



HAL
open science

A Segmented Adaptive Router for Near Energy-Proportional Networks-on-Chip

Maxime France-Pillois, Abdoulaye Gamatié, Gilles Sassatelli

► **To cite this version:**

Maxime France-Pillois, Abdoulaye Gamatié, Gilles Sassatelli. A Segmented Adaptive Router for Near Energy-Proportional Networks-on-Chip. ACM Transactions on Embedded Computing Systems (TECS), 2022, 21 (4), pp.1-27/40. 10.1145/3529106 . hal-03724047

HAL Id: hal-03724047

<https://hal.science/hal-03724047v1>

Submitted on 15 Jul 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A Segmented Adaptive Router for Near Energy-Proportional Networks-on-Chip

MAXIME FRANCE-PILLOIS, LIRMM, Univ Montpellier, CNRS, Montpellier, France

ABDOULAYE GAMATIÉ, LIRMM, Univ Montpellier, CNRS, Montpellier, France

GILLES SASSATELLI, LIRMM, Univ Montpellier, CNRS, Montpellier, France

A Network-on-Chip (NoC) is an essential component of a chip multiprocessor (CMP) which however contributes to a large fraction of system energy. The unpredictability of traffic across a NoC frequently involves an expensive over-sizing of NoC resources which in turn leads to a significant contribution to the CMP power consumption. There exists a body of work addressing this issue, however so far solutions fall short when aiming for power reduction whilst maintaining high NoC performance. This paper proposes to combine router architecture optimizations with smart resource management to overcome this limitation. Based on a fully segmented architecture, we present an online adaptive router adjusting its active routing resources to meet the current traffic demand. This enhanced power-gating strategy significantly decreases both static and dynamic power consumption of the NoC, up to 70% for synthetic traffic patterns and up to 58% for real traffic workloads, while preserving NoC latency and throughput. Thanks to these adaptive power-saving mechanisms the proposed segmented NoC router provides near energy-proportional operation across the range of used benchmarks.

CCS Concepts: • **Networks** → **Network on chip**; **Network dynamics**; • **Computer systems organization** → *Interconnection architectures*; *Real-time system architecture*.

Additional Key Words and Phrases: Network-on-Chip, Energy Efficiency, Dynamic Power Management

ACM Reference Format:

Maxime France-Pillois, Abdoulaye Gamatié, and Gilles Sassatelli. 2022. A Segmented Adaptive Router for Near Energy-Proportional Networks-on-Chip. *ACM Trans. Embedd. Comput. Syst.* 1, 1, Article 1 (January 2022), 27 pages. <https://doi.org/10.1145/3529106>

1 INTRODUCTION

Whether homogeneous or heterogeneous, chip multiprocessors (CMPs) are at the heart of the majority of electronic devices, in desktop/server-class systems or embedded devices alike. The observed momentum around edge computing further boldens that shift from single-core microcontrollers to CMPs in meeting computing needs. These parallel processing architectures put a significant pressure on the interconnect subsystem which, unless properly sized, may constitute a performance bottleneck. Network-on-Chips (NoCs) [13] therefore emerged as the defacto communication architecture template for ensuring low-latency/high-bandwidth communications between CMP cores and the memory subsystem.

Authors' addresses: Maxime France-Pillois, maxime.france-pillois@lirmm.fr, LIRMM, Univ Montpellier, CNRS, Montpellier, Montpellier, France; Abdoulaye Gamatié, abdoulaye.gamatie@lirmm.fr, LIRMM, Univ Montpellier, CNRS, Montpellier, Montpellier, France; Gilles Sassatelli, gilles.sassatelli@lirmm.fr, LIRMM, Univ Montpellier, CNRS, Montpellier, Montpellier, France.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.

1539-9087/2022/1-ART1 \$15.00
<https://doi.org/10.1145/3529106>

Given the wide variety of possible traffic patterns, typical design decisions often rely on worst-case assumptions leading to significant NoC overdesign, which in turn incurs a significant power consumption overhead even under light traffic conditions. In many application scenarios, however, CMPs handle rather sporadic (event-driven) processing requests, which makes idle/light load power consumption account for a predominant fraction of the total energy consumption. As this particular consideration is not captured by the unit energy efficiency displayed for most NoCs, we advocate for the design of energy-proportional NoCs, i.e. NoCs for which power consumption is meant to be proportional to the traffic [6].

It has been reported that up to 28% of the tile power consumption [29] can be attributed to the NoC alone [4, 15], which motivated a number of investigations aiming at reducing NoC power consumption/increase NoC energy efficiency, without putting a specific focus on energy-proportionality.

These optimization efforts can be classified according to the following taxonomy:

- Router microarchitectures: specific optimizations proposed for the router hardware microarchitecture. [16, 18, 21]
- Advanced power management techniques: fine or mid-grain power gating, dynamic voltage and frequency scaling (DVFS). [11, 12, 45, 48]
- Routing/Flow control strategies: energy is saved by either bypassing router stages or entire routers in a fast-lane fashion. [33, 38]

In most cases, however, energy savings are either attained at the expense of performance degradation or achieve moderate benefits under asymmetric light traffic, which therefore does not satisfactorily contribute to decreasing the total energy consumption for systems whose typical average load is low.

Problem formulation. Given the aforementioned limitations of the approaches in energy-efficient NoC router literature, we advocate for the design of NoC routers capable of matching "on-the-fly" the amount of active (i.e. powered) routing resources to the traffic requirements. We promote this concept for enabling significant reductions in the leakage power that remains under no-load conditions, as well as similar savings under the temporal fluctuations and spatial heterogeneity of traffic demand, such as depicted in Figure 1, which shows some features of the *Fluidanimate* PARSEC benchmark [8].

Figure 1(b) represents the number of packets received by each core, arranged in an 8x8 grid, during the Region Of Interest (ROI) of this benchmark. We note eight hotspot nodes that concentrate most of the NoC traffic. They correspond to the location of the off-chip memory controllers. Figure 1(a) shows the temporal fluctuations of the Packet Injection Rate (PIR) during the ROI for select cores as well as the total PIR, which illustrates the optimization potential we intend to exploit through on-the-fly control of active router resources.

Proposed solution. We analyze the Roundabout router architecture introduced by Effiong *et al.* in [18] as having suitable concepts for devising the intended fine-grained control of router resources. The Roundabout router template has unique properties that we can exploit for our own purpose: purely distributed control per lane, no central shared resource such as a crossbar, buffers that are shareable across flows. We here consider routers as being made of *segments*, each of which is conceptually a minimal sufficient resource capable of performing the routing function i.e. routing an incoming packet to any output port.

We further devise a custom distributed power management strategy to reduce both static and dynamic power consumption whilst removing the energy consumption related to the structural components such as crossbar and arbiter.

This paper makes the following contributions:

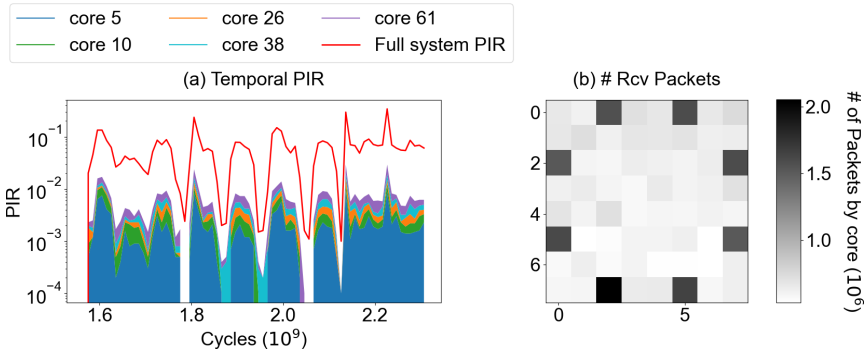


Fig. 1. (a) Temporal variation of the Packet Injection Rate (b) Number of packets received by cores for Fluidanimate benchmark running on 64 cores

- Temporal analysis of the PARSEC benchmarks workload [8] and proposal of guiding principles to shape NoCs for meeting realistic workload requirements.
- Design of an adaptive smart power gating solution on ‘self-sufficient’ router segments, enabling routers to dynamically adapt their active resources to traffic.
- Design of an adaptive dynamic segment assignment solution, capable of dynamically attaching segments to input ports where traffic pressure is high.
- Demonstration of NoC energy gains of up to 70% on synthetic traffic patterns and up to 58% on PARSEC benchmarks’ workloads compared with a conventional router offering similar performance, i.e. latency and throughput.

Outline. The rest of this paper is organized as follows: Section 2 discusses related works in the domain of NoC power optimization; Section 3 analyzes realistic workload properties that motivate our design decisions; Section 4 revisits the principles of the Roundabout router; Section 5 introduces the ‘self-sufficient’ router segment concept and describes the Segment Power-Gating Management service and Dynamic Segment Assignment approaches; Section 6 presents and analyzes the experimental results, and Section 7 draws conclusions.

2 RELATED WORK

A significant body of work on NoC power reduction exists. Three main approaches can be identified: 1) optimization/simplification of the router microarchitecture, 2) application of power management strategies to NoC and 3) alternative flow control/routing techniques, often referred to as bypass techniques. The literature review proposed in this section focuses on the first two approaches. The third [33, 38] is not discussed further due to being orthogonal to our work.

2.1 Router architecture optimization

Conventional routers are composed of four basic blocks [13]: 1) Input buffers, 2) Arbiters, 3) Crossbars and 4) Output buffers. Input buffers are typically favored over output buffers, as they are often better exploited by through-traffic early in the router pipeline. Most of the buffering capacity is therefore pinned to input ports, which incurs significant leakage and dynamic power consumption.

The second most power-consuming element is the crossbar, which sometimes accounts for more than 50% of the overall router power under high traffic conditions (see Section 6).

One popular approach for decreasing router power consumption lies in reducing input buffer power consumption. Input buffer size heavily influences high-load NoC performance [48]. However, application traffic patterns typically are temporally and spatially correlated, meaning that at any given time only a subset of input buffers are heavily used. This leads to an overall significant under-utilization of input buffers as reported in [21]. Several works attempt to increase buffer usage by means of buffer sharing. The RoShaQ router [43] implements central buffers or shared queues to store incoming flits. However, to provide an acceptable packet latency delay, this solution requires inserting a small buffer for each input. When traffic load is high, incoming flits have to cross two crossbars and are stored twice: first in the small input buffer, then in the shared queues. This results in a power overhead. The same issue was seen in other proposals implementing a central shared buffer [26, 40].

Inspired by the Dynamically Allocated Multi-Queue buffer, originally presented by Tamir *et al.* in [42], ViChaR [37] proposes input port buffers shared between all Virtual Channels (VCs). This solution makes better use of the input buffer free slots so that buffer size can be reduced without drastic performance degradation. Nevertheless, the logic required to manage these unified buffers is quite complex and consumes a significant amount of power itself. As a consequence, with an identical buffer size, ViChaR consumes more power than a conventional solution. However, when buffer size is reduced, ViChaR becomes competitive from a power consumption point of view, but the router throughput is impacted. More recently, Farrokhbakht *et al.* proposed UBERNoC [21], a router composed of shared input buffers with a single crossbar. This proposal does not introduce specific routing restrictions, which is often the downside of shared buffer routers. But, as observed with its predecessors, this solution induces router throughput and latency penalties when traffic load increases.

Another body of work targets the crossbar, which is the second most power-consuming component in a conventional NoC router. Some works propose the reduction of the crossbar's complexity [3, 16]. With the same perspective, Das *et al.* [14] divided the initial 5-to-5 crossbar into two 2-to-2 crossbars according to the row-column directions. This new crossbar, with an optimized switch allocation technique, improves the energy efficiency of the router by up to 20% under high traffic loads. In [32], Kim introduced priority in the router arbitration to decrease both arbitration and crossbar complexities, reducing the crossbar power consumption of about 30%. While these works exposed significant dynamic power reduction, they did not target the leakage power which is a prominent source of energy consumption with the sporadic traffic workloads.

Abad *et al.* [2] reviewed the classical router architecture and proposed removing the crossbar and global arbiter from the router in favor of a ring-based router. Incoming packets are injected into a buffered ring and travel on it until reaching their output port. This approach shares some similarities with the Roundabout router [18], however, the pure circular flow of packets is prone to deadlock, which requires advanced flow control techniques. Furthermore, this NoC uses store-and-forward flow control, which requires large buffers capable of storing at least one entire packet, contrary to most other NoCs that rely on wormhole. More recently, Effiong *et al.* [18] developed a concept for routers based on wormhole flow control. The Roundabout architecture removes the crossbar and global arbiter of conventional routers while benefiting from shared buffers. Since it only implements open-rings, this solution does not require specific flow control to avoid deadlock with an XY routing strategy. For that reason, we selected Roundabout as a suitable template for our investigations.

2.2 Smart power management

A body of work promotes DVFS [27, 34] for decreasing dynamic energy consumption in NoCs. However, variations in voltage and frequency incur a significant penalty on packet latency. DVFS

control engines therefore require to take into account application-level performance requirements [45]. To overcome this issue, some recent studies exploit machine learning techniques for deciding in a predictive manner voltage and frequency levels for meeting performance requirements [12, 23]. Beyond the dynamic power consumption addressed with DVFS, it is necessary to reduce the static power as well. To fill this demand, Clark *et al.* [12] combine power gating and DVFS. They report notable power savings in NoCs, but the resulting throughput is deteriorated due to wake-up latency.

Power gating is frequently used to reduce static power consumption in NoCs. This technique is applied at a component level as well as a router level. When power gating is performed at the router level, the wake-up delay incurs a NoC performance overhead, even when strategies are adopted for alleviating this issue [10, 11, 20, 22, 41]. The Power Punch solution [11] claims to completely hide the router's wake-up latency by sending wake-up signals up to three hops in advance. This is efficient albeit costly as it requires a dedicated wake-up wiring network. Proactive router power-gating solutions [39, 41] use by-passing and re-routing to avoid waking up power-gated routers. However, the new packet route may cross more routers, increasing the dynamic power of these routers [47].

A more flexible approach consists in performing power gating at a finer grain, on router components themselves. In [35], Matsutani organizes a router in micro-power domains that can be independently power-gated. Since the number of logic gates is limited in each power domain, the wake-up delay is very reduced. Fine-grain power gating has also been leveraged to ensure a minimum service by only switching off the extra buffer entries [31, 48] which makes for an interesting contribution towards achieving energy-proportionality.

Since we aim at power-gating router segments mainly composed of small buffers, this last technique shares some similarities with the present work. However, our approach differs in three fundamental aspects: 1) it relies on a fully segmented architecture, i.e. Roundabout, enabling to negate the power consumption related to structural components such as arbiters and crossbars; 2) power management decisions are taken at the router level and neither require extra inter-router links nor use bandwidth for transferring load monitoring information; 3) our strategy for segment activation reduces static as well as dynamic power.

In conclusion, the literature review reveals the difficulty in designing a router that is truly power-efficient, approaching energy proportionality, in that the overwhelming majority of power-optimized routers undergo performance degradation compared to their own baselines. Our proposal combines architectural optimizations together with traffic-adaptive power management to design power-efficient, near energy proportional routers that preserve router throughput and latency.

3 REALISTIC WORKLOADS ANALYSIS

NoC performance, e.g. latency, throughput, is frequently evaluated with synthetic traffic patterns such as uniform, transpose, or bit reverse. Although these patterns exhibit different characteristics, they do not account for the transient fluctuating nature of many real application traffics [5, 24, 46]. Typical NoCs are either 1) bandwidth-optimized, when targeting specific embedded systems running dataflow applications, or 2) latency-optimized for general-purpose CMP in which most traffic is cache-management related: cache misses, invalidations, etc. For NoC targeting general-purpose CMP, traffic workload is unpredictable and often implies router oversizing, which incurs additional area and energy consumption.

Some previous studies have already dealt with real application features [5, 24, 46] though none of them clearly exhibits the traffic variation during the application execution. We therefore performed this analysis so as to decide which features are key to ensure both high-performance and energy savings across the whole range of traffic patterns and intensities.

We analyzed traces from the Netrace utility tool [28]. These traces are extracted from the execution of the PARSEC benchmarks [8] on the M5 simulator [9] for a 64-cores shared-memory

CMP system. The CMP configuration is detailed in Table 2. This traffic mostly relates to cache-management operations i.e. cache-miss and invalidation traffic which is highly latency-sensitive. *Simmedium* input sets is used for PARSEC benchmarks except for *Bodytrack* and *Swaption* that use *simlarge* input sets.

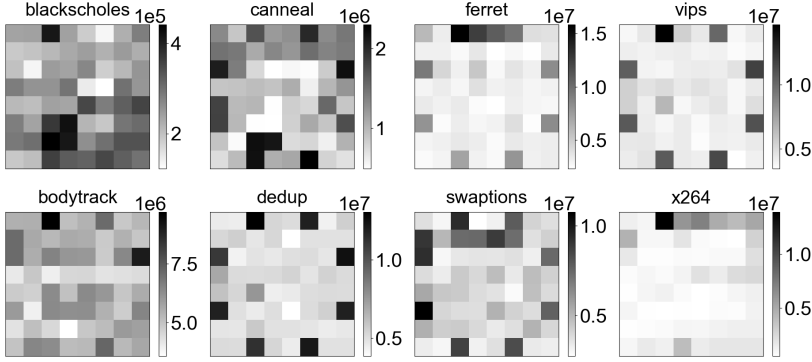


Fig. 2. Number of packets received by cores for 8 benchmarks of the PARSEC suite running on 64 cores (8x8 array)

Figure 2 shows the number of packets received along the ROI for eight PARSEC benchmarks. The 64 cores are organized onto an 8-by-8 Mesh, with their IDs translated into X-Y coordinates. Whichever benchmark we consider, a relatively large imbalance in the number of packets received by cores is observed, with no obvious symmetry in the traffic patterns. These structural differences reveal that some routers are more stressed than others, which makes for overall inefficient use of routing resources of the homogeneous NoC template.

Figure 3 shows stacked plots of temporal variations in the Packets Injection Rate (PIR). Throughout the ROI of eight PARSEC benchmarks, we plotted the PIR of five randomly selected cores and the PIR of the full 64-core system represented by the red curve. The PIR is processed using a 100-cycle sliding window. Plotted PIR are averaged over 100K samplings so as to filter out jitter and improve readability. Clear traffic patterns emerge along the benchmark’s execution and the following observations can be made:

- (1) From one benchmark to another, the Full system PIR (i.e. the sum of the PIR’s contributed by all cores for a benchmark) significantly differs.
- (2) Within each benchmark, the Full system PIR fluctuates during the ROI. For instance, in *bodytrack* and *swaptions*, some remarkable phases are observed. However, in *vips* the variation of the Full system PIR is somehow chaotic. At a fine-grain level, each core-related PIR also fluctuates: the number of packets injected by each core in the NoC varies over the time. For instance, this trend is noticeable for core 5 in the *blackscholes* benchmark, where the corresponding PIR varies between 0% and 4% during the ROI. The PIR of this core even increases to 6% for the *canneal* benchmark.

These qualitative observations highlight the interest of having adaptive routers able to match their internal active resources to the time-changing traffic requirements. Moreover, the spatial and temporal bursty behavior of realistic workloads calls for fine-grain NoC adaptive techniques, ideally at the granularity of router ports. Unfortunately, the usual power management techniques, e.g. DVFS, hardly permit devising such solutions in a practical manner. Indeed, the synchronization cost between different voltage/frequency islands is significant and typically results in coarse grain implementations i.e. with islands comprising several routers. Fine-grain NoC adaptivity however

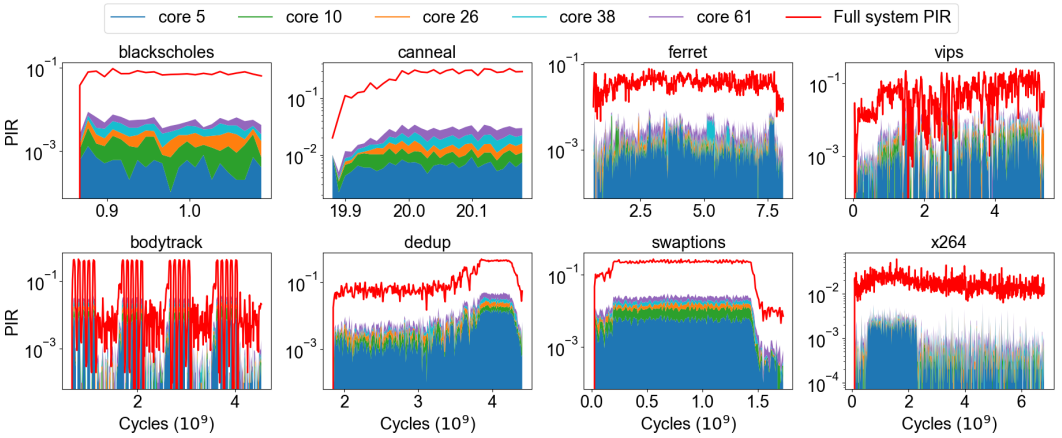


Fig. 3. Stacked plots of temporal variations in the packet injection rate for 8 benchmarks of the PARSEC suite

remains attractive from an energy efficiency standpoint, as conceptually allows activating maximum routing capacity solely *wherever* needed and *whenever* necessary.

4 THE ROUNDABOUT ROUTER

This section presents the basic principles of Roundabout as introduced by Effiong *et al.* in [18]. This router is inspired by real-life multi-lane roundabouts where cars go in one lane and then switch to a high-priority lane should they miss their exit.

Figure 4 illustrates the Roundabout architecture. Lanes are partitioned into primary and secondary lanes. Input ports inject packets into primary lanes, while secondary lanes only convey packets already inside the router. The number of lanes is a design parameter that can be decided according to the performance requirements, e.g. the maximum desired throughput. In this paper, we model an 11-lane router: five primary lanes represented by thick black lines in the figure, and six secondary lanes denoted by gray lines. This configuration provides two secondary lanes per router input port, which enables a fine grain power management of the router. Unlike the initial router principle, our model does not implement the notion of priority between primary and secondary lanes – both have an equal chance of obtaining shared resources (typically output ports). The priority concept enables the reduction of the router’s maximum latency by giving precedence to packets that experience larger latency delays inside a router, i.e. packets flowing through many lanes. However, this choice almost prevents packets from going out directly from the primary lane when the traffic load is high. From a power efficiency point of view, the longer the packet travels, the higher the energy consumed. Therefore, output controllers operate a round-robin policy to ensure fairness in the output port access granting. Note that this strategy may potentially increase the worst-case packet latency by $N-1$ cycles, where N is the number of primary lanes in a router.

Another strength of Roundabout lies in the intrinsic buffer sharing. In our implementation, two input ports can share the same lane and therefore the buffers that make up this lane. Lane sharing is, however, as enabling flows to share the same path, prone to deadlock. A comprehensive study of deadlock-freeness is therefore required and discussed in section 4.2. Physical implementations are extensively discussed in [19]. Roundabout exploits low-level flow control at lane-level, between buffers. This particular approach is such that it requires specific hardware constructs for performing handshaking between adjacent buffers, resulting in a form of back-pressure in the flit path. Two implementations are discussed, a pure asynchronous Roundabout version and a synchronous-elastic

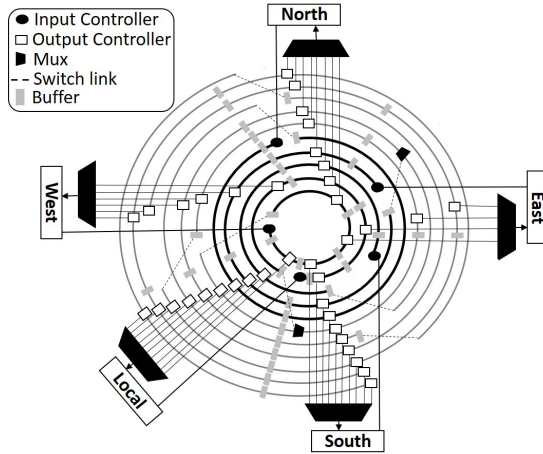


Fig. 4. Architecture of an 11-lane Roundabout (5 primaries, 6 secondaries)

[18]. Both have been studied using conventional physical synthesis design flows. They achieve absolute lower silicon area footprint and higher performance to area ratios compared to most other reported proposals extracted from the literature, including the seminal Hermes NoC router [36]. The distributed nature of the Roundabout router, including handshaking, makes for a lesser pressure on router-level wiring and thereby enables dense implementations.

4.1 Packet handling principle

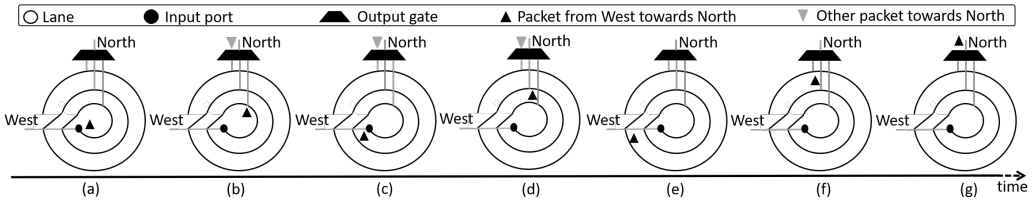


Fig. 5. Illustration of Roundabout principle

Here, we illustrate Roundabout functioning through an example in which an incoming packet destined for the North port enters from the West port. Figure 5 depicts the packets flowing on a simplified Roundabout architecture made of one primary lane and two levels of secondary lanes. The packet first enters the Roundabout into primary lane 0 (the innermost lane) from the West port input controller (see Figure 5(a)). Flits of the packet then flow in the lane until the head flit reaches the output controller corresponding to its destination, here the North port (see Figure 5(b)). The output controller checks the output availability and grants access to the output if available, i.e. not already in use by another packet coming from another lane. In case the output port is busy, the packet carries on in the same lane until its end where it then switches to the first-level secondary lane (see Figure 5(c)). It then circulates on that lane until it gets another chance to exit the router at the same port, where the same availability check takes place (see Figure 5(d)). Should the output port be once more busy the packet stays on the lane until it reaches the next secondary lane, which is the last-level secondary lane in that particular case (see Figure 5(e)). Since this lane is the last

opportunity for the packet to leave the router, when it reaches the output controller corresponding to its destination, the packet stops moving and waits until it obtains access to the output gate (see Figure 5(f) and (g)).

4.2 Deadlock-freeness

As stated in Section 2, Roundabout avoids deadlock thanks to its open-ring lanes and the use of the XY routing algorithm. However, when several input ports are injected into the same lanes (shared buffer strategy), deadlock can occur at a NoC level. To prevent deadlocks, the proposed Roundabout configuration should not contain cyclic resource dependencies. A preliminary “hand-made” analysis allows us to empirically derive some design rules in lane-sharing policy to avoid cyclic dependencies:

- (1) At most, two input ports can share a lane.
- (2) Two symmetric ports, i.e. South-North or East-West pairs, cannot share the same lane.
- (3) The local input port can only share a lane with another port performing moves along the West-East axis.

When applying the above rules, we obtain the router configuration depicted in Figure 4, where the West-Local ports and South-East ports share their secondary lanes, while the North port is alone in its lanes.

Then, we guarantee the deadlock-freeness of this router configuration by constructing its Control-Dependence Graph (CDG) for nine routers according to a 2D Mesh 3-by-3 grid. A 3x3 configuration is the minimal yet sufficient configuration for capturing all possible dependencies, due to the presence of a central router having 4 neighbors. As the CDG exposes all possible dependencies between resources, the absence of cycles formally ensures deadlock-freeness [13, 17].

This cyclic dependency analysis is performed using the *NetworkX* python library [1], as a manual analysis on such a large graph is both tedious and error-prone. Such an analysis is enough to ensure the deadlock-freeness for XY routing on 2D Mesh when simple flow control is employed. However, more complex control flows, such as the Worm-Bubble Flow Control, can be used at the input controller level to enable arbitrary routing algorithms and topologies.

5 A SEGMENTED ROUTER

We propose two complementary approaches, based on Roundabout router architecture, to make the power consumption closely proportional to the data traffic. Our first approach, called *Power Saver*, leans towards low-power NoC devices. It aims at drastically reducing both static and dynamic power by means of adaptive power-gating. Our second approach, called *Dynamic Lane Assignment*, looks at the improvement of the throughput offered by a NoC. It thus dynamically assigns router resources of a Roundabout router to the input port that is most in demand.

The rationale of the approach lies in regarding a lane as a functionally self-sufficient routing tile, that we refer to as a *segment*. The amount of active segments influences the performance mostly due to the segment-level buffering. Therefore, smart on-the-fly management of active segments is the intended strategy for avoiding performance penalty and achieving energy-proportional operation.

5.1 Power Saver approach

At the initial state, only the primary lanes of the router are powered. They ensure an elementary service, allowing the routing of all incoming packets to their output ports. Some specific traffic load conditions must be satisfied before the extra lanes get “unlocked”. This “expansion” process takes place online, while packets travel across active lanes.

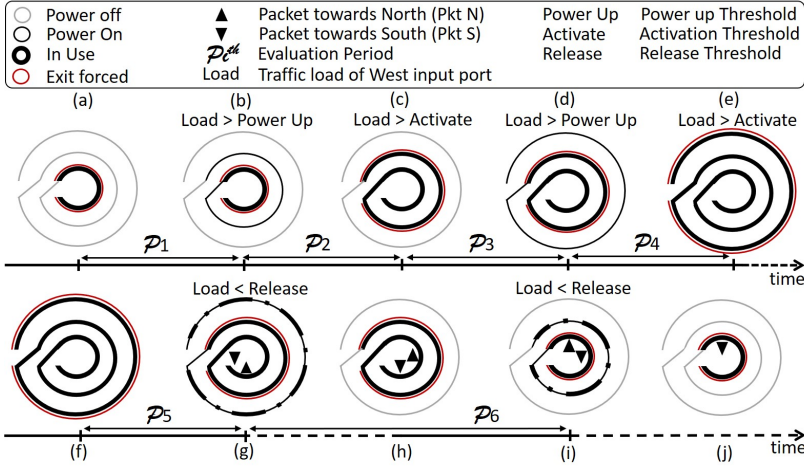


Fig. 6. Illustration of *Power Saver* lanes activation (a, b, c, d, e) / deactivation (f, g, h, i, j) scenario

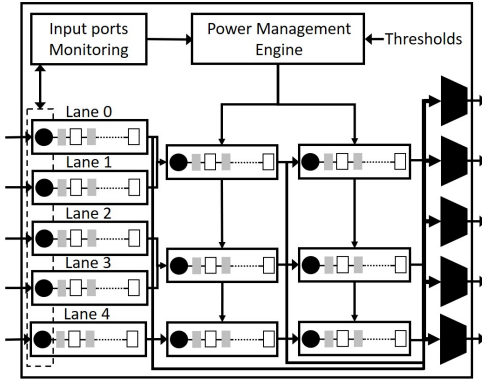


Fig. 7. Microarchitecture of the Roundabout with *Power Saver* engine

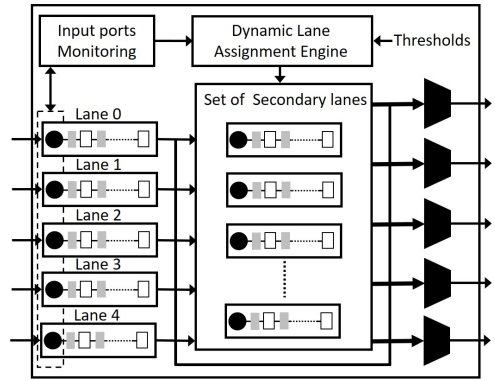


Fig. 8. Microarchitecture of the Roundabout with *Dynamic Lane Assign.* engine

Figure 6 shows a scenario of lanes activation/deactivation for a two-level secondary-lane router, similar to our 11-lane router configuration. For the sake of clarity, only the three-lane stages of the West input port are depicted in this figure. Initially, flits flow only in the primary lane (see Figure 6(a)). We can select the activated lanes online, by forcing packets to exit at a particular lane, illustrated by the red-color circles in Figure 6. When a packet reaches a busy output port: if the current lane is not tagged as "Exit Forced", the packet keeps traveling in the current lane until it switches to the next lane. However, if the current lane is tagged as "Exit Forced", the packet stalls until its output port becomes available, then exits the router. Afterwards, if the traffic load on an input port exceeds a predefined *Power Up* threshold, the first secondary lane is powered on (see Figure 6(b)). Then, if the load stress on the input port keeps increasing until it reaches a predefined *Activate* threshold (see Figure 6(c)), packets are authorized to switch to this extra lane. If the load pressure on the input port remains higher than the predefined *Power Up* threshold, the next level

secondary lane associated with this port is turned on (see Figure 6(d)). This lane then becomes active if the input stress remains higher than the predefined *Activate* threshold (see Figure 6(e)).

Conversely, when traffic load on an input port decreases under a predefined *Release* threshold, the last lane shutdown process is initiated. Subsequent packets are forced to exit from the previous lane while in-flight flits (see Figure 6(g)) are routed according to the regular policy, until the lane is drained, then the lane is powered down (see Figure 6(h)). In the figure, dashed lines depict this transitional state where lanes are partially used.

However, under heavy traffic, some flits may get blocked past their output port in the preceding lane, which may result in a deadlock. This phenomenon is illustrated in Figure 6 (g, h, i, j). Here, the black triangles represent two packets flowing from the West to the North ports, i.e. *Pkt N*, and coming from the West toward the South ports, i.e. *Pkt S*. In Figure 6(h), *Pkt S* reaches its destination port, but this port is busy. Since *Pkt S* flows in a lane that is not tagged as 'Exit Forced', it keeps moving on in the lane. But when *Pkt N* reaches its output port, it has to wait for the port to release since the current lane is now tagged as "Exit Forced". Next, in Figure 6(i), the traffic load on the West input port leads to the releasing of the last lane. Finally, *Pkt S* is blocked in the lane.

Afterwards, when *Pkt S* is able to move again, it has no possibility to leave the router, as the following lane has been since released (see Figure 6(j)). In fact, the secondary lane was considered free of flits, despite the incoming flit(s). We refer to this phenomenon as the *orphan flit* issue. We resolve the issue by adding a buffer at the end of each lane. If a flit enters this buffer while the next lane is off, the next lane is immediately powered on. Then, the flit can flow in this lane until it reaches its output port. The next lane is then again powered off once the *orphan flits* have been processed. This trick implies a small latency overhead for flits before they reach their output port, i.e. the lane wake-up time. However, this is only marginally detrimental to the NoC performance since *orphan flits* are rare. We evaluated the *orphan flit* occurrence during the execution of PARSEC benchmarks. *Orphan flits* appeared only during the execution of the x264 benchmark, at an average rate of 6.8 flits per million.

To evaluate the stress on an input port, we added a monitoring module inside routers. This module assesses a *busy rate* by counting the number of cycles the input controller is busy. The assessment is done over a predefined time slice, referred to as Evaluation Period ' Pi^{th} ' in Figure 6. This *busy rate* gives an estimate of the load at the input port. The decisions taken by the power management engine depend on this load.

Note that the two-stage activation: power on then start usage, enables the hiding of the lane wake-up latency, assuming the Evaluation Period is longer than the wake-up delay. It takes about 8 cycles to wake-up a router in 45nm [11].

Figure 7 depicts the microarchitecture of the Roundabout router implementing the *Power Saver* approach. As shown in the figure, input ports are connected to primary lanes. The monitoring engine tracks input port loads and reports corresponding *busy rates*. The power management engine takes power-up, activation, release decisions for all lanes linked to a given input port. These decisions are based on: current *busy rates*, predefined thresholds, and the state of the lanes. Note that each lane has direct access to the output gates, to route packets to the intended output.

5.2 Dynamic Lane Assignment approach

The realistic workload analysis, presented in Section 3, exposes large temporal variances in node injection rates, and also disparities in the destination nodes. We extended the flexibility of the Roundabout segmented architecture by dynamically assigning the extra lanes. The secondary lanes are no longer statically linked as depicted in Figure 4, but rather dynamically selected from a set of free lanes and attached where decided. This process is driven according to input port *busy rate*. If the *busy rate* of the most stressed input port is above a predefined threshold while a lane remains

in the set of free secondary lanes, the Dynamic Lane Assignment engine assigns this lane to this input port. It then proceeds to the powering and activation of the lane as described in the previous section. Next, when a lane is released (i.e. freed and powered off), the control engine registers it back into the set of free secondary lanes.

Figure 8 shows the microarchitecture of the Roundabout router implementing the *Dynamic Lane Assignment* approach. The input stage is similar to that of the *Power Saver* (see Figure 7). The main difference resides in the *Dynamic Lane Assignment* engine which decides how many lanes (from the secondary lane pool) are assigned to each primary lane. Power management (power-up, activate, release) is performed in a direct manner according to the lane assignment decisions, i.e. depending whether or not a lane is allocated.

To enable dynamic lane assignment, additional links must be added to connect each primary lane to all secondary lanes, and all secondary lanes to each other. This results in a small area and power consumption overheads (see Sections 6.3 and 6.4).

While the static lanes assignment ensures the availability of a next lane to evacuate the *orphan flits*, the *Dynamic Lane Assignment* may not provide this possibility. When all secondary lanes have been already (re)assigned before unblocking the *orphan flits*, a flit can enter into the additional buffer inserted at the end of the lane. This can happen while there are no remaining free lanes to give the flit the opportunity to exit. Since the deadlock-freeness requirement prevents from merely re-injecting the *orphan flits* into the router, which requires to guarantee at least one "rescue" lane per input port group is available. With the XY-routing strategy, three rescue lanes have to be reserved from the set of dynamic secondary lanes (see Section 4.2). Nevertheless, if the traffic load requires the use of all the resources, the rescue lanes can obviously be leveraged to drain the traffic. Note that the rescue lanes are still secondary lanes in that they are not powered by default, but only on-demand.

Thereby our 11-lane configuration is made of five primary lanes, three rescue secondary lanes, and three dynamic secondary lanes. The static configuration allows an input port to benefit from one to three lanes. The dynamic approach offers up to five lanes to drain the traffic load from an input port.

6 EVALUATION

In this section, we evaluate the two proposed approaches: *Power Saver* and the *Dynamic Lane Assignment* (*Dyn. Assign.* in short). We show the benefits of these approaches against conventional and state-of-the-art routers models.

6.1 Methodology

To assess the performance and power consumption of our approaches for synthetic traffic patterns and realistic application workloads, we use HNOCS [7], a fast and high-level discrete-event NoC simulator based on *Omnet++* [44]. This framework enables the simulation of full custom heterogeneous NoCs. However it has the following limitations: 1) only the classical XY routing algorithm is supported, 2) only uniform traffic patterns are supported, and 3) no power estimation is implemented.

Traffic injection. We enhanced the HNOCS framework to evaluate our approaches over broader traffic patterns. First, we implemented additional synthetic traffic patterns: *transpose*, *bit reverse*, *shuffle*, *hotspot* and *bit complement* [13]. Second, we extended the framework with the *Netrace* library [28] to enable the evaluation of realistic application workloads, beyond the aforementioned synthetic traffic patterns. *Netrace* provides traces of NoC communications occurring during the execution of the PARSEC benchmark suite on a 64-core CMP architecture. Hence, we exploit textitNetrace

to perform trace-based simulation in the HNOCS simulator. Since this library enforces packets dependencies, the application characteristics are preserved over NoC's performance changes.

Power estimation method. We implemented the analytical power model used by the *Orion3* library in the HNOCS simulation framework to perform both static and dynamic power estimations. Given the parameters of the target hardware technology, and the simulated toggle rate, this power model estimates the power consumption of a device with a small error: under 10% compared to a Register-Transfer Level model [30]. Table 1 summarizes the main technological parameters used in our study.

Table 1. Technology parameters

Manufacturing	Vdd	Frequency	Crossbar type
45nm	1.0V	650MHz	Matrix

We performed our comparative study of different router models using the same power-estimation library. We assume the link and clock energy consumption are similar for all evaluated router models. Therefore, we only modeled the power consumption of buffers, arbiters and crossbars. Since the Roundabout architecture does not implement any crossbar, we only estimate the power of buffers (in their in-lane configurations for Roundabout) and arbiters. Nevertheless, regarding the Roundabout version with *Dynamic Lane Assignment*, we need to consider the high number of added inter-lane connections (see Section 5.2). We therefore modeled these connections as small crossbars to fairly account for their corresponding contribution to the router power consumption.

Roundabout modeling. The HNOCS simulator only provides basic modules: links, sources, sinks and routers. Since a Roundabout lane is quite similar to a ring NoC topology with heterogeneous routers, we described each basic element of Roundabout, e.g. input controller, buffer, output controller, with fine-tuned HNOCS routers. This modeling enabled us to implement a buffer-level handshaking mechanism reflecting that of elastic buffers. Inter-router timing has been adapted to better adhere to the original Roundabout router, made of logic gates and elastic buffers [18]. Thanks to the flexibility of this modeling strategy, several router configurations can be easily modeled.

Consistency of Roundabout modeling. We checked our model accuracy against the performance and the power results given by Effiong in [19] for the RTL Roundabout implementation. We first modeled the exact same Roundabout configurations as that of [19]. Regarding the latency, in [19], Effiong gives for a 4-lane Roundabout the number of clock cycles for several paths. We tuned our model to feature these exact same latencies for properly accounting for the Roundabout specifics. Our model, therefore, mimics very closely the throughput and latency results given in that paper for synthetic traffics. Our power consumption estimations were consistently about 15% less than the post-synthesis accurate results displayed by Effiong, which is easily explained by our decision to leave link and clock power consumption aside. We therefore consider this accuracy enough for these investigations in which most analysis is further performed in a relative manner.

6.2 Simulation set-up

For our NoC model simulations, we used the parameters shown in Table 2. The CMP configuration support is provided by the *Netrace* library which enables to inject real application traffics. Note that all simulations are performed with a single Virtual Network (VN) and a single Virtual Channel (VC), as Roundabout provides no support for VCs and this is not required for our purposes. We avoid cache coherence protocol deadlocks by using the buffer over-sizing strategy for sink nodes [25]. This proves sufficient for this study, albeit not reflecting a realistic implementation which

would require devising several VNs. Roundabout could easily be extended to provide several VNs through lane replication, which is beyond the scope of this paper.

Table 2. Simulation Parameters

Cores	64 cores, in-order, Alpha ISA, 2GHz
L1 Cache	32KB Ins + 32KB Data, 4-way set, 3-cycle latency
L2 Cache	64 bank shared, 16MB, 8-way set, 8-cycle latency
Coherency	MESI coherence protocol
Memory	150-cycle latency, 8 on-chip memory controllers
Topology	8x8 2D-Mesh
Link Band.	32 bits/cycle
Packet size	10 flits
Flit size	4 Bytes
Flow Control	Wormhole
Virtual Chan.	1-VN, 1 VC/VN
Routing Algo.	XY

We compared our approaches, i.e. *Power Saver* and *Dynamic Lane Assignment*, against: two conventional routers, a shared buffer router, and a fine-grain power-gating router.

For conventional routers, we benefit from the HNOCS library that models a state-of-the-art three-stage pipeline router [7]. The depth of the VC buffer is usually chosen to cover the so-called Round-Trip Time (RTT) credit. A typical buffer size found in NoCs is four flits. Modeling Roundabout elements by routers imposes the presence of a minimal 1-flit buffer for each module composing lanes. The 11-lane Roundabout configuration thus implements at least 60-flit buffers. We then assessed the NoC performance for a common 4-flit input buffer HNOCS router, and a 12-flit input buffer router equivalent to our C11 Roundabout router (60-flit buffers).

The minimal Roundabout configuration introduced in [18] is a state-of-the-art optimized shared buffer router. We therefore compare our proposals to this router.

In Section 2, we stated that power gating approaches at the buffer entry level [31, 48] present similarity to our *Power Saver* approach. They make it possible to hide wake-up latency. Thus, we modeled a perfect power-gating solution on the buffer entries of a conventional HNOCS router for fairly accounting for state-of-the-art.

Table 3 summarizes the simulated router configurations.

Additional parameters are required for the online resource management in our approaches. When the rate of traffic load on an input port drops below 40%, i.e. the *Release* threshold, the last lane is freed then powered off. When this rate goes up to 60%, i.e. the *Power up* threshold, the next secondary lane assigned to this input port is powered on. When this rate increases above 80%, the powered, yet unused, lane is activated and then contributes to increasing the router buffering. The rate of the input load is evaluated over a period of 300 simulation cycles, in a sliding-window fashion. This value is chosen w.r.t. the mean time a packet needs to cross a router: about 30 *cycles* = 3 *cycles* * 10 *flits in a packet*, which corresponds to evaluating the traffic load over 10 consecutive packets. Thresholds and averaging window size are fixed for this study but can be tuned for meeting specific system requirements, though this is beyond the scope of this paper.

6.3 Performance evaluation

We first analyze benefits of our approaches for synthetic traffic patterns. It enables to analyze latency, power and energy efficiency across a wide range of Injection Rates.

Table 3. Router Parameters

Router Type	Flow control	Flit storage
HNOCS (ref)	3 stage-pipeline	12-flit/VC Total 60 flits
HNOCS small buffer (conv)	3 stage-pipeline	4-flit/VC Total 20 flits
HNOCS Power gating	3 stage-pipeline Power gating on buffer entries	12-flit/VC Total 60 flits
Roundabout Minimal (C0)	5 lanes (Prim. 3, Sec. 2)	Total 60 flits
Roundabout Full Power (baseline - C11)	11 lanes (Prim. 5, Sec. 6)	Total 60 flits
Roundabout (our approaches - C11) <i>Power Saver,</i> <i>Dynamic Lane Assignment</i>	11 lanes (Prim. 5, Sec. 6) Eval. Period: 300 cycles Release thld: 40% Power up thld: 60% Activ. thld: 80%	Total 60 flits

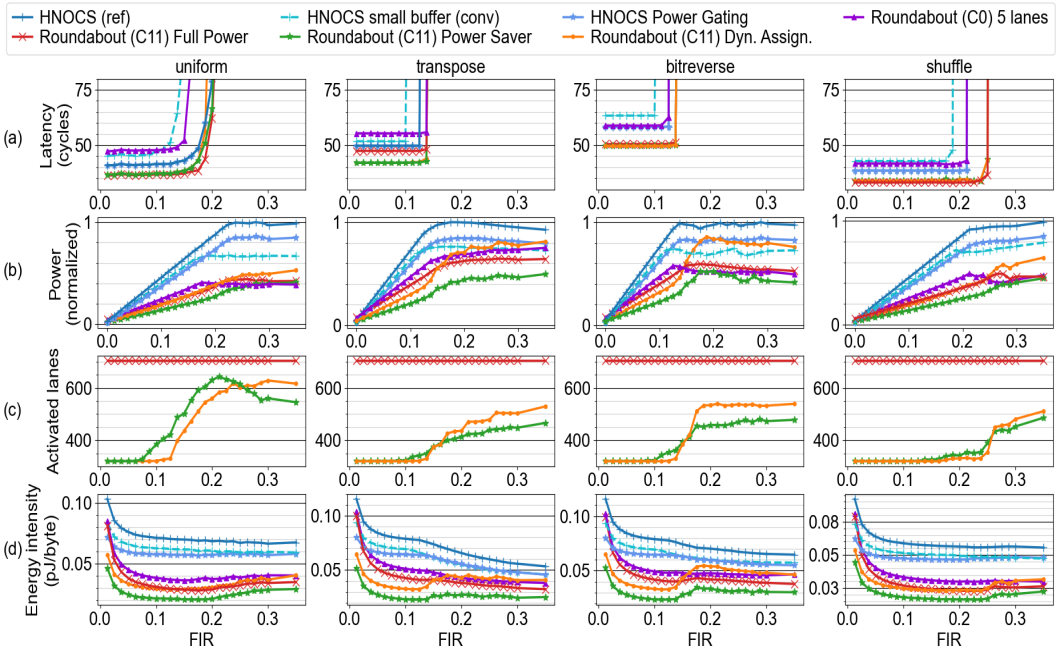


Fig. 9. (a) Latency, (b) Power, (c) Number of activated lanes, (d) Energy intensity for a 64-node 2D-Mesh NoC for synthetic traffic patterns

Synthetic traffic patterns. Figure 9 shows the performance evaluation results in terms of (a) latency, (b) power, (c) number of activated lanes and (d) energy-efficiency, according to the Flit Injection Rate (FIR) for the NoC configurations given in Table 3.

The performance of NoCs is reported in Figure 9(a), where the maximum offered throughput is given by the NoC saturation point. Our proposals offer a maximum throughput similar to that of the reference HNOCS for uniform and transpose traffic patterns. Slightly better results are achieved for bit reverse and shuffle traffic patterns. The C0 Roundabout configuration reveals that the sharing of primary lanes does not make it possible to efficiently handle the symmetric traffic patterns, i.e. uniform and shuffle. Performance degradation is however less with non-uniform traffic patterns, i.e. transpose and bit reverse patterns.

Figure 9(b) gives the power consumption for the full 64-router NoC, which are normalized w.r.t. the power consumption of the reference HNOCS configuration. Our proposal reduces the zero-load power by about 22%, 45% and 50%, compared to the reference HNOCS, Roundabout C0 and Roundabout C11 in Full Power configuration, respectively. When the traffic load increases, the *Power Saver* router (green curve) shows a larger power reduction compared to all other routers. In the operating region, i.e. about 2% below the FIR saturation threshold, both approaches reduce power consumption by at least 50% compared to the HNOCS conventional router with equivalent buffer size, labeled HNOCS (ref) in the figure. For uniform and shuffle traffic patterns, the *Power Saver* approach decreases the power by about 40%, compared to the Roundabout C0 configuration and more than 30% compared to the Roundabout C11 Full Power router. For transpose and bit reverse traffic patterns, the reduction is even more significant, with 50% and 40% compared to the Roundabout C0 and C11 configurations in Full Power mode, respectively. Regarding the *Dynamic Lane Assignment* approach, we observe a visible increase in power consumption resulting from the added inter-lane interconnect, while the NoC saturation point is not improved. This finds roots in the fact that this approach is tailored for bursty/fluctuating traffic requirements and not stationary synthetic traffic patterns.

Figure 9(c) depicts the number of active lanes according to the FIR for Roundabout-based configurations. For both adaptive approaches the number of active lanes grows in accordance with the injection rate around the saturation points. Due to the specific features of traffic patterns and the node-asymmetry phenomenon (routers at the edge of the mesh have unconnected ports), the number of active lanes never reaches the upper bound of 704 lanes, i.e. $11 \text{ lanes} * 64 \text{ routers}$. We notice a weak reactivity for Dynamic Lane Assignment approach caused by the lane granting policy. Indeed, only the most stressed primary lane is granted additional resource in each evaluation period, whereas the Power Saver approach may allocate many lanes at once, i.e. in a single evaluation period.

Figure 9(d) shows the energy intensity of the different NoC configurations. Roundabout-based routers consume significantly less energy than conventional routers thanks to their efficient shared buffer microarchitecture. The flatter the plot, the more energy-proportional the configuration, due to a constant energy cost for data transport. Table 4 shows the absolute plot slopes of the energy intensity. The smaller the slope, the more energy-proportional the solution is. Since we distinguish three distinct phases in these plots, we compute the slope for each phase. The first phase consists of a low traffic load, i.e. FIR under 5%. The second phase is the usual operating region for FIR ranging from 5% to 14%. The third phase corresponds to temporary traffic burst for FIR above 14%. Unsurprisingly, the most energy-proportional router is the ideal power-gating router which cuts power under no-load conditions. Our solutions exhibit a slight slope compared to other routers in the operating region, which indicates a near energy-proportional functioning in this region for all traffic patterns.

Table 5 summarizes the improvement of NoC efficiency achieved by our approaches.

Power savings. Figure 10 shows the power distribution for a 64-router NoC with uniform traffic patterns for: 60-flit buffer HNOCS conventional routers, 11-lane Full Power Roundabout and 11-lane *Power Saver* Roundabout. The power consumption is plotted for three key FIRs: at zero-load, at the

Table 4. Absolute slopes of the energy intensity plots for the three functioning phases

	Traffic	Low	Medium (operating region)	High
		FIR 1% – 5%	FIR 5% – 14%	FIR 14% – 35%
HNOCS (Ref)	Uniform	0.721	0.051	0.012
	Transpose	0.834	0.083	0.102
	Bit reverse	0.834	0.078	0.047
	Shuffle	0.759	0.056	0.008
HNOCS Power Gating	Uniform	0.317	0.019	0.004
	Transpose	0.365	0.030	0.082
	Bit reverse	0.363	0.023	0.040
	Shuffle	0.339	0.022	0.005
Roundabout Full Power	Uniform	1.150	0.081	0.031
	Transpose	1.311	0.088	0.044
	Bit reverse	1.311	0.088	0.015
	Shuffle	1.189	0.087	0.010
Roundabout Power Saver	Uniform	0.575	0.027	0.042
	Transpose	0.659	0.031	0.082
	Bit reverse	0.658	0.029	0.011
	Shuffle	0.596	0.043	0.028
Roundabout Dyn. Assign.	Uniform	0.635	0.038	0.056
	Transpose	0.725	0.008	0.001
	Bit reverse	0.724	0.015	0.013
	Shuffle	0.657	0.048	0.041

Table 5. Energy-efficiency gain of evaluated routers for Uniform and Transpose traffic patterns

	Traffic	Low	Medium	High
		FIR 1%	FIR 10%	FIR 20%
HNOCS (Ref)		—	—	—
HNOCS Power Gating	Uniform	31%	19%	16%
	Transpose	31%	18%	15%
Roundabout C11 Full Power	Uniform	21%	56%	58%
	Transpose	13%	46%	42%
Roundabout C11 Power Saver	Uniform	55%	70%	69%
	Transpose	55%	71%	59%
Roundabout C11 Dyn. Assign.	Uniform	44%	58%	55%
	Transpose	44%	59%	34%

middle of the operating region (FIR=0.1) and at the saturation point (FIR=0.2). Plot (a) gives the static power for each router component, while plot (b) shows the dynamic power.

In Figure 10(a), the static power of our *Power Saver* approach increases according to the FIR. This is caused by the power-gating technique, which cuts leakage power in the non-powered lanes. The static power consumption of our both approaches is also much less than the initial Roundabout model (i.e. Full Power model) at the saturation point. This is explained by the node-asymmetry phenomenon: routers at the edge of the Mesh have unconnected ports and, as such, unused lanes: the *Power Saver* approach maintains those lanes powered down as now traffic is detected.

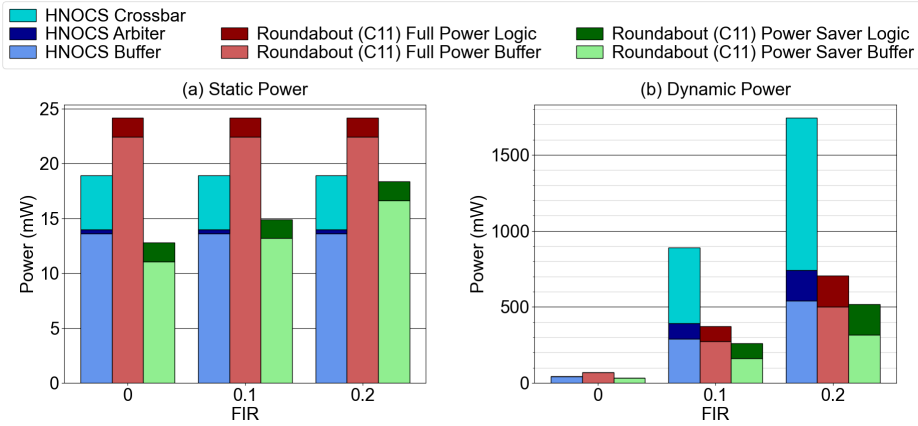


Fig. 10. (a) Static and (b) Dynamic power for a 64-node 2D-Mesh NoC composed of 60-flit buffer routers

Figure 10(b) shows a steep increase in dynamic power with FIR. Note the significant, expected, increase in the crossbar power consumption. These results also highlight the major contribution of dynamic power within the NoC power consumption when the traffic increases. This explains why a direct power-gating technique provides limited benefits, further motivates more elaborated resources management techniques to address both static and dynamic power consumption.

Quality of Service. The multi-lanes Roundabout architecture with no priority (see Section 4) may imply the delaying of some flits. This flit delaying would extend the packet latency. Hence, the monitoring of the packet latency deviation w.r.t. the flit latency allows us to empirically check this behavior. Figure 11 plots the average number of cycles between the flit latency and the packet latency for the four previously considered synthetic traffic patterns. We observe no significant variations between conventional router architectures, i.e. HNOCS routers, and the Roundabout routers, i.e. Roundabout C11 configurations. Regarding *bitreverse* traffic pattern, the average deviation tends to decrease for Roundabout routers due to the traffic pattern nature. These results, therefore, indicate the limited impact of the lack of priority for the Roundabout architecture for typical packet length.

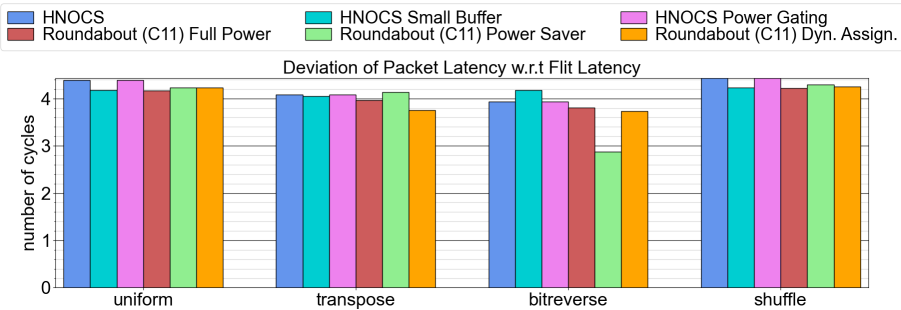


Fig. 11. Deviation of Packet Latency w.r.t. Flit Latency for a 64-node 2D-Mesh NoC for synthetic traffic patterns

Realistic traffic workloads. We use the *Netrace* library [28] to inject realistic workloads from PARSEC benchmarks into the NoC models made of: HNOCS reference routers, HNOCS ideal Power

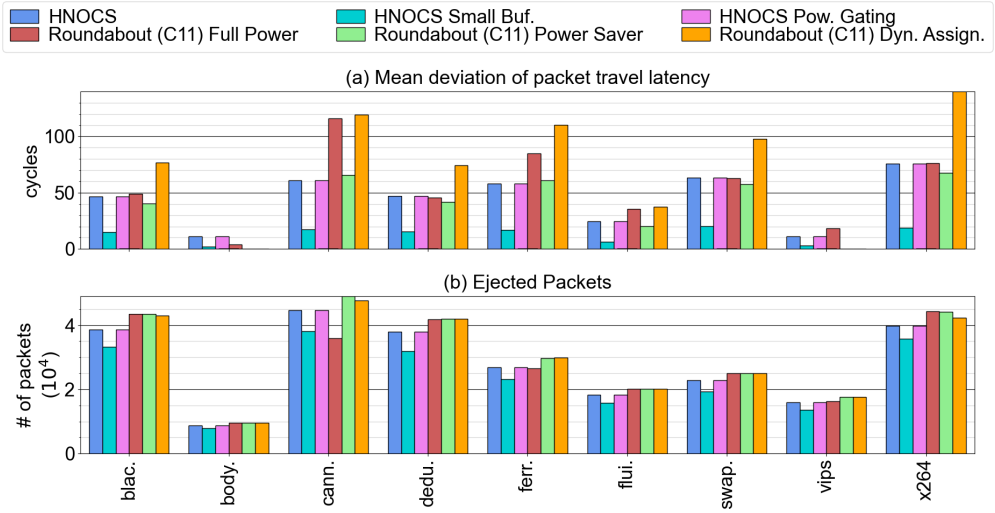


Fig. 12. Packet statistics for 64-node NoC over 100K cycles of PARSEC benchmarks

gating routers, Roundabout C11 Full Power routers, Roundabout C11 *Power Saver* and Roundabout C11 *Dynamic Lane Assignment*. Since the simulation of billions of cycles for each benchmark is too time-consuming, we focus on the first 100K cycles of each benchmark's ROI.

Figure 12(a) presents the mean deviation of packet travel latency between the no-load (minimal) latency and the actual packet latency for the PARSEC benchmarks. We observe a higher latency deviation for the initial Roundabout Full Power configuration than for the *Power Saver* approach. The latter adds restrictions on the used secondary lanes. While a single transient collision on an output gate involves switching to secondary lanes with the initial Roundabout scheme, this approach promotes short router crossing paths flowing on primary lanes only for moderate traffic workloads. We also observe a higher deviation latency for the *Dynamic Lane Assignment* approach since packets cross more secondary lanes when the traffic load increases. The *Power Saver* approach introduces on its side less delay than the conventional router for most benchmarks.

Figure 12(b) shows the number of packets processed by the NoC during the 100K first cycles of the ROI. Due to a possibly smaller router-crossing delay on short paths (e.g. local port to South port) and a better offered throughput, the Roundabout routers handle more packets than the conventional router for almost all benchmarks. Since the *Power Saver* policy promotes short paths inside routers compared to the initial Roundabout Full Power configuration, it manages to process more packets than other routers. The *Dynamic Lane Assignment* solution is here no better than *Power Saver*, due to the overall longer packet travel time incurred by the use of additional lanes.

Figure 13 gives the normalized energy breakdown of the evaluated NoCs during the same period. These results are normalized against HNOCS in term of number of packets, such that average normalized energies displayed relate to a similar number of conveyed packets. Therefore, we used the number of packets handled by the HNOCS solution as a reference to compute the energy consumption of the other configurations. Except for *bodytrack* benchmark, Roundabout routers outperform the conventional router of equivalent buffer size thanks to the removal of the crossbar and its expensive power consumption. For *bodytrack*, the obtained power is not far from zero-load power, since the amount of packets sent in the first 100K cycles is small (see Figure 12).

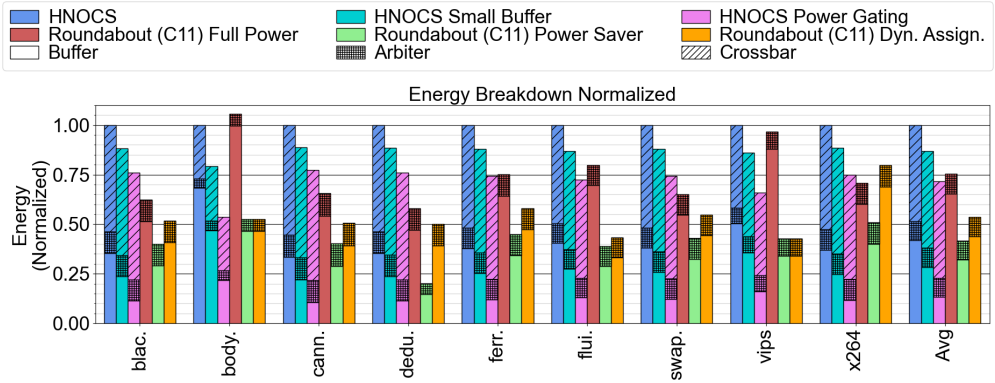


Fig. 13. Normalized energy breakdown for 64-node NoC over 100K cycles of PARSEC benchmarks

As a whole, significant energy consumption reductions are achieved by the proposed *Power Saver* and *Dynamic Lane Assignment* approaches. *Power Saver* achieves on average 58% energy savings across the whole PARSEC benchmarks, as summarized in Table 6. Similar to the behavior observed with synthetic traffic patterns, the *Dynamic Lane Assignment* approach consumes more energy than the *Power Saver* router due to its extra inter-lanes links, yet still outperforms both the conventional router and Full Power Roundabout.

Table 6. Energy reduction for PARSEC benchmark

Router Type	Average Reduction
HNOCS (Ref)	—
HNOCS Small Buffer	−13.1%
HNOCS Power Gating	−28.4%
Roundabout (C11) Full Power	−23.3%
Roundabout (C11) Power Saver	−58.5%
Roundabout (C11) Dyn. Assign.	−46.3%

Figure 14 illustrates the number of active lanes for the two proposed approaches during the 100K first cycles of the PARSEC benchmarks' ROI. Since the simulated period corresponds to the beginning of the ROI, we merely observe the warm-up phase with many activations and few deactivations. Nevertheless, We note that the lane utilization changes from one benchmark to another with various activation timestamps corresponding to specific software phases or events. Even though the variation tendencies are the same for both approaches, we remark that the dynamic assignment strategy offers more lanes to drain the traffic issued by software events.

6.4 Implementation overhead

As explained in Section 5, the dynamic resource management mechanisms applied in both approaches requires dedicated hardware. We here discuss this hardware overhead.

A counter of $\log(\text{Evaluation Period})$ bits is needed for each primary lane to monitor its load (*Number of Primary Lanes* \times 1 counter). Another counter is used to track the number of flits in a lane for each secondary lane (*Number of Secondary Lanes* \times 1 counter). The size of this counter depends on the number of flits able to co-exist in a lane (i.e. the buffer size). Approximately, $(\text{Total Number of Lanes} \times 5) + (\text{Number of Secondary Lanes} \times \text{Total Number of Lanes} - 1)$ registers

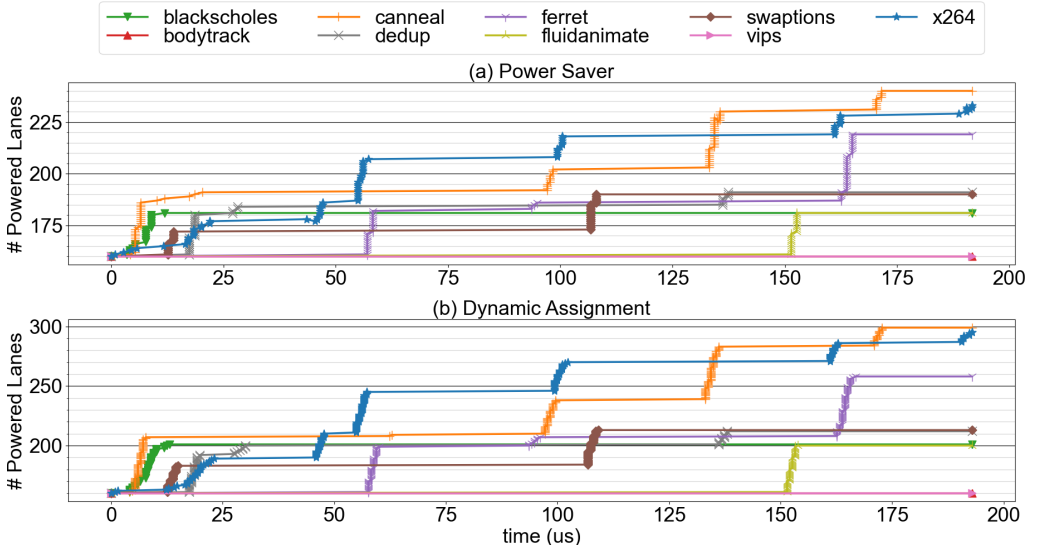


Fig. 14. Lanes usage over 100K cycles of PARSEC benchmarks

of $\log(\text{Total Number of Lanes})$ bits are required to store the current router state. In addition, some basic logic is also needed to implement a 5-stage Finite-State Machine. The *Dynamic Lane Assignment* approach is more expensive since it requires i) a table of free lanes with *Number of Dynamic Lanes* entries, ii) inter-lane links and an additional multiplexer for each secondary lane (with *Number of Primary Lanes*+*Number of Dynamic Lanes* inputs and one output port). Table 7 summarizes the additional hardware required by our proposal. We synthesized the 11-lane Roundabout configuration (C11) and the required supports for our two approaches with the Cadence Genus RTL compiler for 28nm FDSOI cell library. Table 8 displays the estimated area of these modules. The *Power Saver* approach has a minor impact on the total router area. The *Dynamic Lane Assignment* approach requires extensive modifications of the router microarchitecture to enhance its flexibility. Therefore, these improvements incur a larger area overhead.

Table 7. Hardware summary of our proposal

Router	# Entities	Type
Roundabout (C11) Power Saver	5	6-bits counters
+	6	4-bits counters
Roundabout (C11) Dyn. Assign.	105	4-bits registers
Roundabout (C11) Dyn. Assign. only	3	4-bits registers
	6	Mux (8 \rightarrow 1)
	24	Inter-lane links

In addition, the fine-grain power gating technique usually add a "sleep" transistor to logic cells to control the path to Vdd. This incurs a significant area overhead, but it can be mitigated by leveraging entire gate clusters such as lanes in our segmented router. In [35], Matsunatni *et al.* divided a classical router into 35 power domains, where each domain is driven by a separate power gating system. They reported that the power gating management hardware increases the area by less than 5%. In our case, we have only seven power domains. We therefore assume that the expected

Table 8. Hardware overhead of our proposal

Hardware modules	Area (in 28nm FDSOI)	Overhead
Roundabout (C11)	16972 μm^2	—
Power Saver engine	1401 μm^2	+8.26%
Dynamic Lane Assignment engine	3077 μm^2	+18.13%

overhead implied by the power gating system will be less than 5% for our proposal. Besides, Figure 10 shows that the dynamic power of the router dominates the static part. So, a less intrusive clock gating technique could also be considered.

6.5 Scalability

In this section, we investigate the scalability of our proposal. Since uniform and shuffle patterns are both symmetric non-deterministic traffic workloads, we limited the scalability study to only one of them. Similarly, transpose and bit-reverse patterns are both asymmetric and deterministic, so we applied the same restriction. Therefore, we assume that uniform and transpose patterns are representative enough to study scalability concerns.

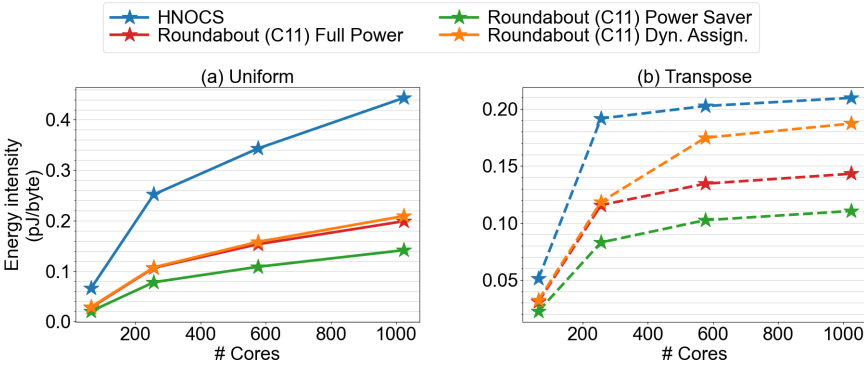


Fig. 15. Evolution of the minimum energy intensity according to the number of cores for (a) Uniform and (b) Transpose traffic patterns

6.5.1 Large Mesh topologies. Figure 15 shows the most energy-efficient working points for NoC topologies made of HNOCS and Roundabout routers. For the uniform traffic, the *Power Saver* approach scales better for large network sizes. Its corresponding energy intensity is lower than that of the Full Power router. Thus, it provides more energy savings. Regarding the transpose traffic, the *Power Saver* approach remains the best.

6.5.2 Importance of lane configuration. We evaluate and compare the scalability of our proposed power management approaches to baseline Roundabout configurations integrating more secondary lanes. The following router architectures are considered: C11 (5 primary + 6 secondary lanes), C14 (5 primary + 9 secondary lanes) and C17 (5 primary + 12 secondary lanes). Evaluations are performed on uniform traffic.

Figure 16(a) shows the NoC offered throughput for a 64-node 2D-Mesh NoC. Solid lines represent the throughput measured for the baseline Roundabout (i.e. Full Power) routers. Dashed lines show the results obtained for the Roundabout configurations while the *Power Saver* approach is applied.

We first notice the maximum throughput increases with the number of secondary lanes. For C14 and C17 configurations, the baseline and the *Power Saver* approach offer similar throughput. Baseline configuration C11, however, displays a throughput drop, which could be explained by the extended latency occurring at high traffic loads, though this would require more investigations.

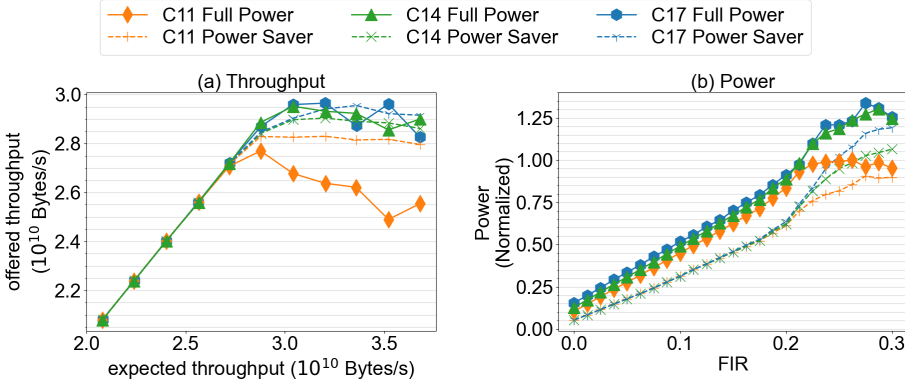


Fig. 16. Impact of Roundabout configuration on (a) Throughput, and (b) Power for a 64-node 2D-Mesh NoC for Uniform synthetic traffic

Despite the above observation on the baseline C11, we can argue that the *Power Saver* solution does not deteriorate the NoC throughput regardless the Roundabout configuration.

Note that the small variations of the throughput for the different Roundabout configurations results from the limited size of buffers chosen in this work. Indeed, larger buffer sizes would provide significant performance improvements as demonstrated by Effiong *et al.* in [18].

Figure 16(b) shows the power consumption normalized w.r.t. the C11 baseline (Full Power) consumption. It highlights the increase in efficiency of the *Power Saver* solution for configurations including a large number of secondary lanes. Note the power offset at 0 FIR for all baseline configurations which is removed in all *Power Saver* solutions, thereby underlining the near energy-proportional operation of the solution. Power consumption plots are further overlapping for C11, C14 and C17 in most of the operating region since the power gating is able to power down all unused lanes until communication pressures becomes high. These results confirm the scalability and near energy-proportionality of the *Power Saver* solution.

6.5.3 Effects of Additional Buffers. As mentioned in Section 4, the flexibility of the Roundabout architecture allows the addition of buffers along lanes. We evaluated our proposal for routers with additional buffers. Figure 17 plots the principal metrics impacted by the buffer resizing: the maximum offered throughput, the maximum power, and the mean power efficiency. We still consider a 64-node 2D-Mesh NoCs running uniform and transpose traffic patterns. The evaluated NoCs are composed of routers implementing buffers able to store 60-flit, 80-flit, and 110-flits. The 60-flit buffers configurations are the routers studied in the previous sections. The HNOCS 80-flit and 110-flit configurations are obtained by respectively allocating 16 and 22 flits to conventional input buffers. We model the Roundabout 80-flit configuration by uniformly adding 20-flits buffers to the primary lanes of the Roundabout C11 configuration. The Roundabout 110-flit configuration is achieved by uniformly adding 40-flits buffers to the primary and the secondary lanes of the Roundabout C11 configuration.

Regarding the offered throughput, the handled traffic load increases with the NoC storage capacity for uniform traffic. The ability of the uniform traffic to benefit from the major part of the

storage capacity explains this behavior. On the contrary, the transpose traffic pattern only uses and saturates a reduced number of router ports, i.e. buffers. Therefore, the increase of the NoC storage capacity does not significantly improve the offered throughput. This specific feature of the transpose pattern justifies the significant power reduction achieved by the power management techniques for this traffic compared to the uniform traffic. We also note that the Roundabout 110-flit configuration consumes less energy than the 80-flit configuration. The buffers of the former configuration are better spread out along the primary and the secondary lanes. Hence, the storage capacities are better balanced, and our power management approaches can operate on more buffers. As for the NoC power efficiency, we notice that Roundabout routers still over-perform conventional routers with and without Power Gating management technique for all considered buffer sizes. Therefore, unlike conventional routers, the balanced increase of storage capacity does not decrease the NoC power efficiency for Roundabout routers benefiting from our power management techniques.

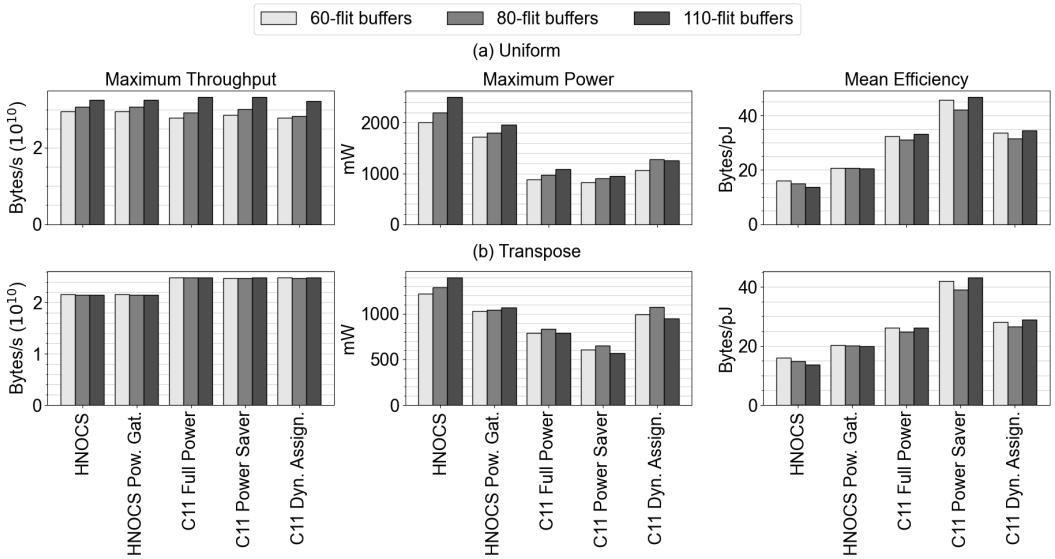


Fig. 17. Impact of Buffer Sizes on Throughput, Power and Efficiency for a 64-node 2D-Mesh NoC for (a) Uniform and (b) Transpose synthetic traffic patterns

7 CONCLUSION

Typical NoC router designs often fail to suitably address the spatial disparities and temporal fluctuations in communication demand resulting from the execution of applications on CMPs. This is partly due to active underused buffering resources, thereby resulting in limited energy efficiency. In this paper we present two resource-adaptive Roundabout router features to improve the energy-efficiency of NoCs toward energy-proportional operation. We specifically extended the router ability to adapt to bursty traffic workloads, by dynamically activating or attaching extra lanes (segments).

The first approach, called *Power Saver*, provides a highly energy efficient router configuration activating "on-the-fly" segments according to the current traffic demand. The second approach, i.e. *Dynamic Lane Assignment*, enhances the router flexibility by dynamically attaching segments to arbitrary input ports where traffic load is high.

The resulting segmented adaptive router achieves notably reduced static and dynamic power consumptions while preserving NoC offered throughput. Our proposed *Power Saver* NoC approach outperforms conventional routers such as HNOCS by improving energy-efficiency by up to 70% for synthetic traffics. It also reduces the energy consumption by 58% on average for real traffics, observed on ROI slices of the PARSEC benchmarks.

REFERENCES

- [1] [n. d.]. NetworkX – NetworkX documentation. <https://networkx.org/>
- [2] Pablo Abad, Valentin Puente, José Angel Gregorio, and Pablo Prieto. 2007. Rotary Router: An Efficient Architecture for CMP Interconnection Networks. In *Proceedings of the 34th Annual International Symposium on Computer Architecture (ISCA '07)* (San Diego, California, USA). Association for Computing Machinery, New York, NY, USA, 116–125. <https://doi.org/10.1145/1250662.1250678>
- [3] Jung Ho Ahn, Young Hoon Son, and John Kim. 2013. Scalable High-Radix Router Microarchitecture Using a Network Switch Organization. *ACM Trans. Archit. Code Optim.* 10, 3, Article 17 (Sept. 2013), 25 pages. <https://doi.org/10.1145/2512433>
- [4] Multicore Architectures, Sheng Li, Jung Ho Ahn, Jay Brockman, and Norman Jouppi. 2009. McPAT 1.0: An Integrated Power, Area, and Timing Modeling Framework for Multicore Architecture. (01 2009).
- [5] Mario Badr and Natalie Enright Jerger. 2014. SynFull: Synthetic traffic models capturing cache coherent behaviour. In *2014 ACM/IEEE 41st International Symposium on Computer Architecture (ISCA)*. 109–120. <https://doi.org/10.1109/ISCA.2014.6853236>
- [6] L. A. Barroso and U. Hölzle. 2007. The Case for Energy-Proportional Computing. *Computer* 40, 12 (2007), 33–37. <https://doi.org/10.1109/MC.2007.443>
- [7] Y. Ben-Itzhak, E. Zahavi, I. Cidon, and A. Kolodny. 2012. HNOCS: Modular open-source simulator for Heterogeneous NoCs. In *2012 International Conference on Embedded Computer Systems (SAMOS)*. 51–57. <https://doi.org/10.1109/SAMOS.2012.6404157>
- [8] Christian Bienia. 2011. *Benchmarking Modern Multiprocessors*. Ph. D. Dissertation. Princeton University.
- [9] N. L. Binkert, R. G. Dreslinski, L. R. Hsu, K. T. Lim, A. G. Saidi, and S. K. Reinhardt. 2006. The M5 Simulator: Modeling Networked Systems. *IEEE Micro* 26, 4 (2006), 52–60. <https://doi.org/10.1109/MM.2006.82>
- [10] L. Chen and T. M. Pinkston. 2012. NoRD: Node-Router Decoupling for Effective Power-gating of On-Chip Routers. In *2012 45th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. 270–281. <https://doi.org/10.1109/MICRO.2012.33>
- [11] L. Chen, D. Zhu, M. Pedram, and T. M. Pinkston. 2015. Power punch: Towards non-blocking power-gating of NoC routers. In *2015 IEEE 21st International Symposium on High Performance Computer Architecture (HPCA)*. 378–389. <https://doi.org/10.1109/HPCA.2015.7056048>
- [12] M. Clark, Y. Chen, A. Karanth, B. Ma, and A. Louri. 2020. DozzNoC: Reducing Static and Dynamic Energy in NoCs with Low-latency Voltage Regulators using Machine Learning. In *2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. 1–11. <https://doi.org/10.1109/IPDPS47924.2020.00011>
- [13] William Dally and Brian Towles. 2003. *Principles and Practices of Interconnection Networks*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
- [14] C. R. Das, M. S. Yousif, V. Narayanan, D. Park, C. Nicopoulos, and J. Kim. 2006. A Gracefully Degrading and Energy-Efficient Modular Router Architecture for On-Chip Networks. In *33rd International Symposium on Computer Architecture (ISCA'06)*. 4–15. <https://doi.org/10.1109/ISCA.2006.6>
- [15] B. K. Daya, C. O. Chen, S. Subramanian, W. Kwon, S. Park, T. Krishna, J. Holt, A. P. Chandrakasan, and L. Peh. 2014. SCORPIO: A 36-core research chip demonstrating snoopy coherence on a scalable mesh NoC with in-network ordering. In *2014 ACM/IEEE 41st International Symposium on Computer Architecture (ISCA'14)*. 25–36. <https://doi.org/10.1109/ISCA.2014.6853232>
- [16] G. Dimitrakopoulos, N. Georgiadis, C. Nicopoulos, and E. Kalligeros. 2013. Switch folding: Network-on-Chip routers with time-multiplexed output ports. In *2013 Design, Automation Test in Europe Conference Exhibition (DATE)*. 344–349. <https://doi.org/10.7873/DATE.2013.081>
- [17] J. Duato. 1995. A necessary and sufficient condition for deadlock-free adaptive routing in wormhole networks. *IEEE Transactions on Parallel and Distributed Systems* 6, 10 (1995), 1055–1067. <https://doi.org/10.1109/71.473515>
- [18] Charles Effiong, Gilles Sassatelli, and Abdoulaye Gamatie. 2017. Distributed and Dynamic Shared-Buffer Router for High-Performance Interconnect. In *Proceedings of the Eleventh IEEE/ACM International Symposium on Networks-on-Chip (NOCS '17)* (Seoul, Republic of Korea). Association for Computing Machinery, New York, NY, USA, Article 2, 8 pages. <https://doi.org/10.1145/3130218.3130223>

- [19] Charles Emmanuel Effiong. 2017. *Exploration of multicore systems based on silicon integrated communication networks*. Phd Thesis. Université Montpellier. <https://tel.archives-ouvertes.fr/tel-01944111>
- [20] Hossein Farrokhbakht, Hadi Mardani Kamali, Natalie Enright Jerger, and Shaahin Hessabi. 2018. SPONGE: A Scalable Pivot-Based On/Off Gating Engine for Reducing Static Power in NoC Routers. In *Proceedings of the International Symposium on Low Power Electronics and Design (ISLPED '18)* (Seattle, WA, USA). Association for Computing Machinery, New York, NY, USA, Article 17, 6 pages. <https://doi.org/10.1145/3218603.3218635>
- [21] Hossein Farrokhbakht, Henry Kao, and Natalie Enright Jerger. 2019. UBERNoC: Unified Buffer Power-Efficient Router for Network-on-Chip. In *Proceedings of the 13th IEEE/ACM International Symposium on Networks-on-Chip (NOCS '19)* (New York, New York). Association for Computing Machinery, New York, NY, USA, Article 1, 8 pages. <https://doi.org/10.1145/3313231.3352362>
- [22] H. Farrokhbakht, M. Taram, B. Khaleghi, and S. Hessabi. 2016. TooT: an efficient and scalable power-gating method for NoC routers. In *2016 Tenth IEEE/ACM International Symposium on Networks-on-Chip (NOCS)*. 1–8. <https://doi.org/10.1109/NOCS.2016.7579326>
- [23] Q. Fettes, M. Clark, R. Bunescu, A. Karanth, and A. Louri. 2019. Dynamic Voltage and Frequency Scaling in NoCs with Supervised and Reinforcement Learning Techniques. *IEEE Trans. Comput.* 68, 3 (2019), 375–389. <https://doi.org/10.1109/TC.2018.2875476>
- [24] Paul V. Gratz and Stephen W. Keckler. 2010. Realistic Workload Characterization and Analysis for Networks-on-Chip Design. In *The 4th Workshop on Chip Multiprocessor Memory Systems and Interconnects (CMP-MSI)*.
- [25] Andreas Hansson, Goossens Kees, and Rădulescu Andrei. 2007. Avoiding Message-Dependent Deadlock in Network-Based Systems on Chip. *VLSI Design 2007* (04 2007). <https://doi.org/10.1155/2007/95859>
- [26] S. M. Hassan and S. Yalamanchili. 2013. Centralized buffer router: A low latency, low power router for high radix NOCs. In *2013 Seventh IEEE/ACM International Symposium on Networks-on-Chip (NOCS'13)*. 1–8. <https://doi.org/10.1109/NoCS.2013.6558397>
- [27] Robert Hesse and Natalie Enright Jerger. 2015. Improving DVFS in NoCs with Coherence Prediction. In *Proceedings of the 9th International Symposium on Networks-on-Chip (NOCS '15)* (Vancouver, BC, Canada). Association for Computing Machinery, New York, NY, USA, Article 24, 8 pages. <https://doi.org/10.1145/2786572.2786595>
- [28] Joel Hestness, Boris Grot, and Stephen W. Keckler. 2010. Netrace: Dependency-Driven Trace-Based Network-on-Chip Simulation. In *Proceedings of the Third International Workshop on Network on Chip Architectures (NoCArc '10)* (Atlanta, Georgia, USA). Association for Computing Machinery, New York, NY, USA, 31–36. <https://doi.org/10.1145/1921249.1921258>
- [29] Y. Hoskote, S. Vangal, A. Singh, N. Borkar, and S. Borkar. 2007. A 5-GHz Mesh Interconnect for a Teraflops Processor. *IEEE Micro* 27, 5 (2007), 51–61. <https://doi.org/10.1109/MM.2007.4378783>
- [30] A. B. Kahng, B. Lin, and S. Nath. 2015. ORION3.0: A Comprehensive NoC Router Estimation Tool. *IEEE Embedded Systems Letters* 7, 2 (2015), 41–45. <https://doi.org/10.1109/LES.2015.2402197>
- [31] G. Kim, J. Kim, and S. Yoo. 2011. FlexiBuffer: Reducing leakage power in on-chip network routers. In *2011 48th ACM/EDAC/IEEE Design Automation Conference (DAC)*. 936–941.
- [32] J. Kim. 2009. Low-cost router microarchitecture for on-chip networks. In *2009 42nd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. 255–266. <https://doi.org/10.1145/1669112.1669145>
- [33] A. Kumary, P. Kunduz, A. P. Singhx, L. Pehy, and N. K. Jhay. 2007. A 4.6Tbits/s 3.6GHz single-cycle NoC router with a novel switch allocator in 65nm CMOS. In *2007 25th International Conference on Computer Design*. 63–70. <https://doi.org/10.1109/ICCD.2007.4601881>
- [34] Li Shang, L. Peh, and N. K. Jha. 2002. Power-efficient Interconnection Networks: Dynamic Voltage Scaling with Links. *IEEE Computer Architecture Letters* 1, 1 (2002), 6–6. <https://doi.org/10.1109/L-CA.2002.10>
- [35] H. Matsutani, M. Koibuchi, D. Ikebuchi, K. Usami, H. Nakamura, and H. Amano. 2010. Ultra Fine-Grained Run-Time Power Gating of On-chip Routers for CMPs. In *2010 Fourth ACM/IEEE International Symposium on Networks-on-Chip (NOCS'10)*. 61–68. <https://doi.org/10.1109/NOCS.2010.16>
- [36] Fernando Gehm Moraes, Aline Mello, Leandro Möller, Luciano Ost, and Ney Laert Vilar Calazans. 2003. A Low Area Overhead Packet-switched Network on Chip: Architecture and Prototyping. In *IFIP VLSI-SoC 2003, IFIP WG 10.5 International Conference on Very Large Scale Integration of System-on-Chip, Darmstadt, Germany, 1-3 December 2003*, Manfred Glesner, Ricardo Augusto da Luz Reis, Hans Evekling, Vincent John Mooney III, Leandro Soares Indrusiak, and Peter Zipf (Eds.). Technische Universität Darmstadt, Insitute of Microelectronic Systems, 318–323.
- [37] C. A. Nicopoulos, D. Park, J. Kim, N. Vijaykrishnan, M. S. Yousif, and C. R. Das. 2006. ViChar: A Dynamic Virtual Channel Regulator for Network-on-Chip Routers. In *2006 39th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. 333–346. <https://doi.org/10.1109/MICRO.2006.50>
- [38] Iván Pérez, Enrique Vallejo, and Ramón Bevide. 2019. SMART++: Reducing Cost and Improving Efficiency of Multi-Hop Bypass in NoC Routers. In *Proceedings of the 13th IEEE/ACM International Symposium on Networks-on-Chip (NOCS '19)* (New York, New York). Association for Computing Machinery, New York, NY, USA, Article 5, 8 pages.

<https://doi.org/10.1145/3313231.3352364>

- [39] Iván Pérez, Enrique Vallejo, and Ramón Bevide. 2021. S-SMART++: A Low-Latency NoC Leveraging Speculative Bypass Requests. *IEEE Trans. Comput.* 70, 6 (2021), 819–832. <https://doi.org/10.1109/TC.2021.3068615>
- [40] R. S. Ramanujam, V. Soteriou, B. Lin, and L. Peh. 2010. Design of a High-Throughput Distributed Shared-Buffer NoC Router. In *2010 Fourth ACM/IEEE International Symposium on Networks-on-Chip (NOCS'10)*. 69–78. <https://doi.org/10.1109/NOCS.2010.17>
- [41] A. Samih, R. Wang, A. Krishna, C. Maciocco, C. Tai, and Y. Solihin. 2013. Energy-efficient interconnect via Router Parking. In *2013 IEEE 19th International Symposium on High Performance Computer Architecture (HPCA)*. 508–519. <https://doi.org/10.1109/HPCA.2013.6522345>
- [42] Yuval Tamir and Gregory L. Frazier. 1992. Dynamically-Allocated Multi-Queue Buffers for VLSI Communication Switches. *IEEE Trans. Comput.* 41, 6 (June 1992), 725–737. <https://doi.org/10.1109/12.144624>
- [43] A. T. Tran and B. M. Baas. 2011. RoShaQ: High-performance on-chip router with shared queues. In *2011 IEEE 29th International Conference on Computer Design (ICCD)*. 232–238. <https://doi.org/10.1109/ICCD.2011.6081402>
- [44] András Varga and Rudolf Hornig. 2008. An Overview of the OMNeT++ Simulation Environment. In *Proceedings of the 1st International Conference on Simulation Tools and Techniques for Communications, Networks and Systems and Workshops (Marseille, France) (Simutools '08)*. ICST (Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering), Brussels, BEL, Article 60, 10 pages.
- [45] Yuan Yao and Zhonghai Lu. 2016. DVFS for NoCs in CMPs: A thread voting approach. In *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. 309–320. <https://doi.org/10.1109/HPCA.2016.7446074>
- [46] Jieming Yin, Onur Kayiran, Matthew Poremba, Natalie Enright Jerger, and Gabriel H. Loh. 2016. Efficient synthetic traffic models for large, complex SoCs. In *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. 297–308. <https://doi.org/10.1109/HPCA.2016.7446073>
- [47] Di Zhu, Yunfan Li, and Lizhong Chen. 2019. On Trade-off Between Static and Dynamic Power Consumption in NoC Power Gating. In *2019 IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*. 1–6. <https://doi.org/10.1109/ISLPED.2019.8824936>
- [48] Davide Zoni, Andrea Canidio, William Fornaciari, Panayiotis Englezakis, Chrysostomos Nicopoulos, and Yiannakis Sazeides. 2017. BlackOut: Enabling fine-grained power gating of buffers in Network-on-Chip routers. *J. Parallel and Distrib. Comput.* 104 (2017), 130–145. <https://doi.org/10.1016/j.jpdc.2017.01.016>