

**On the power efficiency, low latency, and quality of service in networkon-chip** Wang, P.

Citation

Wang, P. (2020, February 12). On the power efficiency, low latency, and quality of service in network-on-chip. Retrieved from https://hdl.handle.net/1887/85165

Version:	Publisher's Version
License:	<u>Licence agreement concerning inclusion of doctoral thesis in the</u> <u>Institutional Repository of the University of Leiden</u>
Downloaded from:	https://hdl.handle.net/1887/85165

Note: To cite this publication please use the final published version (if applicable).

Cover Page



# Universiteit Leiden



The handle <u>http://hdl.handle.net/1887/85165</u> holds various files of this Leiden University dissertation.

Author: Wang, P. Title: On the power efficiency, low latency, and quality of service in network-on-chip Issue Date: 2020-02-12

# Chapter 5

# **EVC-based Power Gating Approach**

**Peng Wang**, Sobhan Niknam, Sheng Ma, Zhiying Wang, Todor Stefanov, "EVC-based Power Gating Approach to Achieve Low-power and High Performance NoC" *in Proceedings of the Euromicro Conference on Digital System Design (DSD)*, Chalkidiki, Greece, 2019.

**I** which corresponds to **Contribution 3** introduced in Section 1.4. By using our EVC-based power gating approach, not only the packet latency increase caused by power gating can be further reduced, but also the power consumption can be reduced at high traffic workloads. The remainder of this chapter is organized as follows. Section 5.1 further elaborates on the problem of the low efficiency in the bypass path mentioned in Section 4.7. Section 5.2 summarizes the main contributions in this chapter. Section 5.3 briefly introduces the express virtual channel scheme (EVC). Then, we provide an overview of the related bypass-based power gating approaches in Section 5.4. It is followed by Section 5.5, which elaborates our EVC-based power gating approach. Section 5.6 introduces the experimental setup and presents experimental results. Finally, Section 5.7 concludes this chapter.

# 5.1 **Problem Statement**

As we have introduced in Chapter 4, bypass-based power gating approaches are more comprehensive to reduce the power consumption and the packet latency increase caused by the power gating. However, in many bypass-based approaches, there are only a few bypass latches to temporarily store packets on a bypass path. Before bypassing powered-off routers, packets have to be blocked until there are available bypass latches, which significantly undermines the efficiency of the bypass paths. As a result, in most of the bypass-based approaches, the bypass paths are not very efficient to transfer packets. For example, the bypass paths in [CP12, FTKH16, ZL18] and in our D-bypass power gating approach [WNM<sup>+</sup>19a] (Chapter 4) cannot continuously transmit packets via bypass powered-off routers. Even though the approach in [BHW<sup>+</sup>17] can continuously transmit packets via bypass powered-off routers, it has significant timing overhead and hardware overhead to recover the routing information that is lost in the powered-off routers. As a consequence, all aforementioned bypass-based approaches still have significant packet latency increase caused by the power gating. Furthermore, like most of the coarse-grained power gating approaches, these bypassbased power gating approaches cannot fully utilize the idle time of each component in a router to reduce the power consumption. When the traffic workload becomes high, most of the routers become busy and cannot be powered off to reduce the power consumption of a NoC. As a consequence, these bypass-based power gating approaches are effective in reducing the power consumption only at low traffic workloads.

# 5.2 Contributions

In order to overcome the aforementioned drawbacks, we propose an express virtual channel based (EVC-based) power gating approach. In our approach, multiple virtual bypass paths are pre-defined at design time. Packets can take these virtual bypass paths to bypass intermediate routers that can be powered-on or powered-off. When a packet takes a virtual bypass path, the sink router of the virtual bypass path is powered-on. There is sufficient amount of buffers in sink routers to hold packets. Thus, packets can continuously go through a virtual bypass path. Furthermore, compared with other bypass-based approaches [CP12, FTKH16, BHW<sup>+</sup>17, ZL18, WNM<sup>+</sup>19a] in which the packets can only bypass powered-off routers, in our EVC-based approach, packets can bypass powered-on routers as well. Thus, even at a high traffic workload, our approach also can reduce the power consumption by reducing the dynamic power. The specific novel contributions of this chapter are the following:

• We propose a specific distribution of virtual bypass paths on a NoC, which allows more packets to take the virtual bypass paths compared to the conventional EVC scheme [KPKJ07]. More importantly, we extend the router structure to guarantee that a virtual bypass path cannot be blocked by powered-off routers. Thus, by allowing packets to go through the virtual bypass paths without blocking, these packets can avoid experiencing the wake up process at the intermediate routers. Furthermore, based on our extended router structure, a certain transmission ability of the powered-off/being charged routers is kept to transfer packets going through the normal paths. In this way, the packet latency

increase caused by the power gating is further reduced. We also propose an effective power gating scheme to control the power switching of routers. Finally, we propose an approach to freeze virtual bypass paths in order to resolve starvation, which is a common issue in EVC-based NoCs [KPKJ07].

• By experiments, we show that our EVC-based power gating approach can effectively reduce the power gating negative impacts on the performance and power consumption. Thanks to the efficient virtual bypass paths, our EVC-based power gating approach can achieve lower network latency than the related approaches [MKWA08, WNWS17, ZL18, WNM<sup>+</sup>19a], even lower than a NoC without power gating. However, our EVC-based power gating approach causes unbalanced resource allocation, which results in slight performance penalty in the total execution time of real applications. Compared with a NoC without power gating, our EVC-based power gating approach achieves notable reduction of the power consumption, which is comparable with the related approaches [MKWA08, WNWS17, ZL18, WNM<sup>+</sup>19a]. Furthermore, by allowing packets to bypass powered-on routers as well, our approach consumes less power than the related approaches [MKWA08, WNWS17, ZL18, WNM<sup>+</sup>19a] under high traffic workloads.

# 5.3 Background

In order to better understand the novel contributions of this chapter, in this section, we introduce the conventional EVC [KPKJ07] scheme that allows packets to bypass intermediate routers along a virtual bypass path.

The express virtual channel (EVC) scheme [KPKJ07] is a classical virtual bypass technique. As shown in Figure 5.1(a), the virtual bypass paths (red dashed lines) are pre-defined on a NoC topology. These virtual bypass paths are implemented without the need for extra physical wires, but based on the virtual channels in a router to share the existing wires. The basic EVC router architecture is shown in Figure 5.1(b). Compared with the conventional router in Figure 2.5, in each input port, one EVC latch is added, and the virtual channels are partitioned into two groups, normal virtual channels (N-VCs) and express virtual channel (E-VCs). N-VCs are used to accept packets from neighbor upstream routers. E-VCs in the sink routers of the virtual bypass paths.

By allocating E-VCs to packets, the source router in a virtual bypass path can determine if the packet takes the virtual bypass path. For example, in Figure 5.1, assume that a packet has to be sent form *Router*00 to *Router*04. Based on the transmission distance, *Router*00 is aware that by taking the virtual bypass path from *Router*00 to *Router*03, the packet will have lower latency. So, *Router*00 treats this packet as



Figure 5.1: Express virtual channel.

an E-packet (the packet going through a virtual bypass path) and allocates one E-VC in *Router*03 for this packet. When the packet reaches *Router*01 and *Router*02, this packet is temporarily held in the EVC latch with the highest priority. Then, this packet is directly sent without experiencing the pipeline stages in *Router*01 and *Router*02, and reaches *Router*03. When this packet reaches *Router*03, this packet is stored at the allocated E-VC. *Router*03 knows this packet should go to the normal path to its destination *Router*04, and treats this packet as a N-packet (the packet going through the normal path between routers) and allocates a N-VC in *Router*04 for this packet. After experiencing the pipeline stages in *Router*03, this packet is sent to its destination *Router*04.

By taking virtual bypass paths, E-packets do not need to experience the pipeline stages in the intermediate routers. This implies that most of the components in the intermediate routers are unnecessary to transfer E-packets. This characteristic is attractive and promising for realizing a power gating NoC to allow packets to bypass powered-off routers. We exploit effectively this characteristic in this thesis to realize our EVC-based power gating approach.

### 5.4 Related Work

Several approaches propose a bypass-based power gating NoC. In NoRD [CP12], a virtual ring is pre-defined on a NoC, which works as a backup NoC. When a packet is blocked by a powered-off router, it can go along this virtual ring to bypass the powered-off router. However, limited by the low efficiency and poor scalability of

the virtual ring, packets may be detoured for a long distance to their destinations. As a consequence, NoRD has significant packet latency increase and is not suitable for large NoCs. In contrast, in our approach, we pre-define multiple virtual bypass paths, which are separately distributed on the whole NoC. Packets go along their shortest routing path and separately take these virtual bypass paths to bypass the powered-off routers. Thus, our EVC-based power gating approach has lower packet latency and better scalability.

In Turn-on on Turn (TooT) [FTKH16], a bypass path is pre-defined in the horizontal (X + /X -) and vertical (Y + /Y -) directions. Thus, packets can bypass a powered-off router if the packets do not need the powered-off router to change the transmission direction or to eject from the NoC. So, TooT does not need to frequently power on the powered-off routers and can more efficiently reduce the static power consumption. However, limited by a few bypass latches on a bypass path, packets have to be blocked until there are available bypass latches. As a consequence, the bypass paths are inefficient to transmit packets in order to bypass the powered-off routers and TooT still has significant packet latency increase. In contrast, in our EVC-based power gating approach, when a packet goes through a virtual bypass path, the sink router is powered on. Thus, there are more buffers to be used to hold packets and packets can continuously go through the virtual bypass path. As a consequence, bypass paths in our approach are more efficient than TooT in terms of transmitting packets, therefore the packet latency increase is reduced.

Similar to TooT, Fly-over [BHW<sup>+</sup>17] also allows packets to bypass powered-off routers in the horizontal (X + /X -) and vertical (Y + /Y -) directions but Flyover does not need to block packets to wait for available bypass latches between the neighbor routers. This is because Fly-over dynamically realizes the credit-based flow control [DT04] between the source router and the sink router on a bypass path to guarantee that there is no buffer overflow. When a source router transmits packets to bypass the intermediate powered-off routers, the sink router must be powered-on. Thus, there is sufficient amount of buffers available to be used to hold packets and Fly-over can continually transmit packets. However, Fly-over has to utilize a complex mechanism to realize the credit-based flow control between the source router and the sink router, which causes significant timing and hardware overhead. In contrast, in our EVC-based power gating approach, the virtual bypass paths are (static) pre-defined. Thus, our EVC-based approach has no such extra timing overhead.

In contrast to TooT and Fly-over, the bypass path in EZ-bypass [ZL18] is dynamically built to allow packets to bypass the powered-off routers in any direction. Thus, a packet can bypass a powered-off router, even when this router is required to change the transmission direction. As a result, EZ-bypass is more flexible and can be more efficient to reduce the power consumption. However, in EZ-bypass, when a packet bypasses powered-off routers, this packet has to stay in the powered-off routers for multiple clock cycles to experience the pipeline stages of routers. As a consequence, the bypass latch is occupied by one packet for a long time and the bypass path is frequently blocked, which undermines the efficiency of the bypass path. In contrast, in our EVC-based power gating approach, when a packet bypasses intermediate routers, this packet does not experience the router pipeline stages. Thus, our EVC-based power gating approach can achieve lower packet latency than EZ-bypass. Furthermore, compared with NoRD [CP12], TooT [FTKH16], Fly-over [BHW<sup>+</sup>17], and EZ-bypass [ZL18] in which the packets can bypass only powered-off routers, in our EVC-based approach, packets can bypass powered-on routers as well. Thus, even at a high workload traffic, our approach also can reduce the power consumption by reducing the dynamic power.

# 5.5 Our EVC-based Power Gating

In this section, we present our novel power gating approach which uses and extends the conventional EVC scheme to allow packets to bypass powered-off routers. First, in Section 5.5.1, we propose a distribution of the virtual bypass paths to allow more packets to take the virtual bypass paths. Then, in Section 5.5.2, we extend the EVC router structure to guarantee that the virtual bypass paths are not blocked by powered-off routers. Thus, packets can always take a virtual bypass path to bypass the intermediate routers that may be powered-off. Furthermore, based on our extended router structure, a powered-off router has certain transmission ability to transfer also packets that take the normal paths. So, even though some packets do not take a virtual bypass path, they can avoid as much as possible to be blocked by powered-off routers. In Section 5.5.3, we describe our power gating scheme used in our EVC-based power gating approach, and in Section 5.5.4, we use an example to illustrate our power gating scheme. Finally, in Section 5.5.5, we propose an approach to resolve the starvation which may occur when using our EVC-based power gating approach.

### 5.5.1 Distribution of Virtual Bypass Paths

In the EVC scheme, packets can bypass the intermediate routers only when they take virtual bypass paths. So, in order to allow packets to bypass the intermediate routers that may be powered-off, we have to allow more packets to take the virtual bypass paths. To achieve this goal, in each direction, we pre-define one virtual bypass path between each two routers with three hops. As shown in Figure 5.2(a), in the X+ direction, we set one virtual bypass path between Router00 and Router03, Router01 and Router04 and so on. The virtual bypass paths in the X-, Y+, and Y- directions have similar settings, but are not shown in Figure 5.2(a) for the sake of clarity.



Figure 5.2: Extended EVC-based power gating approach.

Compared with the conventional distribution of the virtual bypass paths [KPKJ07] in Figure 5.1(a), the packets in Figure 5.2(a) have higher probability to take a virtual bypass path. For example, in a  $8 \times 8$  2D mesh NoC, there are in total 4032 routing paths from one source router to a destination router. Based on the distribution of the virtual bypass paths in Figure 5.1(a), the average number of virtual bypass paths in a routing path is 0.56, while, based on our distribution of the virtual bypass paths in Figure 5.2(a), the average number of a routing path is 1.13.

In our EVC-based power gating approach, routers always try to send packets to a virtual bypass path. Only when there is no virtual bypass path available, the packets are sent along the normal path between routers.

#### 5.5.2 Extended Router Structure

We have extended the basic EVC router in Figure 5.1(b) to enable and support our novel power gating scheme. As shown in Figure 5.2(b), one power control (ctrlr) unit is added in the router. Handshaking control signals WU (wakeup) and PG (power gating) are added between routers. Compared with the conventional power gating, introduced in Section 2.3, extra handshaking control signals,  $WU_{EVC}$  and  $PG_{EVC}$ are added between the source router and the sink router for a virtual bypass path. In each input port, one direct link is added (e.g., the red arrow in Input port 0, shown in Figure 5.2(b)). These direct links are used to build the bypasses in the direction from X+ to X-, X- to X+, Y+ to Y-, and Y+ to Y-. To avoid N-packets to be blocked by the powered-off routers, in our EVC based power gating approach, the EVC latch is also used to hold N-packets when the router is powered-off or **being charged**. When a router is powered off and the EVC latch in an input port is used to hold a N-packet, a bypass path is setup by using the direct link in the input port and the crossbar for E-packets. For example, when a router is powered-off and the EVC latch in the X+ input port holds a N-packet, a bypass path from X+ to X- is built by using the direct link in the X+ input port and the crossbar in this router for E-packets. Then, if an E-packet is coming, it directly goes through this router by taking this directly built bypass in the router. In this way, we guarantee that the virtual bypass path always works for E-packets even when the EVC-latch is occupied by a N-packet. Furthermore, the powered-off router has certain transmission ability to transfer N-packets through the normal paths. In this way, the N-packets have less probability to be blocked by powered-off routers.

To transfer N-packets though a powered-off router, the routing computation unit, the EVC latches, the virtual channel allocator unit, the switch allocator unit, and the crossbar are always powered on to execute the router pipeline stages. The power control (ctrlr) unit only cuts off the power supply of VCs. In this way, even at the powered-off state, the router still keeps a certain ability to transfer packets. Thus, the packets going through the normal paths have less probability to be blocked by the powered-off routers. Furthermore, as these units consume much less power than VCs [WNWS17, ZOG<sup>+</sup>15], our EVC-based power gating approach still can efficiently reduce the static power consumption by powering off the idle VCs.

#### 5.5.3 Power Gating Scheme

In this section, we introduce the conditions which drive our ctrlr unit in Figure 5.2(b) to control the power supply of a router.

#### Powering off a router

When there are no packets left in EVC latches, N-VCs, E-VCs, or the crossbar in a router, and the WU and  $WU_{EVC}$  signals from all its upstream routers are de-asserted, the router goes into the idle state, the  $PG_{EVC}$  and PG signals are asserted to all upstream routers, but at this moment, the power supply is not cut off yet. After waiting  $T_{idle\_detect}$  clock cycles, the ctrlr unit asserts the sleep signal (Figure 5.2(b)) and cuts off the power supply. If there is any WU or  $WU_{EVC}$  signals asserted during  $T_{idle\_detect}$ , the ctrlr unit immediately de-asserts the  $PG_{EVC}$  and PG signals. By delaying  $T_{idle\_detect}$  clock cycles to cut off the power supply, we can avoid non-beneficial power gating caused by short idle time of routers, which causes frequent power gating and additional power consumption.

#### Powering on a router

The process of powering on a router in our EVC-based power gating approach is an extension of the wakeup process shown in Figure 2.9 explained in Section 2.3. If a source router determines that a packet should take the virtual bypass path to the sink router, this source router asserts the corresponding  $WU_{EVC}$  to power on the sink router. If a router determines that a packet should take the normal path to the downstream router, this router asserts WU to power on the downstream router. Once the powered-off router receives the  $WU_{EVC}$  signal or the WU signal, the powered-off router starts to charge and goes into the wakeup state. After  $T_{wakeup} - MARGIN_{EVC}$  clock cycles, the router de-asserts  $PG_{EVC}$  and the source router can send packets to this router using the virtual bypass path. After  $T_{wakeup} - MARGIN$  clock cycles, the router de-asserts PG and the upstream router can send packets to this router using the normal path. By setting properly  $MARGIN_{EVC}$  and MARGIN, a router can send packets before the powered-off router is fully charged, but it is guaranteed that when a packet reaches the powered-off router, this router is just fully charged. In this way, we can hide part of the wakeup delay and optimize the power gating process. It should be noted that  $MARGIN_{EVC}$  is larger than MARGIN. This is because by taking virtual bypass paths, E-packets have more time on the transmission via multiple hops than N-packets taking the normal path to transfer over a single hop. This implies that the wakeup delay has less negative impact on the virtual bypass paths. Thus, it is more beneficial for packets to take the virtual bypass paths to avoid the negative impact of the power gating.

#### 5.5.4 Example of Our Power Gating Approach

In this section, we use the example shown in Figure 5.3 to clearly illustrate our EVCbased power gating approach.

In Figure 5.3(a), at time T = 0, Router0 and Router1 are powered-on and Router2 and Router3 are powered-off. Router0 is going to send an E-packet (the red blocks in Figure 5.3) to Router3 by using the virtual bypass path, so Router0 asserts the  $WU_{EVC}$  signal to wakeup Router3. Router1 is going to send one packet to Router3, but there is no virtual bypass path available, so Router1 treats this packet as a N-packet (the blue blocks in Figure 5.3) and sends it by using the normal path to Router2 first. So, Router1 has to asserts the WU signal to wakeup Router2.

At time T = 1, Router2 and Router3 receive the WU and  $WU_{EVC}$  and begin to power on, respectively. At time T = 0, 1, 2, 3, Router1 executes the router pipeline stages for its N-packet. The head flit of the N-packet leaves Router1 at time T = 3. At time T = 4, this head flit is going through the link, as shown in Figure 5.3(b). At time T = 2, Router2 and Router3 de-asserts the  $PG_{EVC}$  signals, but the E-packet



Figure 5.3: An example of our power gating approach.

is still blocked for one clock cycle at *Router*0. So, at time T = 4 (Figure 5.3(b)), the E-packet has not been sent yet.

In Figure 5.3(c), at time T = 5, the head flit of the N-packet reaches *Router2* and *Router2* holds this head flit at its EVC latch. At the same time, in *Router2*, one bypass path is setup by using the direct link and the crossbar. The head flit of the E-packet leaves *Router0* and is traversing the link.

In Figure 5.3(d), at time T = 6, as *Router*2 has to execute the router pipeline stages for the N-packet. The head flit of the N-packets has to occupy the EVC latch for multiple clock cycles. For the E-packet, the head flit reaches *Router*1 and is held at the EVC latch. The tail flit of the E-packet also leaves *Router*0.

In Figure 5.3(e), at time T = 7, the head flit of the E-packet leaves *Router*1 and the tail flit of the E-packet is held at the EVC latch of *Router*1.

In Figure 5.3(f), at time T = 8, the head flit of the E-packet directly goes through the directly built bypass path in *Router*2, and is traversing the link from *Router*2 to *Router*3. The tail flit of the E-packet is traversing the link from *Router*1 to *Router*2.

In Figure 5.3(g), at time T = 9, the head flit of the N-packet leaves *Router2* and the bypass path in *Router2* is demolished. For the E-packet, the head flit reaches its destination *Router3*. *Router3* is just fully charged and stores this flit into the allocated E-VC. The tail flit of the E-packet is held at the EVC latch in *Router2*.

In Figure 5.3(h), at time T = 10, the head flit of the N-packet is stored in Router3

and the tail flit of this N-packet is stored in *Router2*. As *Router2* and *Router3* are already full charged. These flits are stored in the corresponding N-VCs.

This example clearly shows that, by temporarily holding the packets in the EVC latches, the powered-off/ being charged routers can keep certain transmission ability to transfer N-packets. Thus, the N-packet can avoid as much as possible to be blocked by the powered-off/being charged routers. Furthermore, this process does not block the virtual bypass paths at all.

#### 5.5.5 Resolving Starvation

Starvation is a common issue in EVC-based NoCs [KPKJ07]. When an E-packet goes through an intermediate router along one virtual bypass path, the E-packet has the highest priority and the intermediate router has to send it first. If the source router continuously transfers E-packets through the virtual bypass path, the N-packets in the intermediate router cannot get a chance to be sent and starvation occurs. In order to resolve the starvation, we use the approach provided in [KPKJ07] to detect the starvation and then we use our novel approach to temporarily freeze the related virtual bypass paths. For example, considering Figure 5.2(a), if Router01 continuously sends E-packets to Router04 or Router02 continuously sends E-packets to Router05, Router03 cannot send packets to its downstream Router04. Once such starvation occurs, *Router*03 needs to freeze these two virtual bypass paths. To simplify the control between routers, we use two different ways to freeze these two virtual bypass paths: 1) To freeze the virtual bypass path from Router01 to Router04, Router03 informs the sink Router04 to assert  $PG_{EVC}$  in the direction X-. In this way, Router01 cannot send E-packets to Router04; 2) At the same time, to freeze the virtual bypass path from Router02 to Router05, Router03 informs the source Router02 to stop allocating E-VCs in the X+ direction to packets. In this way, *Router*02 cannot send E-packets to *Router*05 and the virtual bypass path is freezed. Thus, as all the virtual bypass paths through Router03 are freezed, no E-packets prevent *Router*03 to send its packets, thereby resolving the starvation. When the packets, initially affected by the starvation, leave Router03, then Router03 informs Router04 to de-assert the  $PG_{EVC}$  signal as well as Router03 allows Router02 to allocate E-VCs to packets. In this way, the frozen virtual bypass paths are activated and can be used again.

# **5.6 Experimental Results**

In order to evaluate our EVC-based power gating approach in terms of performance and power consumption, we have implemented our approach using the full-system simulator called Agate [CZPP16]. Agate is based on the widely used full-system

Network topology	$8 \times 8$ mesh
Router	4-stage pipeline
Virtual channel	(1 N-VC, 1 E-VC)/VN, 3 VNs,
Input buffer size	1-flit/ ctrl VC, 5-flit / data VC
Routing algorithm	X-Y DoR
Link bandwidth/delay	128 bits/cycle, 1 clock cycle
Wakeup delay	8 clock cycles
Break even time	10 clock cycles
$T_{idle\_detect}$	8 clock cycles
$MARGIN_{EVC}$ / $MARGIN$	6/4 clock cycles
Private I/D L1\$	32 KB
Shared L2 per bank	256 KB
Cache block size	16 Bytes
Coherence protocol	Two-level MESI
Memory controllers	4, located one at each corner

Table 5.1: Parameters used in experiments.

simulator GEM5 [BBB<sup>+</sup>11] and Agate supports the simulation of the key items in NoC power gating techniques. The NoC model and power model used in Agate are based on Garnet [AKPJ09] and Dsent [SCK<sup>+</sup>12], respectively. The key parameters used in our experiments are shown in Table 5.1. We choose a four-stage pipeline router. There are three virtual networks (VNs): two data VNs and one control VN. In each input port, there is one N-VC and one E-VC for each VN. The value of the wakeup delay and break even time (BET) are set according to the related works [CZPP15] and [CP12]. Based on the NoC configuration, we set  $T_{idle\_detect}$ ,  $MARGIN_{EVC}$ , and MARGIN such that we keep the correctness of the NoC.

For comparison purpose, we have implemented the following power gating approaches: (1) NO\_PG: the baseline NoC without power gating; (2) Conv\_PG: conventional power-gating NoC, which is deeply optimized by sending WU (Look ahead [MKWA08]) and de-asserting PG signals [CZPP16] in advance, thus 6 clock cycles of the wakeup delay are hidden in our experiments; (3) DB\_PG [WNWS17]: the NoC with our DB-based power gating approach introduced in Chapter 3. In each input port of a router, a one-flit size duty buffer is added to implement the Duty Buffer approach. (4) EZ-bypass [ZL18]: the power gating process. Compared with other bypass-based related approaches [CP12, FTKH16, BHW<sup>+</sup>17], EZ-bypass is more flexible to allow packets to bypass the powered-off routers. (5) D\_bypass: the NoC with our D-bypass power gating approach introduced in Chapter 4. (6) EVC\_PG: the NoC with our EVC-based power gating approach.



Figure 5.4: Latency across different injection rates.

#### 5.6.1 Evaluation on Synthetic Workloads

In order to explore the behaviour of our EVC\_PG, in this section, we evaluate the performance and power consumption of our EVC\_PG approach under synthetic traffic patterns. We select three synthetic traffic patterns: 1) Uniform random: packets' destinations are randomly selected; 2) Bit-complement: packets from source router (x, y) are sent to destination router (N-x, N-y), N is the number of routers in the X and Y dimensions of a NoC; 3) Transpose: packets from source router (x, y) are sent to destination router (y, x).

#### Effect on NoC Network Latency

Figure 5.4 shows the average packet latency under different injection rates. Compared with NO\_PG, Conv\_PG, DB\_PG, EZ-bypass, and D\_bypass, our EVC\_PG has the lowest average packet latency. These results indicate that our EVC\_PG can effectively reduce the negative impact of the wakeup delay and can be used to achieve low latency communication. On the other hand, our EVC\_PG has lower saturation points than NO\_PG, Conv\_PG, EZ-bypass, and D\_bypass for the Uniform random and Transpose patterns, but has higher saturation point for the Bit-complement pattern. The lower saturation points indicate that our EVC\_PG causes some throughput loss. This is because, in order to support the EVC scheme, the VCs in our EVC\_PG are partitioned into E-VCs and N-VCs, which may undermine the flexibility and effectiveness of VCs. Since, Conv\_PG, EZ-bypass, and D\_bypass are based on NO\_PG, they have the same saturation points as NO\_PG. However, the impact caused by the partition of E-VCs and N-VCs highly depends on the traffic pattern. Thus, for Bit-complement, our EVC\_PG achieves higher saturation point.



Figure 5.5: Power consumption across different injection rates.

#### **Effect on NoC Power Consumption**

Figure 5.5 shows the power consumption normalized to NO\_PG under different injection rates. When the injection rate is around 0.001 packets/node/cycle, our EVC\_PG has slightly higher power consumption than Conv\_PG, EZ-bypass, and D\_bypass, but much lower than NO\_PG. This is because, in order to avoid packets to be blocked by powered-off routers, we always keep some components powered on in the powered-off routers, which causes extra power consumption but this power consumption is rather low. When the injection rate increases, more and more routers become busy and cannot be powered off. The power reduction in Conv\_PG, DB\_PG, EZ\_PG, D\_bypass, and EVC\_PG becomes lower and lower, but DB\_PG has much higher power reduction than the other approaches. This is because DB\_PG can separately power off VCs in each input port of routers whereas Conv\_PG, EZ-bypass, D\_bypass, and EVC\_PG can power off a router only when all of the input ports of the router are idle. Thus, DB\_PG fully utilizes the idle time of each input port to reduce the power consumption.

When the injection rate is higher than 0.02 packets/node/cycle in Figure 5.5(a), 0.02 packets/node/cycle in Figure 5.5(b), and 0.03 packets/node/cycle in Figure 5.5(c), Conv\_PG and EZ-bypass become ineffective in reducing the power consumption, while DB\_PG, D\_bypass, EVC\_PG still can effectively reduce the power consumption. The power reduction in our EVC\_PG is due to the fact that packets can also bypass powered-on routers, which saves some dynamic power.

When the injection rate further increases, the dynamic power takes higher and higher portion of the total power consumption. Our EVC\_PG reduces more the dynamic power consumption, which causes the curves for our EVC\_PG in Figure 5.5(a), Figure 5.5(b), and Figure 5.5(c) to decline. As a result, when the injection rates are higher than 0.04 packets/node/cycle in Figure 5.5(a), 0.03 packets/node/cycle in Figure 5.5(b), and 0.04 packets/node/cycle in Figure 5.5(c), our EVC\_PG consumes less power than D\_bypass. When the injection rates are higher than 0.07 packets/node/cycle



Figure 5.6: Execution time.

in Figure 5.5(a) and 0.05 packets/node/cycle in Figure 5.5(b), our EVC\_PG has lower power consumption than DB\_PG. However, in Figure 5.5(c), DB\_PG has always lower power consumption than our EVC\_PG. This is because DB\_PG and EVC\_PG reach their saturation points at low packet injection rates as shown in Figure 5.4(c). So, the dynamic power consumption takes small portion of the total power consumption. As a consequence, the efficient reduction of the dynamic power consumption in our EVC\_PG does not play a significant role in reducing the total power consumption in this case, whereas DB\_PG more efficiently reduces the static power consumption by separately powering off input ports of routers, leading to better reduction of the total power consumption in this case.

#### 5.6.2 Evaluation on Real Application Workloads

In this section, we use real application workloads to compare the approaches in terms of the application performance, the average network latency, and the NoC power consumption. To do so, we use nine applications from the Parsec [BKSL08] benchmark suite.

#### **Effect on Application Performance**

Figure 5.6 shows the total execution time of the nine applications, which is normalized to the baseline NO\_PG, and the tenth set of bars in Figure 5.6 gives the average results over these nine applications. Our EVC\_PG approach causes less performance penalty (execution time increase) than the related approaches. Compared with the baseline NO\_PG, our EVC\_PG causes, on average, 2.67% performance penalty, which is less than the 28.67% performance penalty in Conv\_PG, 7.24% in DB\_PG, and 5.69% in EZ-bypass, and comparable with the 2.55% performance penalty in



Figure 5.7: Average network latency.

D\_bypass. For benchmarks blackscholes and x264, our EVC\_PG has slightly lower execution time than NO\_PG. In the vips benchmark, our EVC\_PG has its highest performance penalty of 6.17%, which is still lower compared to Conv\_PG, DB\_PG, and EZ-bypass, but higher than D\_bypass. For the ferret benchmark, Conv\_PG, DB\_PG, EZ-bypass and D\_bypass have their highest performance penalty of 47.39%, 21.21%, 19.51%, and 6.03%, respectively.

#### Effect on Average Network Latency

Figure 5.7 shows the average network latency across the nine applications. Compared with NO\_PG across the applications, the average network latency in our EVC\_PG approach is slightly lower, whereas Conv\_PG, DB\_PG, EZ-bypass, and D\_bypass have higher average network latency compared to NO\_PG. As DB\_PG uses a fined-grain power gating scheme, packets in DB\_PG suffer more wake up processes. As a consequence, DB\_PG has much higher average network latency than our EVC\_PG and EZ-bypass. EZ-bypass allows packets to bypass powered-off routers, but packets have to stay at powered-off routers for a long time experiencing the router pipeline stages. In contrast, in our EVC\_PG, the packets can bypass the intermediate routers without the need to experience the router pipeline stages. Thus, our EVC\_PG has lower average network latency than EZ\_PG. While, compared with D\_bypass, which needs extra time to reserve the bypass latch in the powered-off routers, our EVC\_PG is more efficient to transfer packet to bypass the powered-off routers.

Even though our EVC\_PG has a slightly lower average network latency compared to NO\_PG (see Figure 5.7), our EVC\_PG still causes a slightly higher execution time in most of the applications compared to NO\_PG (see Figure 5.6). This is because EVC\_PG causes unbalance NoC resource allocation when E-packets take the virtual bypass paths to bypass intermediate routers and have a higher priority compared to N-packets.



Figure 5.8: Power consumption.

#### **Effect on NoC Power Consumption**

Figure 5.8 shows the breakdown of the NoC power consumption across the nine applications and the tenth set of bars shows the average over these nine applications. The NoC power is broken down into three parts: the power consumption caused by power gating (PG\_overhead) and the static/dynamic power consumption of routers (static/dynamic).

As shown in Figure 5.8, our EVC\_PG approach consumes slightly higher total power than the related approaches Conv\_PG, DB\_PG, EZ\_PG, and D\_bypass. This is because our EVC\_PG needs more components in a router to be always powered on, which causes slightly more static power consumption compared to Conv\_PG, DB\_PG, EZ\_PG, and D\_bypass. As the traffic workloads in real applications are low, the dynamic power consumption is low. As a result, the dynamic power reduction in our EVC\_PG does not play a significant role in reducing the total power consumption. Compared with NO\_PG, our EVC\_PG reduces on average 68.29% of the total power consumption, which is comparable with the 72.94%, 73.56%, 75.30%, and 77.77% reduction of the total power consumption in Conv\_PG, DB\_PG, EZ-bypass, and D\_bypass, respectively.

# 5.7 Discussion

In this chapter, we propose an EVC-based power gating approach. In our approach, packets can take pre-defined virtual bypass paths to bypass intermediate routers that may be powered-on or powered-off. Furthermore, even though some packets do not take a virtual bypass path, our approach tries to ensure that these packets avoid as much as possible blocking in the powered-off routers. As a result, our approach reduces more efficiently the packet latency increase caused by power gating. Furthermore, by allowing packets to bypass powered-on routers to reduce dynamic power

consumption, our approach can achieve lower power consumption under high traffic workloads.