\rho)) \\ q_c(\rho) = 2(q_{c-1}(\rho) - p_{c-1}(\rho)) \end{array} \right\} \quad \text{for } c = 1, 2, \dots, \quad \text{with } \left\{ \begin{array}{l} p_0(\rho) = M\rho \\ q_0(\rho) = 1 \end{array} \right\}. \quad (29)$$

We proceed to derive closed-form expressions for these polynomials using the generating function approach. Let us define the following generating functions:

$$P(z, \rho) \triangleq \sum_{c=0}^{\infty} p_c(\rho) z^c \quad \text{and} \quad Q(z, \rho) \triangleq \sum_{c=0}^{\infty} q_c(\rho) z^c. \quad (30)$$

By making use of (30), (29) yields

$$\left\{ \begin{array}{l} P(z, \rho) - p_0(\rho) = \rho z(3P(z, \rho) + Q(z, \rho)) \\ Q(z, \rho) - q_0(\rho) = 2z(Q(z, \rho) - P(z, \rho)) \end{array} \right\}. \quad (31)$$

Solving the above system of equations for  $P(z, \rho)$  and  $Q(z, \rho)$  gives

$$P(z, \rho) = \frac{\rho[M - (2M - 1)z]}{8\rho z^2 - (2 + 3\rho)z + 1} \quad \text{and} \quad Q(z, \rho) = \frac{1 - (2M + 3)\rho z}{8\rho z^2 - (2 + 3\rho)z + 1}. \quad (32)$$

Let us consider the polynomial of the denominators

$$\Phi(z) \triangleq 8\rho z^2 - (2 + 3\rho)z + 1, \quad (33)$$

and denote by  $z_1$  and  $z_2$  its two roots given by (8). Note that for  $\rho$  in the interval  $[0, \frac{2}{9})$  the two roots are real numbers, with  $0 < z_2 < z_1$ , whereas for  $\rho$  in the interval  $(\frac{2}{9}, 1]$  the discriminant  $\Delta$  is negative and therefore the two roots are complex conjugates given by

$$z_{1,2} = \frac{2 + 3\rho \pm i\sqrt{-\Delta}}{16\rho} = \frac{1}{\sqrt{8\rho}} e^{\pm i\theta}, \quad \text{with } \theta = \tan^{-1} \left( \frac{\sqrt{-\Delta}}{2 + 3\rho} \right). \quad (34)$$

The generating functions (32) can now be written as follows:

$$P(z, \rho) = \begin{cases} \sum_{i=1}^2 \frac{a_i}{1 - 8\rho z_i z}, & \text{for } \rho \neq \frac{2}{9}, \\ \frac{\frac{2}{9}[M - (2M - 1)z]}{(1 - \frac{4}{3}z)^2}, & \text{for } \rho = \frac{2}{9}, \end{cases} \quad Q(z, \rho) = \begin{cases} \sum_{i=1}^2 \frac{b_i}{1 - 8\rho z_i z}, & \text{for } \rho \neq \frac{2}{9}, \\ \frac{1 - \frac{2(2M+3)}{9}z}{(1 - \frac{4}{3}z)^2}, & \text{for } \rho = \frac{2}{9}, \end{cases} \quad (35)$$

where  $a_i$  and  $b_i$  ( $i = 1, 2$ ) are given by (9).

Inverting (35) yields the result sought:

$$p_c(\rho) = \begin{cases} (8\rho)^c (a_1 z_1^c + a_2 z_2^c), & \text{for } \rho \neq \frac{2}{9}, \\ \frac{2}{9} \left( \frac{3-2M}{4} c + M \right) \left( \frac{4}{3} \right)^c, & \text{for } \rho = \frac{2}{9}, \end{cases} \quad q_c(\rho) = \begin{cases} (8\rho)^c (b_1 z_1^c + b_2 z_2^c), & \text{for } \rho \neq \frac{2}{9}, \\ \left( \frac{3-2M}{6} c + 1 \right) \left( \frac{4}{3} \right)^c, & \text{for } \rho = \frac{2}{9}, \end{cases} \quad c \in \mathbb{N}_0. \quad (36)$$

Substituting (36) into (28), and using (27), yields (7). Also, (6) is a direct consequence of (25) and the fact that the deflection probabilities remain constant for loads exceeding the saturation throughput. ■

### Proof of Lemma 1.

We shall prove the lemma using induction. For  $c = 0$  the inequality holds because according to (26) and (27),  $g_0(\rho_1) - g_0(\rho_2) = (\rho_1 - \rho_2) \frac{A-1}{2} \leq 0$ . Suppose that the inequality holds for  $c = k$ , i.e.  $g_k(\rho_1) \leq g_k(\rho_2)$ . This implies that  $3g_k(\rho_1) + 1 \leq 3g_k(\rho_2) + 1$ ,  $2(1 - g_k(\rho_1)) \geq 2(1 - g_k(\rho_2)) > 0$ , and by making use of (26), we obtain  $g_{k+1}(\rho_1) < g_{k+1}(\rho_2)$ . Thus, the inequality also holds for  $c = k + 1$ . ■

### Proof of Lemma 2.

We shall prove the lemma using induction. For  $c = 1$ , (26) and (27) yield

$$g_1(\rho) - g_0(\rho) = \frac{\rho \Psi(A, \rho)}{2 - (A - 1)\rho}, \quad (37)$$

$$\text{with} \quad \Psi(A, \rho) \triangleq (A+2)(A-1)\rho - 2(A-2). \quad (38)$$

Consequently, by making use of (11), (10) holds because

$$\left\{ \begin{array}{l} \Psi(A, \rho) > 0, \quad \text{for } A = 1, 2 \text{ and } A \geq 3 \text{ with } \rho > \hat{\rho}_A \\ \Psi(A, \rho) = 0, \quad \text{for } A \geq 3 \text{ with } \rho = \hat{\rho}_A \\ \Psi(A, \rho) < 0, \quad \text{for } A \geq 3 \text{ with } \rho < \hat{\rho}_A. \end{array} \right\} \quad (39)$$

and also  $2 - (A-1)\rho > 0$  as  $\rho < u_1 \leq u_0 = \min(1, \frac{2}{A-1})$ . Suppose now that (10) holds for  $c = k$ . From (26) it follows that

$$g_{k+1}(\rho) - g_k(\rho) = \frac{2\rho}{(1 - g_k(\rho))(1 - g_{k-1}(\rho))} (g_k(\rho) - g_{k-1}(\rho)), \quad (40)$$

with the fraction at the right-hand side being positive for  $\rho \in \mathcal{R}_c$ , according to Remark 4. Consequently, (10) also holds for  $c = k + 1$  because the sign of the difference  $g_{k+1} - g_k$  is the same as the one of the difference  $g_k - g_{k-1}$ , which by our assumption is in accordance with (10). ■

#### Proof of Lemma 4.

We shall prove the lemma using induction. From (26) and (27) it follows that  $g_1(\rho; A = 1) = \frac{\rho(3g_0(\rho; A=1)+1)}{2(1-g_0(\rho; A=1))} = \frac{\rho}{2} = g_0(\rho; A = 2)$ . Also,  $g_0(\rho; A = 4) = \frac{3\rho}{2}$  and  $g_1(\rho; A = 2) = \frac{\rho(3\rho+2)}{2(2-\rho)}$  which imply that  $\mathcal{R}_0^{(A=4)} = \mathcal{R}_1^{(A=2)} = [0, \frac{2}{3}]$  given that  $g_0(\rho; A = 4) < 1 \Leftrightarrow g_1(\rho; A = 2) < 1 \Leftrightarrow 0 \leq \rho < \frac{2}{3}$ . Consequently, for  $c = 0$ , (12) holds because  $\frac{\rho(3\rho+2)}{2(2-\rho)} < \frac{3\rho}{2} < 1, \forall \rho \in [0, \frac{2}{3}]$ . Suppose that it holds for  $c = n$ , i.e.  $g_{n+1}(\rho; A = 1) = g_n(\rho; A = 2), \forall \rho \in [0, 1]$ , and  $g_{n+1}(\rho; A = 2) < g_n(\rho; A = 4) < 1, \forall \rho \in \mathcal{R}_n^{(A=4)}$ . We shall show that it also holds for  $c = n + 1$ . By virtue of (26) it follows that  $g_{n+2}(\rho; A = 1) = \frac{\rho(3g_{n+1}(\rho; A=1)+1)}{2(1-g_{n+1}(\rho; A=1))} = \frac{\rho(3g_n(\rho; A=2)+1)}{2(1-g_n(\rho; A=2))} = g_{n+1}(\rho; A = 2)$ . We shall now show that  $g_{n+2}(\rho; A = 2) < g_{n+1}(\rho; A = 4) < 1, \forall \rho \in \mathcal{R}_{n+1}^{(A=4)}$ . By the definition of  $\mathcal{R}_{n+1}^{(A=4)}$ , the right-hand side inequality holds. From Remark 4 we have that  $\forall \rho \in \mathcal{R}_{n+1}^{(A=4)}$ , it holds that  $\rho \in \mathcal{R}_n^{(A=4)}$ . Therefore, from (26) and our assumption it follows that  $g_{n+2}(\rho; A = 2) = \frac{\rho(3g_{n+1}(\rho; A=2)+1)}{2(1-g_{n+1}(\rho; A=2))} < \frac{\rho(3g_n(\rho; A=4)+1)}{2(1-g_n(\rho; A=4))} = g_{n+1}(\rho; A = 4)$ . ■

#### Proof of Lemma 5.

From (26) it follows that the function  $g_\infty(\rho)$  should satisfy the following relation:  $g_\infty(\rho) = \frac{\rho(3g_\infty(\rho)+1)}{2(1-g_\infty(\rho))}$ . Thus,  $g_\infty(\rho)$  is a root of the polynomial  $\Pi(y) \triangleq 2y^2 - (2 - 3\rho)y + \rho$ . Note that the two roots of the polynomial are expressed by  $y_{1,2} = \frac{2-3\rho \pm \sqrt{\Delta}}{4}$ , with  $y_1(\rho)$  and  $y_2(\rho)$  being decreasing and increasing functions in  $\rho$ , respectively, for  $\rho \leq \frac{2}{9}$ . Consequently, the limiting function  $g_\infty(\rho)$  can only exist in an interval  $\mathcal{R}_\infty = [0, u_\infty] \subseteq [0, \frac{2}{9}]$ , with  $g_\infty(\rho) = y_2(\rho) \quad \forall \rho \in [0, u_\infty)$ , owing to Lemma 1. Note also that for  $A \geq 3$ , from (11) and (37) – (40) it follows that  $g_c(\hat{\rho}_A) = \frac{A-2}{A+2}, \forall c \in \mathbb{N}_0$ . Thus,  $g_\infty(\hat{\rho}_A) = \frac{A-2}{A+2}, \forall A \geq 3$ . Moreover,  $y_2(\hat{\rho}_A) = \min(\frac{A-2}{A+2}, \frac{1}{A-1})$ , which is equal either to  $\frac{A-2}{A+2}$  for  $3 \leq A \leq 4$ , or to  $\frac{1}{A-1}$  for  $A \geq 4$ . Thus,  $g_\infty(\hat{\rho}_A) = y_2(\hat{\rho}_A)$  for  $3 \leq A \leq 4$ , and  $g_\infty(\hat{\rho}_A) > y_2(\hat{\rho}_A)$  for  $A \geq 5$ . Therefore, the  $g_\infty(\rho)$  function exhibits a discontinuity at  $\rho = \hat{\rho}_A$  for  $A \geq 5$ , translating into a sharp increase as demonstrated in Fig. 3. Combining the above yields (14) – (16). ■



APPENDIX B  
MAXIMUM THROUGHPUT

**Proof of Theorem 1.**

For ease of reading, we suppress the subscript  $A$  in the following treatment, i.e., we write  $\rho_C^*$  instead of  $\rho_{A,C}^*$ . From Remark 3, and owing to (26) and (27), it follows that the maximum throughput for  $C = 1$  is given as the root of the function  $f_1(\rho) = 1 - \rho - g_0(\rho) = 1 - \rho - M\rho$ , i.e.  $\rho_1^* = \frac{1}{M+1} = \frac{2}{A+1}$ .

First we shall show that

$$\rho_\infty^* < \rho_C^*, \quad \text{for } C = 1, 2, \dots, \quad (41)$$

which according to Remark 3 is equivalent to  $f_C(\rho) > 0$ ,  $\forall C \in \mathbb{N}$ ,  $\rho \in [0, \rho_\infty^*]$ .

Depending on the value of  $A$ , two cases are considered:

*Case 1)  $A \leq 4$ .* According to (24), it suffices to show that for each  $\rho$  such that  $0 \leq \rho \leq \frac{2}{9}$ ,  $f_C(\rho) > 0$ . From Lemma 1 and (7) it follows that  $g_{C-1}(\rho) \leq g_{C-1}(\frac{2}{9}) \leq \frac{1}{3}$ . Consequently, in the interval considered it holds that  $f_C(\rho) = 1 - \rho - g_{C-1}(\rho) \geq 1 - \frac{2}{9} - \frac{1}{3} = \frac{4}{9} > 0$ .

*Case 2)  $A \geq 4$ .* According to (24), it suffices to show that for each  $\rho$  such that  $0 \leq \rho \leq \rho_\infty^* = \hat{\rho}_A$ ,  $f_C(\rho) > 0$ . From (10), (26) and (27), it follows that in the interval considered it holds that  $f_C(\rho) = 1 - \rho - g_{C-1}(\rho) \geq 1 - \rho - g_0(\rho) = 1 - \rho - \frac{A-1}{2}\rho = 1 - \frac{A+1}{2}\rho \geq 1 - \frac{A+1}{2}\hat{\rho}_A = \frac{2A}{(A-1)(A+2)} > 0$ .

To show that  $\rho_{C+1}^* < \rho_C^*$ , it suffices to show that for each  $\rho$  such that  $0 \leq \rho \leq \rho_{C+1}^*$ ,  $f_C(\rho) > 0$ . Note also that for  $A \geq 3$ , (41), (11) and (24) imply that  $\rho_{C+1}^* > \hat{\rho}_A$ . Depending on the value of  $A$  and  $\rho$ , two cases are considered:

*Case 1)  $A \in \{1, 2\}$  and  $0 < \rho \leq \rho_{C+1}^*$ , or  $A \geq 3$  and  $\hat{\rho}_A < \rho \leq \rho_{C+1}^*$ .* From Lemmas 2 and 1 it follows that  $g_{C-1}(\rho) < g_C(\rho) \leq g_C(\rho_{C+1}^*)$ . Thus,  $f_C(\rho) = 1 - \rho - g_{C-1}(\rho) > 1 - \rho_{C+1}^* - g_C(\rho_{C+1}^*) = f_{C+1}(\rho_{C+1}^*) = 0$ .

*Case 2)  $A \geq 3$  and  $0 < \rho \leq \hat{\rho}_A$ .* From Lemma 2 it follows that  $g_{C-1}(\rho) \leq g_0(\rho)$ , and using (26) and (27) we get  $f_C(\rho) = 1 - \rho - g_{C-1}(\rho) \geq 1 - \rho - g_0(\rho) = 1 - \frac{A+1}{2}\rho \geq 1 - \frac{A+1}{2}\hat{\rho}_A = \frac{2A}{(A-1)(A+2)} > 0$ . ■

**Proof of Theorem 2.**

It suffices to show that for any arbitrarily small  $\epsilon (> 0)$ , there exists  $k^*(\epsilon)$  such that the function  $f_{k^*}(\rho)$  has a root in the interval  $(\rho_\infty^*, \rho_\infty^* + \epsilon]$ . Given (41), it suffices to show that  $f_{k^*}(\rho_\epsilon) \leq 0$ , where  $\rho_\epsilon = \rho_\infty^* + \epsilon$ . Depending on the value of  $A$ , two cases are considered:

**Case 1)  $A \leq 4$ .** In this case,  $\rho_\infty^* = \frac{2}{9}$ , such that for  $\rho = \rho_\epsilon = \frac{2}{9} + \epsilon$ ,  $z_{1,2}$  are complex conjugates given by (34). Substituting (34) into (7) and (9), (21) yields after some manipulations

$$f_c(\rho) = V(A, \rho) + \frac{S(A, \rho)}{G(A, \rho) + H(A, \rho) \tan((c-1)\theta)}, \quad \text{for } \rho > \frac{2}{9}, \quad (42)$$

where

$$\left\{ \begin{array}{l} V(A, \rho) = \frac{[2-(A+5)\rho](2-\rho)}{2[2-(2A+1)\rho]}, \quad S(A, \rho) = \rho(2+3\rho)[(A+2)(A-1)\rho - 2(A-2)] \tan(\theta), \\ H(A, \rho) = [2 - (2A+1)\rho]^2, \quad G(A, \rho) = [2 - (2A+1)\rho](2+3\rho) \tan(\theta). \end{array} \right\} \quad (43)$$

It turns out that  $V(A, \rho_\epsilon) > 0$  because  $\rho_\epsilon - 2 < 0$ , and for  $A < 4$  it holds that  $(A + 5)\rho_\epsilon - 2 < 0$ ,  $2 - (2A + 1)\rho > 0$ , and for  $A = 4$  it holds that  $(A + 5)\rho_\epsilon - 2 > 0$ ,  $2 - (2A + 1)\rho < 0$ . Furthermore,  $H(A, \rho_\epsilon) > 0$ , and also  $S(A, \rho_\epsilon) > 0$  because  $(A + 2)(A - 1)\rho_\epsilon - 2(A - 2) > 0$ .

From (42) it now follows that

$$f_c(\rho_\epsilon) \leq 0 \Leftrightarrow -\frac{G(A, \rho_\epsilon)}{H(A, \rho_\epsilon)} - \frac{S(A, \rho_\epsilon)}{V(A, \rho_\epsilon)H(A, \rho_\epsilon)} \leq \tan((c - 1)\theta_\epsilon) < -\frac{G(A, \rho_\epsilon)}{H(A, \rho_\epsilon)}. \quad (44)$$

Let us now define the angles  $\phi_\epsilon$  and  $\omega_\epsilon$  in the interval  $[0, \pi)$  as follows

$$\tan(\phi_\epsilon) = -\frac{G(A, \rho_\epsilon)}{H(A, \rho_\epsilon)} - \frac{S(A, \rho_\epsilon)}{V(A, \rho_\epsilon)H(A, \rho_\epsilon)}, \quad \tan(\omega_\epsilon) = -\frac{G(A, \rho_\epsilon)}{H(A, \rho_\epsilon)}, \quad 0 \leq \phi_\epsilon < \omega_\epsilon < \pi. \quad (45)$$

From (42), it follows for  $A < 4$  that  $-\frac{G(A, \rho_\epsilon)}{H(A, \rho_\epsilon)} = -\frac{2+3\rho_\epsilon}{2-(2A+1)\rho_\epsilon} \tan(\theta_\epsilon) < -\tan(\theta_\epsilon)$ , which in turn implies that  $\theta_\epsilon < \frac{\pi}{2} < \phi_\epsilon < \omega_\epsilon < \pi - \theta_\epsilon$ . For  $A = 4$ , it follows that  $-\frac{G(A, \rho_\epsilon)}{H(A, \rho_\epsilon)} - \frac{S(A, \rho_\epsilon)}{V(A, \rho_\epsilon)H(A, \rho_\epsilon)} = \frac{(2+3\rho_\epsilon)(2-5\rho_\epsilon)}{(9\rho_\epsilon-2)(2-\rho_\epsilon)} \tan(\theta_\epsilon) > \tan(\theta_\epsilon)$ , which in turn implies that  $\theta_\epsilon < \phi_\epsilon < \omega_\epsilon < \frac{\pi}{2}$ . Combining (44) and (45) yields

$$f_c(\rho_\epsilon) \leq 0 \Leftrightarrow \phi_\epsilon \leq (c - 1)\theta_\epsilon < \omega_\epsilon \Leftrightarrow k^*(\epsilon) = \left\lceil \frac{\phi_\epsilon}{\theta_\epsilon} \right\rceil + 1 \leq c \leq \left\lceil \frac{\omega_\epsilon}{\theta_\epsilon} \right\rceil. \quad (46)$$

Thus  $f_{k^*}(\rho_\epsilon) \leq 0$  provided  $k^* \leq \lceil \frac{\omega_\epsilon}{\theta_\epsilon} \rceil$ . Next we shall show that the latter inequality holds. From (45) it follows that

$$\frac{\tan(\omega_\epsilon - \phi_\epsilon)}{\tan(\theta_\epsilon)} = \frac{\tan(\omega_\epsilon) - \tan(\phi_\epsilon)}{(1 + \tan(\omega_\epsilon)\tan(\phi_\epsilon))\tan(\theta_\epsilon)} = \frac{\frac{S(A, \rho_\epsilon)}{V(A, \rho_\epsilon)H(A, \rho_\epsilon)} \frac{1}{\tan(\theta_\epsilon)}}{1 + \frac{G(A, \rho_\epsilon)}{H(A, \rho_\epsilon)} \left( \frac{G(A, \rho_\epsilon)}{H(A, \rho_\epsilon)} + \frac{S(A, \rho_\epsilon)}{V(A, \rho_\epsilon)H(A, \rho_\epsilon)} \right)}. \quad (47)$$

Substituting (34) and (42) into (47), after some manipulations yields  $\frac{\tan(\omega_\epsilon - \phi_\epsilon)}{\tan(\theta_\epsilon)} = \frac{2+3\rho_\epsilon}{2-\rho_\epsilon} > \frac{3}{2}$ . Consequently,  $\omega_\epsilon - \phi_\epsilon > \theta_\epsilon$ , or  $\frac{\phi_\epsilon}{\theta_\epsilon} + 1 < \frac{\omega_\epsilon}{\theta_\epsilon}$ , and therefore  $k^* = \lceil \frac{\phi_\epsilon}{\theta_\epsilon} \rceil + 1 = \lceil \frac{\phi_\epsilon}{\theta_\epsilon} + 1 \rceil \leq \lceil \frac{\omega_\epsilon}{\theta_\epsilon} \rceil$ .

**Case 2)**  $A > 4$ . From (24) it follows that  $\frac{2}{2A+1} < \frac{2}{A+5} < \rho_\infty^* = \hat{\rho}_A < \min(\frac{2}{9}, \frac{2}{A+1})$ . So let us consider  $\rho_\epsilon \in (\rho_\infty^*, \min(\frac{2}{9}, \frac{2}{A+1}))$ . For  $\rho = \rho_\epsilon$ , (8) implies that  $z_{1,2}$  are real numbers, with  $0 < z_2 < z_1$ .

Substituting (7) and (9) into (21) yields

$$f_c(\rho) = \frac{h_1 z_1^{c-1} + h_2 z_2^{c-1}}{b_1 z_1^{c-1} + b_2 z_2^{c-1}}, \quad \text{for } \rho < \frac{2}{9}, \quad (48)$$

$$\text{where } h_i \triangleq (1 - \rho) b_i - a_i = \frac{[2 - (A + 1)\rho](z_i - \xi(A, \rho))}{2(z_i - z_{3-i})}, \quad i = 1, 2, \quad (49)$$

$$\text{with } \xi(A, \rho) \triangleq \frac{4 - (A + 2)\rho}{4[2 - (A + 1)\rho]}. \quad (50)$$

From (33), (50) and (11), it follows that  $\Phi(\xi(A, \rho_\epsilon)) = \frac{\rho_\epsilon(A-1)(A+2)(2-\rho_\epsilon)(\rho_\epsilon-\hat{\rho}_A)}{4[2-(A+1)\rho_\epsilon]^2} > 0$ . Also, using (50) and (8) gives  $\xi(A, \rho_\epsilon) - \frac{z_1+z_2}{2} = \frac{(2-\rho_\epsilon)[(A+5)\rho_\epsilon-2]}{16\rho_\epsilon[2-(A+1)\rho_\epsilon]} > 0$ . Consequently,  $0 < z_2 < z_1 < \xi(A, \rho_\epsilon)$ , which by virtue of (49) implies that  $0 < -h_1 < h_2$ . Similarly, from (33) and (11) it follows that  $\Phi(\frac{A+2}{8}) = \frac{(A-1)(A+2)(\rho_\epsilon-\hat{\rho}_A)}{8} >$

0. Also, using (50) and (8) gives  $\frac{A+2}{8} - \frac{z_1+z_2}{2} = \frac{(2A+1)\rho_\epsilon-2}{16\rho_\epsilon} > 0$ . Consequently,  $0 < z_2 < z_1 < \frac{A+2}{8}$ , which by virtue of (9) implies that  $0 < -b_1 < b_2$ . From the above it now follows that

$$h_1 z_1^{c-1} + h_2 z_2^{c-1} \stackrel{\leq}{=} 0 \Leftrightarrow c \stackrel{\geq}{=} L_h + 1, \quad \text{and} \quad b_1 z_1^{c-1} + b_2 z_2^{c-1} \stackrel{\leq}{=} 0 \Leftrightarrow c \stackrel{\geq}{=} L_b + 1, \quad (51)$$

$$\text{where} \quad L_h \triangleq \frac{\log\left(-\frac{h_2}{h_1}\right)}{\log\left(\frac{z_1}{z_2}\right)}, \quad \text{and} \quad L_b \triangleq \frac{\log\left(-\frac{b_2}{b_1}\right)}{\log\left(\frac{z_1}{z_2}\right)}. \quad (52)$$

By making use of (8), (9), (24) and (50), it can be shown that in the interval considered it holds that  $\frac{z_1}{z_2}\left(-\frac{h_2}{h_1}\right) < -\left(\frac{b_2}{b_1}\right)$  which through (52) implies that

$$L_h + 1 < L_b, \quad \text{and in turn} \quad \lceil L_h \rceil + 1 \leq \lceil L_b \rceil. \quad (53)$$

From (48), (51) and (53) it now follows that

$$-\infty < f_c(\rho_\epsilon) \leq 0 \Leftrightarrow L_h + 1 \leq c < L_b + 1 \Leftrightarrow k^* \triangleq \lceil L_h \rceil + 1 \leq c \leq \lceil L_b \rceil, \quad (54)$$

such that  $f_{k^*}(\rho_\epsilon) \leq 0$ . ■

**Remark 6.** An alternative proof of the theorem can be derived based on the sequence  $\{u_c\}$ . From the definition of  $\mathcal{R}_c$  it follows that  $g_{C-1}(u_{C-1}) = 1$ , such that  $f_C(u_{C-1}) = 1 - u_{C-1} - g_{C-1}(u_{C-1}) = -u_{C-1} < 0$  for  $C \geq 2$ . Therefore, owing to Remark 3, it holds that  $\rho_C^* < u_{C-1}$ ; using (23) and (24) it holds that  $u_\infty < \rho_C^* < u_{C-1}$ . By taking the limit as  $C$  increases and by virtue of (16), we obtain  $u_\infty \leq \lim_{C \rightarrow \infty} \rho_{A,C}^* \leq \lim_{C \rightarrow \infty} u_{C-1} = u_\infty$ , or  $\lim_{C \rightarrow \infty} \rho_{A,C}^* = u_\infty = \rho_\infty^*$ .

#### Proof of Theorem 4.

From Table I it follows that the theorem holds for  $k = 1, 2$ , i.e.  $N = 2, 4$ . We proceed by considering  $k \geq 3$ . For any  $C$  such that  $2 \leq C \leq k - 1$ , it holds that  $A = \frac{N}{2^{C-1}} = 2^{k-C+1} \geq 2^2 = 4$ . Furthermore, for  $A \geq 4$  it also holds that  $\frac{2}{2A+1} \leq \rho_\infty^* = \hat{\rho}_A$ , which by means of (23) and (24) implies that  $\rho_{2A,C-1}^* < \frac{2}{2A+1} \leq \rho_\infty^* = \hat{\rho}_A < \rho_{A,C}^*$ . Furthermore, from Remark 1 and Lemma 5 it follows that  $\rho_{4,k-1}^* < \rho_{2,k}^* = \rho_{1,k+1}^*$ . ■

## REFERENCES

- [1] Q. Yang and K. Bergman, "WDM routing in photonic packet switch," in *Proc. LEOS 2000*, vol. 1 MC2, Rio Grande, Puerto Rico, Nov. 13–16 2000, pp. 31–32.
- [2] B. Small, J. Kutz, W. Lu, and K. Bergman, "Characterizing and simulating the performance of the physical layer of Data Vortex switching nodes," in *Proc. LEOS 2003*, vol. 1 MF5, Tucson, AZ, Oct. 26–30 2003, pp. 59–60.
- [3] Q. Yang, K. Bergman, G. Hughes, and F. Johnson, "WDM packet routing for high-capacity data networks," *IEEE/OSA J. Lightwave Technol.*, vol. 19, no. 10, pp. 1420–1426, Oct. 2001.
- [4] Q. Yang and K. Bergman, "Performances of the Data Vortex switch architecture under non-uniform and bursty traffic," *IEEE/OSA J. Lightwave Technol.*, vol. 20, no. 8, pp. 1242–1247, Aug. 2002.
- [5] Q. Yang, "Optical packet switching for high performance computing," Ph.D. dissertation, Princeton University, Jan. 2002.
- [6] H. Takagi, *Queueing Analysis, Vol. 3: Discrete-Time Systems*. North Holland, Amsterdam: Elsevier Science Publishers, 1993.
- [7] OMNeT++ Discrete Event Simulation System. [Online]. Available: <http://www.omnetpp.org/>
- [8] K. Pawlikowski, H.-D. Jeong, and J.-S. Lee, "On credibility of simulation studies of telecommunication networks," *IEEE Commun. Mag.*, vol. 40, no. 1, pp. 132–139, Jan. 2002.
- [9] R. Rojas-Cessa, E. Oki, and H. Chao, "CIXOB-k: Combined input-crosspoint-output buffered switch," in *Proc. IEEE GLOBECOM 2001*, vol. 4, San Antonio, TX, Nov. 2001, pp. 2654–2660.