

# On the Capacity of Bufferless Networks-on-Chip

Alexander Shpiner, Erez Kantor, Pu Li, Israel Cidon, and Isaac Keslassy, *Senior Member, IEEE*

**Abstract**—Networks-on-Chip (NoCs) form an emerging paradigm for communications within chips. In particular, bufferless NoCs require significantly less area and power consumption, but also pose novel major scheduling problems to achieve full capacity. In this paper, we provide first insights on the capacity of bufferless NoCs. In particular, we present *optimal* periodic schedules for several bufferless NoCs with a complete-exchange traffic pattern. These schedules particularly fit distributed-programming models and network congestion-control mechanisms. In addition, for general traffic patterns, we also introduce efficient greedy scheduling algorithms, that often outperform simple greedy online algorithms and cannot have deadlocks. Finally, using network simulations, we quantify the speedup of our suggested algorithms, and show how they improve throughput by up to 35% on a torus network.

**Index Terms**—Networks-on-Chip, bufferless network, scheduling, collective communication, all-to-all personalized exchange, complete-exchange, interprocessor communication.

## 1 INTRODUCTION

### 1.1 Background

*Networks-on-Chip (NoCs)* form an emerging paradigm for communications within large VLSI systems implemented on a single silicon chip. In a NoC system, modules such as processor cores, memories and specialized blocks exchange data using a network, rather than simple shared busses as in previous systems. A NoC is constructed from multiple point-to-point data links interconnected by routers, such that messages can be relayed from any source module to any destination module over several links, by making routing decisions at the switches.

Interestingly, despite the revolution that the NoC paradigm is causing in computer architecture, *little is known on the capacity region of NoCs*. This is especially surprising because any slight improvement in the capacity of NoCs may have a huge impact. For instance, NoCs are present in most personal computers currently sold around the world [1].

There are many possible NoC topologies. They include simple *line* and *ring* topologies [2], which are widely used in optics-based networks [3]–[5]. For instance, the Intel Sandy Bridge CPU [1], and the IBM Cell Broadband Engine [6] use a ring-based interconnect for their on-chip-network. The *mesh* and the *torus* NoC topologies are also popular [7], [8]. In this paper, we study these four main topologies: line, ring, mesh and torus.

We are especially interested in NoCs that (a) use *bufferless* switches, and (b) carry a *periodic* traffic. Bufferless NoCs are NoCs that eliminate router buffers to reduce NoC cost [9].

- Alexander Shpiner, Israel Cidon and Isaac Keslassy are with the Department of Electrical Engineering, Technion - Israel Institute of Technology. E-mail: shalex@tx.technion.ac.il, {cidon, isaac}@ee.technion.ac.il
- Erez Kantor is with CSAIL, MIT, Cambridge, MA. This work was done when Erez Kantor was a post-doc fellow at the Technion. E-mail: erezk@csail.mit.edu
- Pu Li is with ASML Netherlands B.V. This work was done when Pu Li was a research visitor at the Technion. E-mail: puli@tx.technion.ac.il

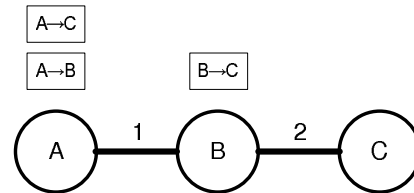


Fig. 1. Simple case with a three-node bufferless NoC sub-network.

As shown in [9]–[14], *bufferless* NoCs and NoCs with small buffers are particularly interesting because they are claimed to offer a lower area and power consumption than those with buffers, in exchange for an increased scheduling complexity and potentially reduced performance. In fact, NoC buffers can consume significant dynamic and static energy, and occupy a large chip area.

Second, we look at applications with a *periodic* traffic. Time-critical periodic applications for embedded devices combine ever-growing computing demands with hard-deadline performance-guarantee requirements. These applications, such as distributed-programming models and network congestion control mechanisms, use a periodic pre-determined traffic pattern with tight deadlines. But the embedded devices often cannot provide the required hard-deadline performance-guarantees. This is because they often rely on NoCs with best-effort communication, which results in an unpredictable network behavior, causing a significant application performance variability.

*Our goal is to devise a periodic schedule algorithm that can provide a capacity-optimal guaranteed service for this traffic pattern in bufferless NoCs.* We also later compare such scheduling algorithms that use the pre-determined nature of the periodic schedule with greedy online algorithms that would not rely on this assumption. We introduce a novel theoretical model of bufferless NoC architectures, and find periodic conflict-free schedules. This is illustrated in the next simple example.

TABLE 1  
Example of non-optimal schedule

time slot	link 1	link 2
1	$A \rightarrow B$	$B \rightarrow C$
2	$A \rightarrow C$	
3		$A \rightarrow C$

TABLE 2  
Example of optimal schedule

time slot	link 1	link 2
1	$A \rightarrow C$	$B \rightarrow C$
2	$A \rightarrow B$	$A \rightarrow C$

**Example 1.** Consider a simple case presented in Figure 1, where the bufferless NoC sub-network consists of three nodes  $A$ ,  $B$  and  $C$  and two links 1 and 2 connecting them. The traffic requirement consists of three packets that need to be sent at each period :  $A \rightarrow B$ ,  $B \rightarrow C$  and  $A \rightarrow C$ . The periodic schedule is created by setting for each packet the time slot at which it is sent in each period. Assume equal propagation times on each link, such that each time slot is sufficient to send a packet on a link. Assume first that the transmissions are scheduled online as shown in Table 1. In this example,  $A \rightarrow B$  is scheduled in the first time slot. Next, packet  $B \rightarrow C$  is also scheduled in the first time slot, because link 2 is free. Finally,  $A \rightarrow C$  is scheduled in the second time slot. It takes two slots to complete the transmission of  $A \rightarrow C$ . Therefore, the length of the total schedule is 3 time slots. As shown in Table 2, the above schedule can be optimized. In the first time slot packet  $A \rightarrow C$  is scheduled on link 1, and in the second time slot on link 2. Link 2 is free in the first time slot, therefore, packet  $B \rightarrow C$  can be scheduled. Finally, packet  $A \rightarrow B$  is scheduled in the second time slot on link 1. The links are utilized all the time, therefore, this schedule is optimal and reaches its full capacity. Its length is 2 time slots.

Note that if the network were buffered, then of course, there would be no need for collision-free scheduling. The collided packets could then be queued at the buffers. However, the buffers would lead to significant power and area costs.

Our suggested periodic transmission of scheduled packets over the NoC can be complemented by best-effort packet services at a strictly lower priority. Therefore, our scheduling algorithms can be used to provide a two-tier solution for both guaranteed-service and best-effort communications in a NoC.

## 1.2 Applications

Many applications and parallel computation models rely on a predetermined periodic traffic. We focus in this paper on a periodic *uniform complete-exchange* communication pattern and on its optimal schedule in different topologies. This uniform complete-exchange pattern, also known as all-to-all personalized communication, is a type of collective communication pattern for Message Passing Interface (MPI) [15], in which all processors need to communicate with all other processors. In this communication pattern, each of the  $N$  processors in the network has a distinct, but equal-size, message to send to each

of the remaining  $N - 1$  processors. It is thus a highly dense communication pattern that can result in many link contentions [16]–[21]. This simple communication pattern enables many numerical parallel algorithms. For example, it includes an increasingly large class of parallel algorithms, such as neural networks, large FFT (Fast Fourier Transform) computations, parallel quick-sort, matrix transpose, array redistribution, distributed table lookup and bitonic sort [22]. More generally, any multi-round algorithm where updates are sent to many randomly-chosen processing elements can be enabled using a topology that allows uniform communications with updates sent to all other processing elements.

It is interesting to note that numerous applications have been shown to allow restructuring that completely separates *the computation logic* from *the communication module*. Such structure complies to a well-established programming model called *Bulk Synchronous Parallel (BSP)* [23]. In BSP, programs are represented as a series of two *super-steps*: first, a computation super-step, then, a communication super-step. BSP has gained popularity for its ability to enable predictable performance on a variety of parallel and distributed platforms, from super-computers to CMPs. Compilers can apply profile-based or static code analysis to determine the pattern in the communication super-step of a BSP application, precompute the communication schedule for a given NoC topology, and augment the compiled code with this information to allow its use at runtime.

The uniform complete-exchange communication pattern could also fit many *end-to-end congestion control* mechanisms in which all nodes exchange information at regular intervals. Such mechanisms can include the estimation of network delays and/or losses in the best-effort network; the periodic acknowledgment of received packets; a generalized hot-spot rate and fairness management; as well as a source-destination queue management [24].

## 1.3 Contributions

The main contribution of this paper is the introduction of *capacity-optimal periodic schedules for uniform traffic over bufferless NoCs, under several topologies*.

The optimality of the schedule is on the capacity utilization of the links, or in other words, the length of the period. We also leverage the fact that the traffic pattern is fixed to compute the schedule only once offline (for instance just after the application compilation). Therefore, by using an optimal offline schedule, we aim to provide a higher capacity utilization, or a shorter period, and hence transmit the same traffic demands with higher throughput using the same bufferless NoC architecture. Using this optimization framework, we prove the existence of several optimal scheduling algorithms on different NoC topologies.

First, we present an algorithm, named Algorithm *DTNS* (Degree-Two NoC Scheduling), for complete-exchange communication in degree-two networks, e.g., line and ring NoC topologies. We prove its optimality and also provide several results on its period length.

Second, we present an algorithm, named Algorithm *TNS* (Torus NoC Scheduling), for complete-exchange in  $N \times N$

torus NoC topologies, and prove its optimality. We later provide lower and upper bounds on the performance of any minimal schedule in the mesh NoC topology. We also provide a constant bound on the ratio between the performance of the TNS algorithm in a mesh and its optimal performance in a torus.

In addition, we also want to provide a more practical algorithm that can deal with general traffic patterns. To do so, we introduce a greedy latency-based scheduling algorithm. In this case the input to our scheduling mechanism is the set of flows with their bandwidths and pre-computed routes [25]–[27]. Then our algorithm allocates the time-slots to the flows in order to minimize the length of the schedule period. This schedule provides both guaranteed throughput and guaranteed latency while using bufferless routers. Note that it could also be used in congestion-free NoCs like optical NoCs; this is, however, beyond the scope of the paper. Finally, using simulations, we show that the proposed schedules outperform approximations of previously-known online algorithms, typically by over 20% in the ring and torus topologies, and 5-15% in the mesh topologies, depending on the number of nodes.

## 2 RELATED WORK

The periodic traffic pattern enables us to rely on a *bufferless NoC*. Providing guaranteed service using a bufferless NoC requires finding a periodic conflict-free schedule. Note that not all bufferless NoC designs meet the required real-time guarantee. For example, bufferless NoCs with deflection routing [9], [10], [28] or dropping [29], [30] are often characterized by a highly-unpredictable network behavior that does not fit our requirements. In fact, [9] shows that buffered NoCs have a better performance, lower cost and complexity than bufferless NoCs with deflection. However, in this paper, we consider periodic collision-free scheduling that is pre-determined and offline, and therefore, does not waste resources on longer routes as in deflection routing.

The works that are closest to ours have been recently realized independently [31], [32]. They also analyze statically scheduled bufferless NoCs with all-to-all complete exchange communication pattern. They formulate the problem of finding an optimal shortest-length schedule for various topologies using an exhaustive search, and provide analytical bounds for the schedule length [31]. In addition, they suggest a heuristic algorithm for generating such schedules [32] and introduce useful and interesting implementation considerations. However, unlike us, they do not attempt to provide a closed-form optimal algorithm. They also assume a slightly different network model, where each node can only inject and absorb a single packet in a time slot, while in our model a node is allowed to inject and absorb a packet to/from each of the four directions in the same time slot.

A bufferless NoC architecture that provides guaranteed service was introduced in Aethereal [33]–[35]. The Aethereal architecture relies on a greedy resource-reservation algorithm that is designed to adapt to changing traffic patterns. In particular, the UMARS algorithm relies on an offline scheduling by ordering the flows by their bandwidth requirements prior to

scheduling them greedily by shortest and less contented route [25]. Other works suggest algorithms for route optimization with multiple paths and shortest latency by keeping in-order arrivals in offline [26], [36] and online [27], [37] modes. However, all these works do not explicitly attempt to *achieve the network capacity* and do not provide an optimal solution for the period length. These papers can nicely extend our work on non-uniform traffic by providing efficient algorithms for routing and mapping, which we do not consider in this paper.

In addition, there have been other works on providing guaranteed service in buffered NoCs. Most suggested architectures, like Nostrum, have relied on router buffers to provide guaranteed service, for instance by using temporally-disjoint networks [38], [39]. Slot allocation in a TDM network-on-chip was introduced in [40]. However it assumes the existence of virtual circuits in the routers, and uses extensive search for providing scheduling for general traffic by validating allocation of each pair of virtual circuits, thus is not optimal and not scalable. Additional papers also focus on all-to-all collective communication patterns in the NoCs [41]–[43]. However, they all assume wormhole routing and buffered NoC. Finally, it is also possible to provide statistical instead of deterministic guarantees [44].

From the theory perspective, bufferless routing was intensively studied under various of names and models. For example, *Direct Routing* [45] defines the problem such that for a given set of packets with corresponding source and destination, the objective is to schedule the injections times of the packets. In different versions of the problem the specified path can be given as an input to the algorithm or defined as another output. The solution requirement is to avoid collisions, and to minimize the schedule period time. [45] presents a randomized  $O(d^2 \log^2 n)$ -approximation algorithm for  $d$ -dimensional Mesh for general traffic pattern. They show that finding an optimal Direct Routing is NP-hard, and for general networks, they show that this problem cannot be easily approximated. Also, algorithms for the complete-exchange pattern for wormhole routing are presented for the torus [46] and mesh [47], [48] topologies. [49] presents a scheduling scheme for complete-exchange in tree networks. Also, in unbuffered optical networks, [50] considers how to schedule packets optimally given several wavelengths, and [51] presents algorithms for a Clos topology.

Another related problem is *hot-potato-routing* studied in [52]–[54], where for a given set of packets, where each packet consists of a source vertex, a destination vertex and an injection time. However, the packets are deflected from their shortest route, therefore the provided solution is not optimal. In addition, packet scheduling in bufferless linear networks was investigated before in offline [55] and online [56] versions. However, no assumption on uniformity or periodicity of the traffic was considered, therefore they could not provide an optimal solution. The baked potato routing algorithm [57] was also among the first to consider switch scheduling in a bufferless network for periodic traffic. However, it only provided a solution for a spanning tree network.

### 3 PROBLEM DEFINITION

In the paper we investigate a bufferless network-on-chip architecture in which the traffic is produced under some periodic traffic demand pattern. The objective is to find a *periodic schedule*  $\Delta$  with *minimal period length*  $S$  needed to service traffic pattern  $\lambda$ , i.e. to maximize the capacity utilization. The periodic schedule  $\Delta$  has to provide congestion-free routing, so that each packet reaches its destination using a *shortest* route.

Formally, consider a network of  $n$  nodes in line, ring, torus or mesh topology. Each node is composed of a router connected to up to 4 neighbor nodes and to a single local module (traffic producer/consumer unit, which, for example, can be a processor core). *The nodes do not buffer, drop or deflect packets.*

Assume that the traffic is produced under some periodic traffic demand pattern. Denote as *uniform* or *complete exchange* a traffic pattern in which each local module of a node, during each period, injects exactly one packet destined to each local module of other nodes in the network. In other words, the normalized traffic demand pattern is  $\lambda_{i \rightarrow j} = 1, \forall i, j$ , where  $i$  and  $j$  are nodes in the network. Further assume that all link capacities are equal to a normalized unit capacity, equivalent to sending one unit-sized packet per time-slot, and denote by  $C_{tot}$  the total link capacity, i.e. the number of links in the network.

Note that larger packets can be subdivided into the unit-size packets that we consider, and which are also sometimes called flits. Then, multi-flit packets can be sent over several schedule periods. These flits can be queued in the nodes before being injected into the network, but once they enter the network, they are not queued in any router before reaching their destination. In the paper, we will use the generic term of *packets* to denote these unit-sized flits.

Each router *must* forward the packets immediately at the next slot. At each time slot, each router can send a single packet to, and can receive a single packet from, each adjacent router. In addition it can send and receive up to 4 packets from its adjacent module.

Note that since the schedule is periodic and pre-determined, there are no deadlocks and no livelocks. Finally, we say that an algorithm  $A$  is *optimal*, if  $A$  produces a minimal period length schedule for the complete-exchange traffic requirements.

A generally-known parameter of a packet *injection rate*  $i$  is defined as the average number of new packets that the nodes send into the network per time unit. The injection rate can be expressed as function of the schedule period length  $S$  as  $i = \frac{n-1}{S}$  per node, or  $i = \frac{n(n-1)}{S}$  for the whole network, since, as defined in the complete-exchange traffic, at each period of length  $S$ , the  $n$  nodes inject  $n-1$  packets each.

### 4 OPTIMAL ALGORITHM FOR COMPLETE EXCHANGE IN DEGREE-TWO NOCS

A degree-two NoC is a NoC topology in which all nodes have a degree of either one or two. There are two types of degree-two NoCs. We denote by *n-Line* a NoC consisting of  $n$  nodes connected by  $n-1$  bi-directional links in a line topology and

by *n-Ring* a NoC consisting of  $n$  nodes connected by  $n$  bi-directional links in a ring topology.

#### 4.1 DTNS Algorithm Description

We start by designing a collision-free scheduling algorithm for complete-exchange in degree-two networks, called *DTNS* (Degree-Two NoC Scheduling). The DTNS algorithm gives a higher priority to the retransmitted packets over new packets at each node, and thus avoids congestions on several types of NoC topologies, as formulated next.

Define a  $\ell$ -hopped packet as a packet whose source-to-destination distance in the shortest route is  $\ell$  links. The algorithm is built in the following way. Each node  $i$  at each time slot  $t$  operates according to the scheme:

- If in time slot  $t-1$  the node receives a packet for forwarding, then it forwards it at time  $t$ .
- Otherwise, it starts transmitting an  $\ell_{max}$ -hopped packet, where  $\ell_{max}$  is the largest number of hops  $\ell$  of all packets left with the node  $i$  as their source node.

For the ring topology, the algorithm assumes a shortest routing. The n-Ring with even number of nodes  $n$ , has two shortest routing alternatives for the  $\frac{n}{2}$ -hopped packets. The choice between the two alternatives can be arbitrary. However, an optimal schedule can be achieved with overlapping periods, by combining the periods into pairs and scheduling the  $\frac{n}{2}$ -hopped packets of one period in the pair at the same time with the  $\frac{n}{2}$ -hopped packets of the next period, over the two alternative routes.

As an example, note that DTNS for  $n=3$  Line produces the schedule presented in Table 2. We obtain the following properties of the DTNS algorithm on degree-two networks.

**Property 1 (n-Line).** *Given an n-Line of  $n$  nodes, the schedule period length  $S_L(n)$  of the DTNS-based schedule is:*

$$S_L(n) = \begin{cases} \frac{n^2}{4} & \text{if } n \text{ is even;} \\ \frac{n^2-1}{4} & \text{if } n \text{ is odd.} \end{cases} \quad (1)$$

*Proof:* The traffic in the n-Line network can be separated into two distinct groups, one for each direction. Therefore, for simplicity, we consider only a single direction in the analysis. By the rules of a DTNS algorithm the packet from node 1 to node  $n$  ( $1 \rightarrow n$ ) with propagation time  $n-1$  slots is transmitted first. It is easy to see that during its transmission, all the packets  $1 \rightarrow a$  and  $a \rightarrow n$ , where  $a = 2, \dots, n-1$  are transmitted. Next, packet  $2 \rightarrow n-1$  is transmitted in  $n-3$  time slots, during which all the packets  $2 \rightarrow a$  and  $a \rightarrow n-1$  where  $a = 3, \dots, n-2$  are transmitted. Summing over all successive packets times, we directly obtain the following schedule length:

$$S_L(n) = \begin{cases} \sum_{i=0}^{\frac{n}{2}-1} (n-2i-1) = \frac{n^2}{4} & \text{if } n \text{ is even;} \\ \sum_{i=0}^{\frac{n-1}{2}-1} (n-2i-1) = \frac{n^2-1}{4} & \text{if } n \text{ is odd.} \end{cases}$$

□

**Property 2 (n-Ring).** *Given an n-Ring of  $n$  nodes, the schedule period length  $S_R(n)$  of the DTNS-based schedule*

TABLE 3  
Example of optimal schedule for 4-node ring.

Clockwise Direction.				
time slot	link 1 → 2	link 2 → 3	link 3 → 4	link 4 → 1
1	1 → 3	2 → 4	3 → 1	4 → 2
2	4 → 2	1 → 3	2 → 4	3 → 1
3	1 → 2	2 → 3	3 → 4	4 → 1
Counter-clockwise Direction.				
time slot	link 2 → 1	link 3 → 2	link 4 → 3	link 1 → 4
1	2 → 1	3 → 2	4 → 3	1 → 4

is

$$S_R(n) = \begin{cases} \frac{n^2}{8} & \text{if } n \text{ is even;} \\ \frac{(n-1)(n+1)}{8} & \text{if } n \text{ is odd.} \end{cases} \quad (2)$$

*Proof Outline:* The proof is very similar to the proof of Property 1, and consists of summing over all successive packet times. We then obtain:

$$S_R(n) = \begin{cases} \sum_{i=1}^{\frac{n}{2}-1} i + \frac{n}{4} = \frac{n^2}{8} & \text{if } n \text{ is even;} \\ \sum_{i=1}^{\frac{n-1}{2}} i = \frac{(n-1)(n+1)}{8} & \text{if } n \text{ is odd.} \end{cases}$$

□

Note that the above property assumes that consecutive periods might overlap. In other words, the transmissions of a new period can be started before the last period was ended. Otherwise, if overlapping is forbidden, we obtain a worse result for even  $n$ 's. This is because overlapping enables an alternate routing of  $\frac{n}{2}$ -hop packets: one period in clockwise direction, and the next period in counter-clockwise direction. Without overlap, we obtain the following result:

**Property 3** (n-Ring without overlap). *Given an n-Ring of  $n$  nodes in which period overlapping is forbidden, the schedule period length  $S_{R,no}(n)$  of the DTNS-based schedule is*

$$S_{R,no}(n) = \begin{cases} \frac{n(n+2)}{8} & \text{if } n \text{ is even;} \\ \frac{(n-1)(n+1)}{8} & \text{if } n \text{ is odd.} \end{cases} \quad (3)$$

*Proof Outline:* The proof is similar again to the proofs above, and simply consists of summing over all successive packet times. We then obtain:

$$S_{R,no}(n) = \begin{cases} \sum_{i=1}^{\frac{n}{2}} i = \frac{n(n+2)}{8} & \text{if } n \text{ is even;} \\ \sum_{i=1}^{\frac{n-1}{2}} i = \frac{(n-1)(n+1)}{8} & \text{if } n \text{ is odd.} \end{cases}$$

□

For the clarity of the algorithm we provide a scheduling example for the 4-node ring in Table 3. The two last slots in the counter-clockwise direction can be used for transmitting the 2-hop packets of the next period. Therefore, the sum of the schedule lengths of the two periods will be 4.

## 4.2 Optimality of DTNS Algorithm

**Theorem 4.** *The DTNS algorithm is optimal on n-Lines and n-Rings with complete-exchange traffic.*

*Proof:* The proof is based on the fact that the bottleneck links are always utilized in the N-Line, and likewise that all the links are always utilized in the n-Ring (they are all bottleneck

links). As a consequence, no other schedule can be more efficient using a smaller period length  $S$ . We prove that the DTNS schedule length is equal to the number of transmissions on the bottleneck link.

First we prove optimality on the  $n$ -Line. We consider two cases.

Case 1:  $n = 2k + 1$  is odd. Let  $l_{mid} = (k, k + 1)$  be the directed link that connects between node  $k$  and node  $k + 1$ . It is easy to verify that  $l_{mid}$  requires to deliver  $k \cdot (k + 1) = \frac{n^2 - 1}{4}$  packets, which is clearly a lower bound of the period length.

Case 2:  $n = 2k$  is even. The number of packets that  $l_{mid}$  requires to deliver is  $(k/2)^2 = n^2/4$ .

Next, we prove optimality on the  $n$ -Ring. Again, we consider two cases.

Case 1:  $n = 2k + 1$  is odd. It is easy to see that each link transmits  $i$  times  $i$ -hopped packet, when  $i = 1, \dots, k$ . Summing the transmission of all the packets we get the length of the optimal period  $S(n) = \frac{(n-1)(n+1)}{8}$  which is equal to the result in Equations (2) and (3).

Case 2:  $n = 2k$  is even. The  $k$ -hopped packets have two shortest paths. Therefore, in each period  $k$ -hopped packets are transmitted only in one direction, while the links on the another direction are idle. Therefore the length of the period consists of the time for the transmission of  $i$ -hopped packets ( $i = 1, \dots, k - 1$ ), which is equal to  $\sum_{i=1}^{k-1} i$  time slots, and the time for the the transition of the  $k$ -hopped packets, which is equal to  $k$  time slots. Summing above, we get the result in Equation (3). Of course, this result is not optimal, because of the idle links. If overlapping of the scheduling between the two adjacent periods is allowed, then the idle links can used to transmit the  $k$ -hopped packets of the next period, and thus the average number of time slots to schedule the  $k$ -hopped packets is equal to  $\frac{k}{2}$ . Thus, we get the result in Equation (2). □

## 5 OPTIMAL ALGORITHM FOR COMPLETE EXCHANGE IN TORUS NOCS

In this section, we present an algorithm, named *TNS* (*Torus NoC Scheduling*), which is designed to provide a periodic schedule over the Torus NoC topology for the complete-exchange traffic pattern. We will demonstrate that TNS is guaranteed to achieve the network capacity in an  $N \times N$  Torus network with  $n = N^2$  nodes, using the same setting assumptions of uniform traffic and unit capacities.

The TNS algorithm provides the *injection time* to the network and the *minimal-length route* for each packet within the period, based on the source and destination of the packet. It also guarantees that there are no packet collisions. Therefore, *the TNS algorithm does not rely on buffering or dropping packets, and it also has no deadlocks.*

### 5.1 TNS Algorithm Description

Consider an  $N \times N$  torus network with  $n = N^2$  nodes, and total allowed capacity  $C_{tot}$  for all links in the network. Each node is composed of a router connected to 4 neighbor nodes and to a single local module (traffic producer/consumer unit). Denote the nodes as a set of tuples  $\{(x, y) \mid x, y \in \{1, \dots, N\}\}$ ,

where the first entity refers to the rows on the torus and the second entity refers to the columns on the torus. Denote by  $DIST(a, b) = \min\{|a - b|, N - |a - b|\}$ , the *distance* between  $a$  and  $b$ , for every  $a, b \in \{1, \dots, N\}$ .

We now present the TNS algorithm in detail. Intuitively, TNS decomposes the period of length  $S$  into several *phases* of unequal lengths. In each phase, it transmits all the packets between all source-destination pairs that have the same maximum distance across dimensions. TNS further decomposes each phase into several sub-phases called *epochs*. In each epoch, it connects nodes that are not only at the same maximum distance, but that also follow a given pattern. We now define the TNS algorithm more formally.

**Phases** — The schedule period is divided into  $\lfloor N/2 \rfloor$  *phases*. We say that a packet belongs to the *envelope* of a square of nodes  $(i + 1) \times (i + 1)$ , if the maximum distance across all dimensions between its source and its destination is equal to  $i$ . In other words, consider a packet from node  $(a, b)$  to node  $(c, d)$ . We denote  $i = \max\{DIST(a, c), DIST(b, d)\}$ . Then, the packet belongs to the envelope of a square of nodes  $(i + 1) \times (i + 1)$ . In *phase*  $i$ , packets in all the envelopes of squares of size  $(i + 1) \times (i + 1)$  are scheduled to be transmitted. For example, the packets on Figure 2 belongs to an envelope of  $3 \times 3$ .

**Epochs** — Each phase consists of several epochs. The number of epochs per phase depends on the phase number. Specifically, phase  $i$  consists of  $2i$  *epochs* for  $1 \leq i < N/2$ . In addition, the number of epochs in phase  $N/2$  depends on whether  $N$  is odd or even. If  $N$  is odd, phase  $N/2$  also consists of  $i = N$  epochs. However, if  $N$  is even, phase  $N/2$  consists of only  $\frac{N}{2} + 1$  epochs.

On epoch  $j \in \{0, 1, \dots, 2i - 1\}$  of phase  $i$ , each node injects four packets to four different destinations by crossing exactly  $i + j'$  links, where  $j' = j$ , if  $j \leq i$  and otherwise,  $j' = j - i$ . First  $i$  links in one direction and the other  $j'$  links in the perpendicular direction, as explained below. The destinations are given by the following three steps, for each of the four packets, one for each direction:

- 1) First, the packet follows a direct traversal walk on  $i$  links;
- 2) Then, in epoch 0 it stops; in the epoch  $(j|0 < j \leq i)$ , it “turns right” in a clockwise direction; in the epoch  $(j|i < j < 2i)$ , it “turns left” in counter-clockwise direction.
- 3) Finally, there is another direct traversal walk through  $j'$  additional links.

For clarity, *all* the packets that are scheduled to be transmitted during the epoch, are injected in the first time slot of the epoch. Since all the packets in the epoch have an equal distance to traverse between the source and destination, the length of the epoch simply equals the number of time slots that a packet takes to reach the destination in the epoch. Thus, epoch  $j$  of phase  $i$  takes  $i + j'$  time-slots. Finally, after completing all the  $N/2$  phases, all the pairs of nodes in the Torus have been connected and have transmitted a packet exactly once.

The example in Figure 2 shows the transmission during phase  $i = 3$  and epoch  $j = 1$  (a) and  $j = 4$  (b). All

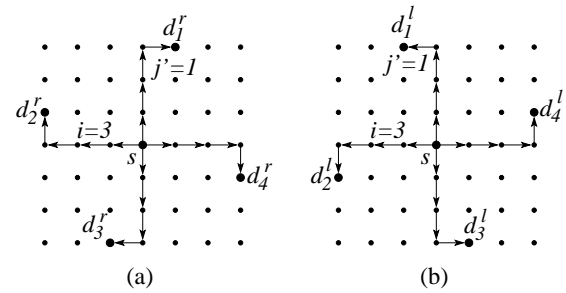


Fig. 2. Two epochs in phase  $i = 3$ . (a) illustrates epoch  $j = 1$ , with the traversal walk and a turn in a clockwise direction. (b) illustrates epoch  $j = 3 + 1 = 4$ , with the traversal walk and a turn in a counter-clockwise direction.

the nodes simultaneously inject the packets in four direction. Every packet is transmitted over  $i = 3$  hops in the direction it was injected. Then, every packet is transmitted  $j' = 1$  more hop left (a) or right (b). The total length of the epoch is  $i + j' = 4$ .

Now that we have formally defined TNS, we show in which epoch each packet is transmitted. First, for a given source node  $(a, b)$  and destination node  $(c, d)$ , denote by  $(a, b) \rightarrow (c, b) \rightarrow (c, d)$  a shorter path that goes from  $(a, b)$  to  $(c, d)$  via  $(c, b)$ . If  $DIST(a, c) \geq DIST(b, d)$ , then we say that this path is a *long-then-short shorter* path. (Note that the TNS algorithm uses only long-then-short shorter paths.) The algorithm sets a unique phase  $i$  and epoch  $j$  for each source-destination pair, denoted by the pair  $(i, j)$ . For instance, consider a packet from source node  $(a, b)$  to destination node  $(c, d)$ . The *phase*  $i$  in which the packet is transmitted is given by:

$$i = \max\{DIST(a, c), DIST(b, d)\}.$$

Moreover, the *epoch* in which the packet is transmitted is either  $j = j^*$  or  $j = j^* + i$ , where

$$j^* = \min\{DIST(a, c), DIST(b, d)\}.$$

The exact epoch is decided by the direction (clockwise or counter-clockwise) of the traversal walk of a long-then-short shorter path between source node  $(a, b)$  and destination node  $(c, d)$ . If there exist two different long-then-short shorter paths (i.e., for the case where  $DIST(a, b) = DIST(c, d)$ , that is  $i = N/2$  for even  $N$ ), then the packet is delivered on epoch  $j = j^*$ .

For example, a packet from  $(3, 4)$  to  $(4, 1)$  is transmitted in phase  $i = \max\{DIST(3, 4), DIST(4, 1)\} = \max\{1, 3\} = 3$ . Then,  $j^* = \min\{1, 3\} = 1$ , and due to the fact that the long-then-short transmission is counter-clockwise, the epoch is  $j = j^* + i = 4$ . The packet will be transmitted in parallel with all other packets destined from  $(3+x, 4+y)$  to  $(4+x, 1+y)$ , with packets from  $(3, 4)$  to  $(6, 5)$ ,  $(2, 7)$ ,  $(0, 3)$  and with all their parallel packets.

## 5.2 Correctness of TNS Algorithm Collision-free Property

In this section we prove the correctness of the collision-free property of the TNS algorithm. Since the propagation

times in all the links are equal, and only a single packet is injected by each node in each direction, there are no collisions during the first step. In the second step, all the packets turn left in epochs ( $j|0 < j \leq i$ ) or all the packets turn right in epochs ( $j|i < j < 2i$ ). Later, in the third step, they all continue straight. Hence the algorithm makes all packets turn in the same way at the precise same time slots, and therefore ensures that there are no collisions during these three steps.

**Theorem 5.** *TNS Algorithm provides collision-free packets scheduling.*

*Proof:* Negatively suppose that two different packets A and B collided, i.e. were transmitted on the same link at the same time slot. The TNS algorithm describes the same traffic pattern at each epoch, at which all the packets pass  $i$  hops in one direction and all the packets turn at the same slot to the same relative direction (left or right). Therefore, packets have the same route in the network. Since, the propagation times in all the links are equal, and the packets that are transmitted at each epoch pass an equal number of hops, packets A and B were injected by the same source to the same direction. That contradicts the rule that each source node injects only a single packet to each direction.  $\square$

### 5.3 Optimality of TNS Algorithm

In this section we prove the optimality of the TNS algorithm. First, we develop expressions for the period length of the TNS schedule in Property 6. Next we prove the algorithm optimality in Theorem 7, by showing that the period length equals the ratio of the traffic requirements by the achievable network capacity.

**Property 6** (TNS schedule period length). *Given an  $n = N \times N$  Torus, the length of the schedule period  $S_T(n = N \times N)$  of the TNS-based schedule is:*

$$\begin{aligned} S_T(n = N \times N) &= \sum_{i=1}^{\lfloor N/2 \rfloor} \left( 2 \sum_{j=i}^{2i} j - 3i \right) \\ &= \frac{N^3 - N}{8} = \frac{n\sqrt{n} - \sqrt{n}}{8}, \end{aligned} \quad (4)$$

for odd  $N$  and

$$S_T(n = N \times N) = \frac{N^3 + 2N}{8} = \frac{n\sqrt{n} + 2\sqrt{n}}{8}$$

for even  $N$ . Furthermore, when overlap is forbidden, the result for even  $N$  becomes

$$S_T(n = N \times N) = \frac{N^3}{8} + N = \frac{n\sqrt{n}}{8} + \sqrt{n}$$

*Proof:* Consider an odd  $N$ . For  $N \times N$  Torus, all connections between all the nodes are covered within  $\lfloor N/2 \rfloor$  phases. It is the maximum distance across all dimensions between the sources and the destinations of all the packets.

The shortest epoch in the phase  $i$  is the one that transmits packets in one dimension, i.e.,  $c = a$  or  $d = b$ . Its length is  $i$  time slots, because packets with latency of  $i$  hops are scheduled within it. The longest epoch in the phase  $i$  is the

one that transmits packets between the corners of the  $(i + 1) \times (i + 1)$  square ( $DIST(a, c) = DIST(b, d)$ ). Its length is  $2i$  time slots. Other epochs are scheduled twice, one for each direction. Therefore, phase  $i$  consists of one epoch for one-dimensional packets, one epoch for corner packets and two epochs for other packets, one in each direction. Thus, the length of phase  $i$  is

$$S_{T,i}(n = N \times N) = 2 \sum_{j=i}^{2i} j - 3i = 3i^2 \quad (5)$$

time slots, for every  $i < N/2$ . For even  $N$ , the last phase (phase  $N/2$ ) has only  $N/2 + 1$  epochs that takes  $\sum_{j=N/2}^N j$  time slots. It is simply counted differently whether or not there is overlap.  $\square$

Note that when overlapping is allowed, it is possible to yield a non-integer period time. For instance, it is possible to pack two periods into three time-slots (for  $N = 2$ ). Actually,  $N = 2$  is an extreme case when a pair of periods is fully overlapped.

**Theorem 7.** *The TNS algorithm is optimal on Torus NoCs topology with complete-exchange traffic.*

*Proof:* Depending on whether  $N$  is odd or even. These two cases are considered separately.

**Case 1:  $N$  is odd:** Denote  $L_{(i,j) \rightarrow (k,l)}$  as the minimum number of hops that a packet passes from node  $(i, j)$  to node  $(k, l)$  ( $i, j, k, l \in \{1, \dots, N\}$ ), and it is given as the sum of the distances on both axes:

$$L_{(i,j) \rightarrow (k,l)} = DIST(i, k) + DIST(j, l). \quad (6)$$

For the uniform traffic, the total number of hops for all packets in a period  $S_t(N)$  is given by:

$$L_t = \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N \sum_{l=1}^N L_{(i,j) \rightarrow (k,l)} = \frac{N^3(N^2 - 1)}{2}. \quad (7)$$

Divide  $L_t$  by  $S_T$  using Equation (4), and get

$$C_{tot} = L_t / S_T = 4N^2, \quad (8)$$

where  $C_{tot}$  is the capacity of the torus, i.e., the total number of available links in the network. It indicates that by this scheduling algorithm, all the links are always utilized during the transmission. Thus, the algorithm is optimal.

**Case 2:  $N$  is even:** Denote *edge packets* as the packets that have two paths with equal distances of  $N/2$  in one of the dimensions, for a shortest routing from source node to the destination node. For example, if the equal distance is in the X axis, the source node is  $(i, k)$ , then the destination node is  $((i + N/2) \pmod{N}, l)$ . We first consider the case when *overlapping* is allowed and calculate the transmission of the *edge packets* in this period twice on account of their transmission in the next period. The duration of transmission of the packets that are not the *edge packets* is the duration of the transmission of phases 1 to  $N/2 - 1$ :

$$\begin{aligned} S_{t, N/2-1} &= \sum_{i=1}^{N/2-1} \left( 2 \sum_{j=i}^{2i} j - 3i \right) \\ &= \frac{N(N-1)(N-2)}{8} \end{aligned} \quad (9)$$

time slots. The sum of latencies of the packets that are not the *edge packets* is:

$$\begin{aligned} L_{t, \text{non-edge}} &= \\ &= 2N(N-1) \sum_i \sum_{k \neq (i+N/2) \bmod N} \text{DIST}(i, k) = \\ &= \frac{N^3(N-1)(N-2)}{2}. \end{aligned} \quad (10)$$

When dividing the sum of latencies  $L_{t, \text{non-edge}}$  by the time slots  $S_{t, N/2-1}$ , we get  $4N^2$ , which is equal to the total capacity of the links  $C_{tot}$ , and indicates that all links are utilized all the time during the transmission. Thus, for the transmission of *non-edge packets*, the scheduling algorithm is optimal. The *edge packets* are transmitted in the phase  $N/2$ . As we seen before, the *edge packets* have two shortest paths. Therefore, they can be transmitted twice during a single *phase*, each on the one of the two shortest paths. Therefore, the duration of transmission of the *edge packets twice* is the duration of phase  $N/2$ :

$$S_{N/2} = \left( 2 \sum_{j=i}^{2i} j - 3i \right) \Big|_{i=N/2} = 2 \sum_{j=i}^{2i} j - \frac{3N}{2} = \frac{3N^2}{4}. \quad (11)$$

By considering all possible cases, we find that the sum of the latencies of the *edge packets* is:

$$\begin{aligned} L_{t, \text{edge}} &= \\ &= 2N^2 \sum_i \sum_{k=(i+N/2) \bmod N} \text{DIST}(i, k) + \\ &+ 2N \sum_i \sum_k \text{DIST}(i, k) = \\ &= 2N^2 \cdot N \cdot N/2 + 2N \left( N \left( \sum_{d=1}^{N/2-1} 2d + N/2 \right) \right) = \\ &= N^4 + N^4/2 = \frac{3N^4}{2}. \end{aligned} \quad (12)$$

Therefore, the following equality holds:

$$C_{tot} = 4N^2 = \frac{2L_{t, \text{edge}}}{S_{N/2}},$$

which means that transmitting the *edge packets* twice during the phase  $N/2$  maximally utilizes the links. One can consider the second transmission of the *edge packets* for the next schedule period. We showed that with uniform traffic and unit link capacities, the links in the Torus are always utilized, hence the algorithm is optimal.  $\square$

## 6 SCHEDULING BOUNDS FOR COMPLETE EXCHANGE IN MESH

We now want to provide some intuition on the schedule period in a NoC Mesh topology. In such a topology, the TNS algorithm is not optimal anymore, and we have not found any characterization of a generally optimal algorithm. Therefore, we will provide instead lower- and upper-bounds on the schedule period, and compare it with the torus schedule period.

### 6.1 Lower Bound for Mesh Schedule Length

We now establish a lower bound on the period length in the Mesh NoC topology.

**Theorem 8.** *Given an  $n = N \times N$  Mesh, the period length  $S_M(n = N \times N)$  of any schedule for complete-exchange traffic satisfies*

$$\begin{cases} S_M(n = N \times N) \geq \frac{N^3}{4} = \frac{n\sqrt{n}}{4} & \text{if } N \text{ is even,} \\ S_M(n = N \times N) \geq \frac{N^3 - N}{4} = \frac{n\sqrt{n} - \sqrt{n}}{4} & \text{if } N \text{ is odd.} \end{cases} \quad (13)$$

*Proof:* As previously, the proof relies on computing the load on the bottleneck links, and using this load for the lower-bound on the period length. In the proof we assume XY routing, however, the result can be generalized for every traffic pattern, as will be shown later. Then, we calculate the load on the link that transmits the largest number of packets (the bottleneck link). The lower bound of the period length is the time which is taken to transmit the packets over the bottleneck link.

According to the properties of the XY routing, each horizontal link  $((i, j) \rightarrow (i+1, j))$  transmits packets from  $j$  nodes  $((i, 1), (i, 2), \dots, (i, j))$  to  $N \cdot (N-j)$  nodes  $(k, l)$ , where  $k = 1, \dots, N$  and  $l = i+1, \dots, N$ .

Thus, the number of packets  $P_{i, j \rightarrow i+1, j}$  that are transmitted on link  $((i, j) \rightarrow (i+1, j))(N)$  in each period is:

$$P_{i, j \rightarrow i+1, j}(N) = j \cdot N \cdot (N-j) \quad (14)$$

Similar calculations can be obtained for the vertical links.

Equation (14) is maximized on links where  $j = \frac{N}{2}$  when  $N$  is even and  $j = \frac{N-1}{2}$  when  $N$  is odd, therefore the most utilized link transmits  $P_{max}(N)$  packets:

$$P_{max}(N) = \begin{cases} \frac{N^3}{4} & \text{if } N \text{ is even;} \\ \frac{N}{4} (N^2 - 1) & \text{if } N \text{ is odd.} \end{cases}$$

Therefore,

$$S_M(n = N \times N) = P_{max}$$

and the result stands.

Initially we assumed XY routing and show that the maximally loaded links are all the central links that connect the left part with the right part of the mesh. Therefore, the traffic load connecting the two parts is equally spread over those central links. We can use this maximal load as the lowest bound for the schedule length. Since the load is equally spread, using any other type of routing cannot decrease the load on the maximally utilized link, and thus cannot decrease the lower bound. Therefore, the final result can be generalized for each type of routing.  $\square$

### 6.2 Upper Bound for Mesh Schedule Length Using TNS

We now want to provide an upper bound on the period length in a mesh. The application of the TNS algorithm in a mesh topology is as follows. Consider an instance of the TNS algorithm on a  $2N \times 2N$  torus. The  $N \times N$  mesh is embedded in a  $2N \times 2N$  torus by an  $N \times N$  subgraph combined from



nodes  $(i, j)$ , where  $N/2 < i, j \leq 3N/2$ . Using the TNS algorithm on the  $2N \times 2N$  torus, the shortest paths between the nodes in this subgraph are routed within the subgraph (contrary to some other shortest paths outside this subgraph that are routed through the boundaries of the torus). For the above  $N \times N$  mesh subgraph we use the TNS scheduling of the full  $2N \times 2N$  torus, transmitting only the packets between the nodes of the  $N \times N$  mesh subgraph, and leaving empty the scheduling slots for the packets outside the subgraph. We can now prove the following result:

**Theorem 9.** *Let  $S_M(n = N \times N)$  denote the minimal schedule length in an  $n = N \times N$  Mesh, and  $S_T(n = N \times N)$  denote the minimal schedule length in an  $n = N \times N$  Torus. Then  $S_M(n = N \times N)$  satisfies:*

$$S_M(n) \leq S_T(4n) \quad (15)$$

*Proof:* From Property 6, we know the schedule length  $S_T(n)$  of TNS in an  $N \times N$  torus. It is easy to see that  $S_T(n = N \times N)$  grows as  $O(N^3)$ , i.e. increases by a factor of 8, for doubling  $N$ . Moreover, we showed that we can schedule an  $N \times N$  mesh using the TNS algorithm applied in a  $2N \times 2N$  torus. The length of this schedule is  $S_T(4n = 2N \times 2N)$ . Therefore, the upper bound of  $S_M(n)$  is  $S_T(4n)$ .  $\square$

The results of Theorems 8 and 9 provide an asymptotic ratio of at most 4 between the upper-bound and the lower-bound, because

$$\begin{aligned} \frac{n\sqrt{n}}{4} &\leq S_M(n = N \times N) \\ &\leq S_T(4n = 2N \times 2N) = \frac{8n\sqrt{n} + 4\sqrt{n}}{8}. \end{aligned}$$

Therefore, we are able to determine the growth rate of the periodic schedule length within a  $4 + \frac{2}{n}$  - constant approximation.

## 7 GREEDY SCHEDULING ALGORITHM FOR GENERAL PERIODIC TRAFFIC PATTERN

The generalized scheduling problem of non-uniform traffic is NP-hard [45]. Therefore, we want to design an efficient heuristic algorithm to schedule an arbitrary traffic pattern in a general NoC topology with a general traffic demand pattern.

**Algorithm Overview:** We introduce the *Latency-based Scheduling Algorithm*. The input of the algorithm is a periodic flow requirement that is given together with a predefined routing. The algorithm outputs a schedule that guarantees periodic service of the traffic without any collisions.

Before defining the algorithm formally, let's first provide some intuition. The algorithm runs offline. It relies on a centralized scheduler that uses given traffic demands and routes to provide a schedule. Its goal is to make the schedule period as small as possible. Intuitively, the packets that are hardest to schedule are those that take the most delay in their transmission, since they occupy more slots and therefore can experience more conflicts. Therefore, to minimize conflicts, in the latency-based algorithm, we *first schedule the packets with the longest latencies*.

Note that in the example of Figure 1, the latency-based scheduling will always produce the optimal schedule that is shown in Table 2, while a random greedy scheduling can produce either of the two schedules. We will show in simulations below that latency-based heuristics also often work better in more complex cases as well.

Also, note that Ref. [32] independently introduced a similar heuristic algorithm for the collision-free scheduling under a slightly different model. However, it is not destined to provide a schedule for any general periodic traffic, but to provide a symmetric schedule for the uniform traffic in symmetric topologies. Since it has a smaller number of scheduling alternatives, it converges faster.

**Algorithm Description:** Let's now formally introduce the latency-based scheduling algorithm. We first consider the given list of traffic demands, and sort all traffic demands by latencies in a decreasing order. Specifically, denote by  $l(d_i)$  the *latency* of a flow demand  $i$ . It is defined as the length of its predefined route between the source and the destination nodes. Assuming for  $n$  flow demands that  $l(d_1) \geq l(d_2) \geq \dots \geq l(d_n)$ , then the ordered set is  $(d_1, d_2, \dots, d_n)$ .

The scheduler considers an empty slot table consisting of links on one axis and time slots on the second axis, for the switching schedule [33] (as illustrated in Tables 1 and 2), and successively attempts to fill it. To do so, the algorithm iterates over the set of unallocated flows, and each time picks a flow demand with the largest latency, which is defined by the number of hops in the flow route. It then attempts to schedule the flow demand by placing it in the earliest available time-slots. Specifically, given the set of flow demands ordered by the latency, from the largest to the smallest,  $(d_1, d_2, \dots, d_n)$ , the scheduler picks demand  $d_1$  and removes it from the set. For all of the node routers on the flow route, it allocates the earliest available time-slots. The injection time of the flow is the allocated time-slot in its first node router on the route (as also described in [25]). When the flow demand set is empty, all the flows are allocated. The period length is set as the latest allocated time-slot on all the routers. For example, the period lengths of the schedules in Tables 1 and 2 are 3 and 2, respectively.

## 8 EVALUATION

### 8.1 Baseline Random-Greedy Algorithm

We define the *random-greedy* algorithm as a baseline in our simulations. Our goal is to approximate the typical algorithms that consider traffic pattern demands in an online fashion, even though we compute our schedule only once at the start in an offline way. To approximate online computations, we consider the demands of the traffic pattern in a random order. To do so, we run a random permutation of all the demands, and then consider them one after the other. In other words, unlike in our latency-based scheduling algorithm, we do not order flows by latency before scheduling them. We skip the ordering by latency step of our latency-based scheduling algorithm. Therefore, the scheduler iterates over a set of unallocated flows, randomly picks one and attempts to schedule it in a greedy manner, by placing it in the earliest available slots.

While the *online* random-greedy algorithm obtains a single schedule, we can obtain a better result in our *offline* setting. To do so, we run the algorithm many times using different random permutations of the demand order, and obtain the schedule with the lowest period length out of all runs. In our simulations graphs we denote the best simulation run by *random-greedy offline*. In the figures 3 we show the distribution of the performance of the different runs of the algorithm.

## 8.2 Period Length with Complete Exchange Traffic

We simulate three NoC architectures: Ring, Torus and Mesh. We assume a complete exchange traffic pattern with a unit-sized packet to send in each period from each source to each destination. The simulations are written in Matlab. Fig. 3 plots the distribution of the period length  $S$  in each of the three NoC topologies, given several algorithms. First, Fig. 3(a) plots the period length distribution in an N-Ring topology with  $N = 16$ . It compares the performances of the *random-greedy* and *latency-based* greedy algorithms, described in Section 7, together with the optimal *TNS* algorithm. We remind that *random-greedy* is meant to approximate the performance of an online algorithm that would service traffic requests in their random order of arrival. On the other hand, the other algorithms attempt to benefit from the fact that the entire traffic demand pattern is known in advance. To obtain an approximate practical distribution, each of the greedy algorithms is run 100 times.

We can see that DTNS manages to schedule all packets within 33 slots; while the random-greedy algorithm needs between 39 and 47 slots with an average of 41.93, and the latency-based algorithm improves upon the random-greedy algorithm by reaching between 35 and 42 slots with an average of 37.94.

Next, Fig. 3(b) plots the period length distribution in an  $N \times N$  Torus topology with  $N = 8$ , i.e.  $N^2 = 64$  nodes. Again, we can see that TNS, the optimal algorithm, performs significantly better than the greedy algorithms, with a meaningful difference with even the best greedy result out of 100 runs.

Last, Fig. 3(c) plots the period length distribution in an  $N \times N$  Mesh topology with  $N = 8$ , i.e.  $N^2 = 64$  nodes. This time, we do not know the optimal algorithm. Therefore, we plot the lower bound achieved in Section 6.1. We see that the distributions of the *random-greedy* and the *latency-based* greedy algorithms are disjoint after 100 runs, indicating that latency-based significantly outperforms random-greedy in this topology and it is relatively close to the lower bound.

## 8.3 Scaling

We now want to study how scaling  $N$  impacts the results. We define *speedup* as the ratio of the average throughput under a specific algorithm to the average throughput under the basic random-greedy algorithm. The higher the speedup, the better the throughput gain. Fig. 4 shows the speedup gained by several algorithms described versus the *average* performance of the random-greedy algorithm achieved in our simulations. The three sub-figures illustrate the three different NoC topologies: the Ring with  $N$  nodes, the  $N \times N$  Torus,

and the  $N \times N$  Mesh. In all the sub-figures,  $N$  varies from 3 to 10. For each  $N$  we run 100 iterations of *random greedy* and *latency-based* algorithms.

In all sub-figures, we show the best speedup result for the random-greedy algorithm (when compared to its average result). We also only show the best result for the latency-based algorithm, since it is run in offline mode, and the developer can run several iterations in order to use the best result. For the Ring and Torus we also compare with the optimal TNS algorithm. In all topologies, we can see that the variances tend to decrease as  $N$  is scaled. We can also notice successive improvements when going *from online to offline mode*, by considering the speedup of the best random-greedy run vs. its average run. Next, we can see a further improvement given the best latency-based run compared to the best random-greedy run. Finally, the optimal algorithm clearly outperforms the greedy algorithms in most cases.

## 8.4 Non-uniform Traffic

Next, we compare the algorithms given non-uniform traffic. Two synthetic patterns are used: *uniform-random* [58] and *bit-permutation* [59]. The *uniform-random* pattern is created by uniformly randomly choosing one destination for each source node. The *bit-permutation* pattern is a pattern where each node sends packet to exactly one other node, and receives packet from exactly one node. Note that even though the choice of destination node in the traffic patterns is uniform, both traffic patterns are non-uniform, since after the traffic patterns are determined, each node is sending packets to a single specific destination only.

Figures 5 and 6 show the distribution of the simulation results with the random-greedy and latency-based algorithms on the  $8 \times 8$  Torus and  $8 \times 8$  Mesh topologies for the *uniform-random* and *bit-permutation* traffic patterns, respectively. In all cases the latency-based algorithm achieves a slightly better distribution.

## 8.5 Analytical Evaluation

We now want to study the tradeoff between the additional capacity needed in the Torus NoC topology when compared to the Mesh NoC topology, and the resulting additional performance. Figure 7 compares the performance ratio and the capacity ratio of the Torus and the Mesh NoC topologies.

- The *performance ratio* is defined as the ratio between the average period length in the random-greedy algorithm given the Torus topology, by the same parameter given the Mesh topology.
- The *capacity ratio* is defined as the ratio between the link capacities of the Torus and those of the Mesh, i.e.  $\frac{\text{Torus capacity}}{\text{Mesh capacity}} = \frac{N^2 + N}{N^2}$ .
- The *link cost ratio* is defined as the ratio between the lengths of the links in the Torus and those of the Mesh. The link length is directly correlated with the chip area required for a link. Assuming the simplified drawn topology, the link cost ratio is equal to 2 for every  $N \geq 2$ .

We can see that as we scale the topologies, the performance ratio seems to stabilize around 1.68, while the capacity ratio

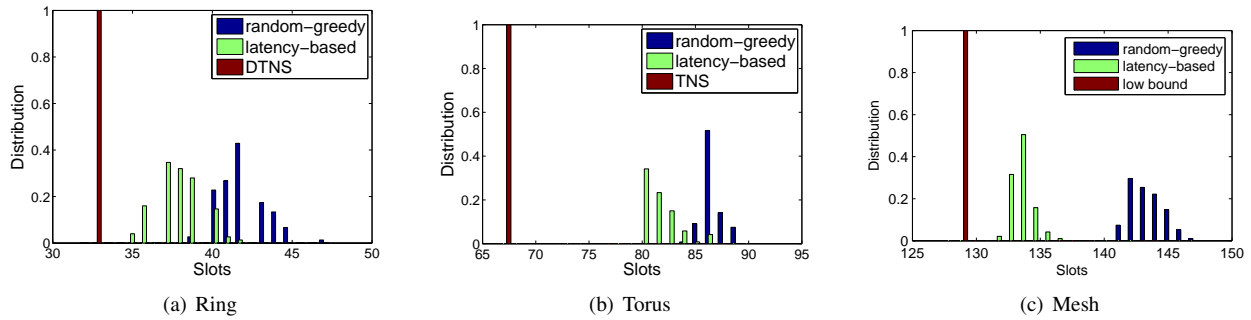


Fig. 3. Period length distribution for different topologies (smaller is better). The latency-based algorithm clearly outperforms the baseline random-greedy algorithm in all topologies. In addition, none of these algorithms can yield the same performance as our optimal DTNS and TNS algorithms in Ring and Torus.

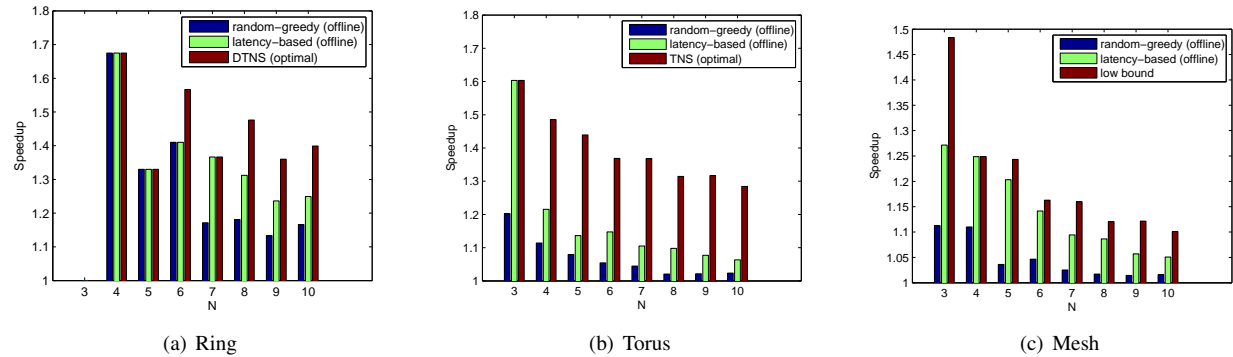


Fig. 4. Speedup for different topologies when compared to the *average* random-greedy online performance. The offline random-greedy algorithm, defined as its best run out of 100, improves upon the average online performance. The offline latency-based algorithm further improves the performance. The optimal algorithms in Ring and Torus provide a speedup above 20% in all cases.

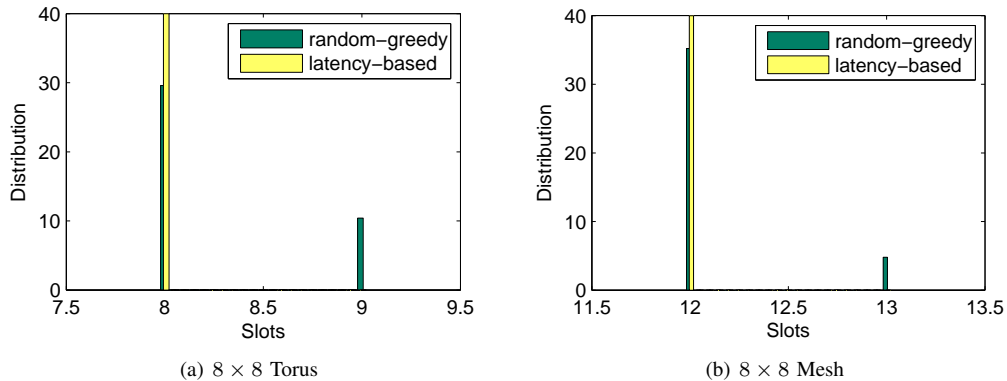


Fig. 5. Period Length Distribution with the *Uniform-Random* traffic pattern. In all cases the latency-based algorithm achieves a slightly better distribution.

converges to 1. At first look, this might indicate that the torus topology is more adapted. However, bear in mind that torus links might need to be longer, and therefore their cost might be higher when compared to mesh links than indicated in this simplistic model. On a more theoretical level, notice that the ratio is indeed always between 1 and 8, as proved in Theorem 9.

Next, we compare the length of the optimal schedule of DTNS and TNS algorithms under the line, ring and torus

topologies as a function of the number of nodes (Figure 8) and a function of the total link capacities (Figure 9). The graphs were obtained using Properties 1 and 2 and Property 6. Note that the impact of potential overlapping periods is negligible with a large number of nodes (and links) and, therefore, we did not refer to it in the plots. The results show the gained speedup of torus over degree-two topologies, and the gained speedup of ring over line.

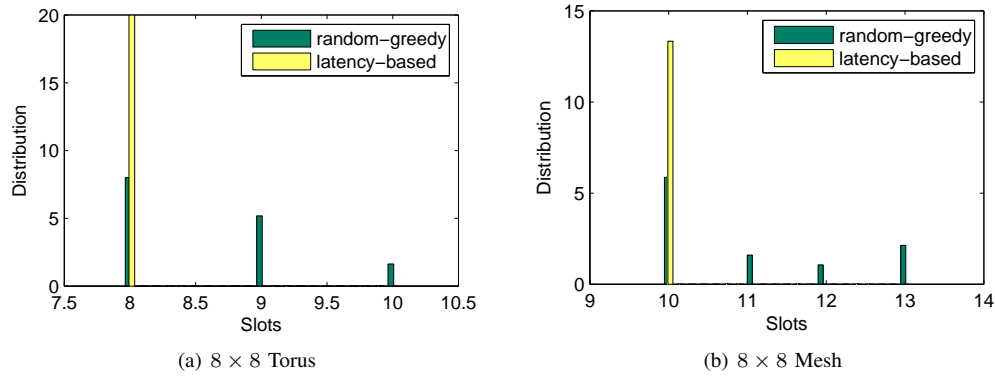


Fig. 6. Period Length Distribution with the *Bit-Permutation* traffic pattern. In all cases the latency-based algorithm achieves a slightly better distribution.

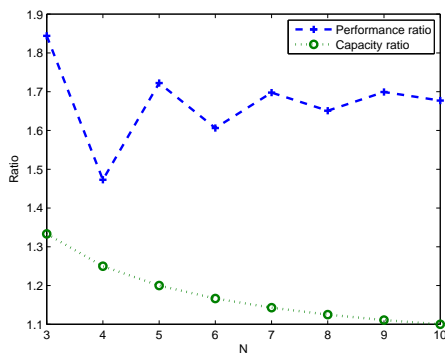


Fig. 7. Comparison of the performance ratio and the capacity ratio of the Torus and the Mesh using average online performance random-greedy algorithm. Note that although the Torus outperforms the Mesh in the bandwidth performance and the capacity, it has also a twice larger *link cost ratio*.

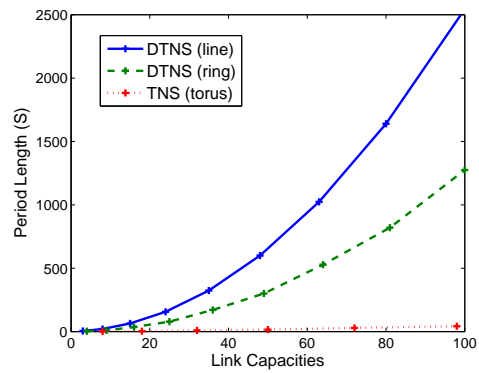


Fig. 9. Period length as function of the total link capacities of the optimal algorithms in line, ring and torus topologies. The capacity of a single link between two adjacent nodes is equal to 1.

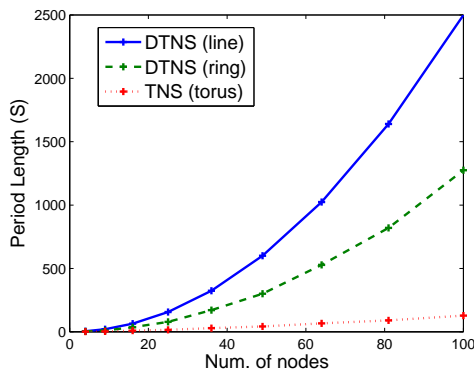


Fig. 8. Period length as function of number of nodes of the optimal algorithms in line, ring and torus topologies.

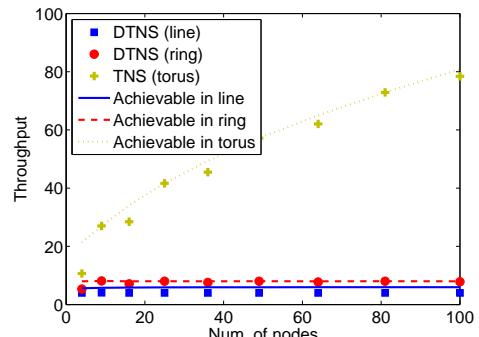


Fig. 10. Comparison of throughput of DTNS and TNS algorithms vs. the achievable capacity as function of the number of nodes.

Next, we compare the throughput of the DTNS algorithm in the line and the ring, and the TNS algorithm in the torus, with the achievable capacity in the networks under varying number of nodes. Figure 10 shows that the DTNS in a line topology is not capable of achieving full capacity, because the links in the

middle of the line transmit more packets than the links closer to the edges of the line. However, DTNS in ring topology achieves full capacity under all number of nodes, as proved in Theorem 4. Finally, TNS in a torus topology achieves full capacity with odd  $N$ , and almost achieves capacity with even  $N$ . Note that it is optimal in the sense that no other schedule can achieve a better throughput in a torus topology, as proved

in Theorem 7. Finally, also note that given a number of nodes, it is the torus topology that scales best, since it grows as  $\Theta(\sqrt{n})$ .

## 9 CONCLUSION

In this paper, we provided several periodic scheduling algorithms for bufferless NoCs, that are designed to meet the hard communication deadlines of real-time applications, applications based on the BSP programming model, and network congestion-control applications. In particular, we introduced DTNS and TNS scheduling algorithms that were proved to be optimal for complete exchange traffic on degree-two and torus networks. We also provided an application of the TNS algorithm on mesh NoCs, and showed that it achieves a constant bounded schedule length compared to the optimal scheduling. The performance-guarantee-based framework for the complete-exchange traffic pattern can also be used for other patterns, which require an almost uniform communication between the nodes, leaving some of the scheduling slots empty. Then, we provided a Latency-based scheduling algorithm for general periodic traffic pattern. We showed that latency-based scheduling algorithm is more efficient than random greedy scheduling on such NoC topologies as rings, tori and mesh.

In future work, we would like to investigate the following open problems: proving that the collision-free scheduling problem for periodic arbitrary traffic is NP-hard, providing improvements to our Latency-based algorithm, and to apply our results to other collision-free networks, like optical networks, which consist of arbitrary topologies.

## ACKNOWLEDGMENTS

The authors would like to thank Mark Silberstein, Ori Rotenstreich and Tsahee Zidenberg for their helpful comments. This work was partly supported by European Research Council Starting Grant No. 210389, the Intel ICRI-CI Center, and by the Technion Funds for Security Research. Alex Shpiner was supported by an Hasso-Plattner-Institut fellowship, and Erez Kantor was supported by an Eshkol fellowship from the Israel Ministry of Science and Technology and by NSF grants Nos. CCF-1217506, CCF-0939370 and CCF-AF-0937274.

## REFERENCES

- [1] A. L. Shimpi, "Intel's Sandy Bridge architecture exposed," [www.anandtech.com/show/3922/intels-sandy-bridge-architecture-exposed/4](http://www.anandtech.com/show/3922/intels-sandy-bridge-architecture-exposed/4), 2010.
- [2] T. Krishna, A. Kumar, P. Chiang, M. Erez, and L.-S. Peh, "NoC with near-ideal express virtual channels using global-line communication," *HotI*, Aug. 2008.
- [3] D. Vantrease, R. Schreiber, M. Monchiero, M. McLaren, N. P. Jouppi, M. Fiorentino, A. Davis, N. Binkert, R. G. Beausoleil, and J. H. Ahn, "Corona: System implications of emerging nanophotonic technology," *SIGARCH Comput. Archit. News*, vol. 36, pp. 153–164, June 2008.
- [4] J. Liu, J. Psota, N. Beckmann, J. Miller, J. Michel, J. Eastep, G. Kurian, L. Kimerling, A. Agarwal, and M. Beals, "ATAC: A manycore processor with on-chip optical network," MIT Dspace, Tech. Rep., 2009.
- [5] S. Le Beux, J. Trajkovic, I. O'Connor, G. Nicolescu, G. Bois, and P. Paulin, "Optical ring network-on-chip (ORNoC): Architecture and design methodology," *DATE*, 2011.
- [6] T. Chen, R. Raghavan, J. Dale, and E. Iwata, "Cell broadband engine architecture and its first implementation," [www.ibm.com/developerworks/power/library/pa-cellperf/](http://www.ibm.com/developerworks/power/library/pa-cellperf/), 2005.
- [7] M. Mirza-Aghatabar, S. Koohi, S. Hessabi, and M. Pedram, "An empirical investigation of mesh and torus NoC topologies under different routing algorithms and traffic models," *DSD*, 2007.
- [8] A. Guerre, N. Ventroux, R. David, and A. Merigot, "Hierarchical network-on-chip for embedded many-core architectures," *NOCS*, May 2010.
- [9] G. Michelogiannakis, D. Sanchez, W. J. Dally, and C. Kozyrakis, "Evaluating bufferless flow control for on-chip networks," *NOCS*, 2010.
- [10] T. Moscibroda and O. Mutlu, "A case for bufferless routing in on-chip networks," *ISCA*, 2009.
- [11] X. Chen and L.-S. Peh, "Leakage power modeling and optimization in interconnection networks," *ISLPED*, 2003.
- [12] A. B. Kahng, B. Li, L.-S. Peh, and K. Samadi, "Orion 2.0: a fast and accurate NoC power and area model for early-stage design space exploration," *DATE*, 2009.
- [13] A. Kodi, A. Sarathy, and A. Louri, "Design of adaptive communication channel buffers for low-power area-efficient network-on-chip architecture," *ANCS*, 2007.
- [14] Y. Hoskote, S. Vangal, A. Singh, N. Borkar, and S. Borkar, "A 5-ghz mesh interconnect for a teraflops processor," *IEEE Micro*, vol. 27, Sep. 2007.
- [15] D. W. Walker and J. J. Dongarra, "MPI: A standard message passing interface," *Supercomputer*, vol. 12, 1996.
- [16] R. Thakur and A. Choudhary, "All-to-all communication on meshes with wormhole routing," *Parallel Processing Symposium*, Apr. 1994.
- [17] Y.-C. Tseng and S. Gupta, "All-to-all personalized communication in a wormhole-routed torus," *Parallel and Distributed Systems*, vol. 7, no. 5, May 1996.
- [18] F. Petrini, "Total-exchange on wormhole k-ary n-cubes with adaptive routing," *IPPS/SPDP*, 1998.
- [19] Y.-J. Suh and S. Valamanchili, "All to-all communication with minimum start-up costs in 2d/3d tori and meshes," *Parallel and Distributed Systems*, vol. 9, no. 5, pp. 442–458, May 1998.
- [20] Y.-J. Suh and K. Shin, "Efficient all-to-all personalized exchange in multidimensional torus networks," *Parallel Processing*, 1998.
- [21] Y.-J. Suh and S. Yalamanchili, "Configurable algorithms for complete exchange in 2d meshes," *Parallel and Distributed Systems*, 2000.
- [22] D. Vainbrand and R. Ginosar, "Network-on-chip architectures for neural networks," *NOCS*, 2010.
- [23] L. G. Valiant, "A bridging model for parallel computation," *Commun. ACM*, vol. 33, 1990.
- [24] A. Singh, W. J. Dally, B. Towles, and A. K. Gupta, "Globally adaptive load-balanced routing on tori," *IEEE Comput. Archit. Lett.*, vol. 3, January 2004.
- [25] A. Hansson, K. Goossens, and A. Radulescu, "A unified approach to mapping and routing on a network-on-chip for both best-effort and guaranteed service traffic," *VLSI Design*, 2007.
- [26] R. Stefan and K. Goossens, "Multi-path routing in time-division-multiplexed networks on chip," *VLSI-SoC*, Oct. 2009.
- [27] O. Moreira, J. J.-D. Mol, and M. Bekooij, "Online resource management in a multiprocessor with a network-on-chip," *SAC*, 2007.
- [28] Z. Zhang, Z. Guo, and Y. Yang, "Bufferless routing in optical gaussian macrochip interconnect," *High-Performance Interconnects (HOTI)*, 2012 *IEEE 20th Annual Symposium on*, pp. 56–63, aug. 2012.
- [29] C. Gomez, M. E. Gomez, P. Lopez, and J. Duato, "BPS: A bufferless switching technique for NoCs," *Workshop on Interconn. Netw. Arch.*, 2008.
- [30] M. Hayenga, N. Jerger, and M. Lipasti, "Scarab: A single cycle adaptive routing and bufferless network," *MICRO*, 2009.
- [31] M. Schoeberl, F. Brandner, J. Sparsø, and E. Kasapaki, "A statically scheduled time-division-multiplexed network-on-chip for real-time systems," *NOCS*, 2012.
- [32] F. Brandner and M. Schoeberl, "Static routing in symmetric real-time network-on-chips," *RTNS*, 2012.
- [33] K. Goossens, J. Dielissen, and A. Radulescu, "Aethereal network on chip: concepts, architectures, and implementations," *Design Test of Computers, IEEE*, vol. 22, no. 5, 2005.
- [34] A. Hansson, M. Coenen, and K. Goossens, "Channel trees: reducing latency by sharing time slots in time-multiplexed networks on chip," *CODES+ISSS*, 2007.
- [35] A. Hansson, M. Subburaman, and K. Goossens, "Aelite: a flit-synchronous network on chip with composable and predictable services," *DATE*, 2009.
- [36] R. Stefan and K. Goossens, "An improved algorithm for slot selection in the Aethereal network-on-chip," *INA-OCMC*, 2011.

- [37] T. Marescaux, B. Bricke, P. Debacker, V. Nollet, and H. Corporaal, "Dynamic time-slot allocation for QoS enabled networks on chip," *Embedded Systems for Real-Time Multimedia*, 2005.
- [38] J. Liu, L.-R. Zheng, and H. Tenhunen, "A guaranteed-throughput switch for network-on-chip," *System-on-Chip*, 2003.
- [39] M. Millberg, E. Nilsson, R. Thid, and A. Jantsch, "Guaranteed bandwidth using looped containers in temporally disjoint networks within the nostrum network on chip," *DATE*, 2004.
- [40] Z. Lu and A. Jantsch, "TDM virtual-circuit configuration for network-on-chip," *IEEE Trans. on VLSI Systems*, 2008.
- [41] V. Dvorak and J. Jaros, "Optimizing collective communications on 2D-mesh and fat tree NoC," *ICN*, 2010.
- [42] J. Jaros, M. Ohlidal, and V. Dvorak, "Complexity of collective communications on NoCs," *Parallel Computing*, 2006.
- [43] H. C. de Freitas, L. M. Schnorr, M. A. Z. Alves, and P. O. A. Navaux, "Impact of parallel workloads on NoC architecture design," *Parallel, Distributed, and Network-Based Processing*, 2010.
- [44] I. Cohen, O. Rottenstreich, and I. Keslassy, "Statistical approach to networks-on-chip," *IEEE Trans. Computers*, vol. 59, no. 6, pp. 748–761, 2010.
- [45] C. Busch, M. Magdon-Ismail, M. Mavronicolas, and P. G. Spirakis, "Direct routing: Algorithms and complexity," *Algorithmica*, 2006.
- [46] Y.-C. Tseng, T.-H. Lin, S. Gupta, and D. Panda, "Bandwidth-optimal complete exchange on wormhole-routed 2d/3d torus networks: a diagonal-propagation approach," *Parallel and Distributed Systems, IEEE Transactions on*, vol. 8, no. 4, pp. 380–396, apr 1997.
- [47] N. S. Sundar, D. N. Jayasimha, and D. K. Panda, "Hybrid algorithms for complete exchange in 2d meshes," *IEEE Trans. Parallel Distrib. Syst.*, vol. 12, no. 12, pp. 1201–1218, 2001.
- [48] B. Juurlink, J. Sibeyn, and P. Rao, "Gossiping on meshes and tori," *Parallel and Distributed Systems, IEEE Transactions on*, vol. 9, no. 6, pp. 513–525, jun 1998.
- [49] A. Faraj, X. Yuan, and P. Patarasuk, "A message scheduling scheme for all-to-all personalized communication on ethernet switched clusters," *Parallel and Distributed Systems, IEEE Transactions on*, vol. 18, no. 2, pp. 264–276, feb. 2007.
- [50] I. Keslassy, M. S. Kodialam, T. V. Lakshman, and D. Stiliadis, "Scheduling schemes for delay graphs with applications to optical packet networks," *HPSR*, 2003.
- [51] Y. Kao and H. Chao, "Design of a bufferless photonic cros network-on-chip architecture," *Computers, IEEE Transactions on*, vol. PP, no. 99, p. 1, 2012.
- [52] A. Ben-Dor, S. Halevi, and A. Schuster, "Potential function analysis of greedy hot-potato routing," *Theory Comput. Syst.*, 1998.
- [53] C. Busch, "Ö(congestion + dilation) hot-potato routing on leveled networks," *SPAA*, 2002.
- [54] U. Feige and P. Raghavan, "Exact analysis of hot-potato routing (extended abstract)," *FOCS*, 1992.
- [55] M. Adler, A. L. Rosenberg, R. K. Sitaraman, and W. Unger, "Scheduling time-constrained communication in linear networks," *Theory of Computing Systems*, 2002.
- [56] J. S. Naor, A. Rosen, and G. Scalosub, "Online time-constrained scheduling in linear and ring networks," *Journal of Discrete Algorithms*, 2010.
- [57] S. Dolev, E. Kranakis, and D. Krizanc, "Baked-potato routing," *Journal of Algorithms*, 1999.
- [58] P. Gratz and S. W. Keelker, "Realistic Workload Characterization and Analysis for Networks-on-Chip Design," *The 4th Workshop on Chip Multiprocessor Memory Systems and Interconnects (CMP-MSI)*, 2010.
- [59] C. Feng, J. Li, Z. Lu, A. Jantsch, and M. Zhang, "Evaluation of deflection routing on various noc topologies," *ASIC (ASICON), 2011 IEEE 9th International Conference on*, pp. 163–166, 2011.



**Alexander Shpiner** is currently working in the system architecture group of Mellanox Technologies. He received a B.S. in computer engineering in 2004 and recently has completed a Ph.D. in electrical engineering at the Technion, Haifa, Israel. His research interests are in the fields of computer networks: congestion control, switch scheduling, network processors and networks on chip.



**Erez Kantor** is a post-doc fellow at the CSAIL, MIT, Cambridge, MA. Previously he was a post-doc fellow at the Electrical Engineering department at the Technion - Israel Institute of Technology. He is a recipient of Eshkol post-doc fellowship from the Ministry of Science and Technology, Israel. Erez received its Ph.D. and M.Sc. from the Weizmann Institute of Science- Israel in 2005 and 2009 respectively.



**Pu Li** was in the Department of Electrical Engineering, Technion, between November 2009 and July 2010. He is now working as New Technology Introduction Engineer in ASML, Netherlands. Previously, he received the B.S. and M.Sc. from Zhejiang University, China, and Technical University Eindhoven, the Netherlands, in 2007 and 2009, respectively. He conducted multiple research projects in the Chinese University of Hong Kong, TNO Netherlands and Philips Research Europe, in the field of communication network. His research interests include wireless technologies, high-performance switching network and network-on-chip.



**Israel Cidon** is a professor of Electrical Engineering at the Technion - Israel Institute of Technology. In 2012, he co-founded Sookasa, a provider of unstructured data security and management for storage cloud services. In 2005-2010 he was the dean of the Electrical Engineering Faculty at the Technion. In 2000, he co-founded Actona Technologies (acquired by Cisco in 2004), which pioneered the technology of Wide Area File Systems (WAFS). This technology, also termed WAN optimization, was

adopted by most world large enterprises. In 1998 he co-founded Viola Networks a provider network and VoIP performance diagnostic suite (acquired by Fluke Networks, 2008). In 1994-5 he founded and managed the high-speed networking group at Sun Microsystems Labs, Mountain View. Between 1985-94 he was the manager of the Network Architecture and Algorithms at IBM T. J. Watson Research Center, NY, where he led several computer networking projects, where he lead several computer networks projects including the first implementations of a packet based multi-media network and IBM first storage area network. For these works, he received twice the IBM outstanding innovation award. In 1981 he co-founded Micronet Ltd., an early vendor of mobile data-entry terminals (Israeli IPO, 2006). He is the co-author of over 170 refereed papers and 27 US patents. His current research involves wire-line and wireless communication networks and on chip interconnects networks.



**Isaac Keslassy** (M'02, SM'11) received his M.S. and Ph.D. degrees in Electrical Engineering from Stanford University, Stanford, CA, in 2000 and 2004, respectively.

He is currently an associate professor in the Electrical Engineering department of the Technion, Haifa, Israel. His recent research interests include the design and analysis of high-performance routers and on-chip networks. He is the recipient of the European Research Council Starting Grant, the Alon Fellowship and the

Yanai Award.