# Investigating and Analyzing the Effect of Router Components on Network Performance on the Chip with Regard to Power Consumption

Farnaz Zogh[1], Azadeh Alsadat Emrani Zarandi[2], Vahid. Sattari Naeini[3]

1- Department of Computer Engineering, Shahid Bahonar University of Kerman, Kerman, Iran.
Email: farnaz.zogh@eng.uk.ac.ir
2- Department of Computer Engineering, Shahid Bahonar University of Kerman, Kerman, Iran
Email: a.emrani@uk.ac.ir
3- Department of Computer Engineering, Shahid Bahonar University of Kerman, Kerman, Iran
Email: vsnaeini@uk.ac.ir

**ABSTRACT:**
A network on the chip is a solution to connection problems compared with traditional-based chip which can fulfil multi-dimensional communication requirements. The router is a key component of the communication network which is referred as its backbone. Since the router occupies the largest area on the chip and it is the most widely used network component, in this paper, the architecture of the router in the network on the chip is examined and its roles with its components and their effect on the performance of the network are investigated, considering the parameters including time (delay), the area and more importantly the power consumption. It is shown that any modification, combination, or correction in any of the component effect on power consumption of the router and hence on the power consumption of the whole chip. To this end, the related works are examined to make an appropriate estimation of their analogy. This article will help researchers who are trying to design an optimal router based on power consumption.

**KEYWORDS:** Network on the Chip, Router Architecture, Power Consumption, Area.

## 1. INTRODUCTION

In the network on the chip, as the number of cores increases, the links needed to communicate also increase exponentially, which results in routing problems and areas. Simply, the components used on the chip are network adapters, routers, and links. Routing nodes route data based on the selected protocols, and Links are used to connect nodes based on the required bandwidth. Links can include one or more logical or physical channels. The algorithm used for routing is one of the important metrics in the network on chips and its efficient performance. The routers duty is sending data packets through the links from one processing unit to another [1]. The architecture of any router is composed of some components, which choosing useful components can lead to performance improvement of the whole network. Therefore, our aim here is investigating available options for each of these components and its impact on the router's performance and also the entire network on the chip.

In the rest of the paper, first the routing components and their impact on network performance are discussed. In the Section 2, router components and roles of each one are introduced. Section 3 briefly describes the role of selecting component types in the router and, the main methods to reduce power consumption of the network on chip are presented in Section 4. Finally, the conclusions are drawn from the topics discussed.

## 2. THE ROUTER COMPONENTS

Each router consists of different parts for storing, deciding, and transmitting data packets. Among distinct types of memory and registers, mainly buffers are used as information storage sections in the router. It is obvious that the more storage space is considered, the more area and power is consumed in the router and also the whole chip, on the other hand the reliability increases and the packet loss and transmission delay are reduced. Accessing this storage space and the targeted guidance of data packets requires decision-making in control units by using a variety of routing algorithms, arbitration and decision making in these units. The other part is the data transmitter used by the two previous partitions to transfer packets. The architecture of a noc router can be seen in Figure 1.
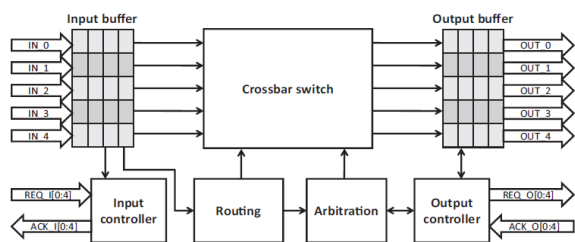
**Fig .1**. The Basic noc router architecture [2].

The network router on the chip usually has five bi-directional ports for transmitting packets across the network on the chip, the most commonly used names are routing ports: North, South, East, West and the local port. The router, block is the main structure of a grid on a chip, and should be available on each chunk node on the grid. A grid on a two-dimensional chip with mesh connectivity is shown in Fig. 2. As it can be seen in the figure, the local port is a connection between the router and the PE[1] network processing unit on its own chip, While other ports are the connection between neighboring routers.
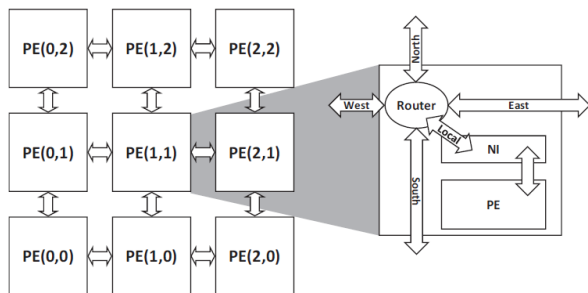


**Fig. 2**. Network router on chip, PE and NI in mesh connectivity [2].

Generally, each sub-block in the router performs one of the following main tasks:
1. Crossbar switch: This component connects all input ports with all output ports along a multiplexer, which its select pins are based on the routing, decision making and arbitration units [3].
2. Input / Output buffers: Flits which comes from the network or local node are entered to the input buffers. Also, these flits are saved in output buffers after they are processed by the router and before they are sent to the network [3].
3. Input / Output Controller: The implementation of manual shaking or controlling for the required data flow at the flit level between neighboring routers is the duty of these controllers to prevent overload or low load in the input or output buffers.

4. Routing Unit: The static or dynamic routing algorithm which selects the output port in order to send packets is implemented by this section.
5. Arbiter unit: Controls the connection of input ports to outputs and also, when multiple ports request for routing algorithms to access the same port, preventing bottlenecks and congestions by prioritizing requests that have already been made.

## 3. INVESTIGATING THE EFFECT OF COMPONENT SELECTION ON ROUTER PERFORMANCE

In this part of the article, the router performance is analyzed by selecting the different components. As the first part, switch is considered. Each switch is composed of two parts: connections and arbitration. In the crossbar section, each of the input port are connected to all the output port of the switch, but arbitration makes these connections available. In fact, all connections can be connected provided that there is enough space in the buffer Unless several ports simultaneously request to connect a same links at the same time, the arbitration is used to implement algorithm to manage or prioritize the sending of packets. shows the RTL view with using the VHDL code from the crossbar to the multi selector.
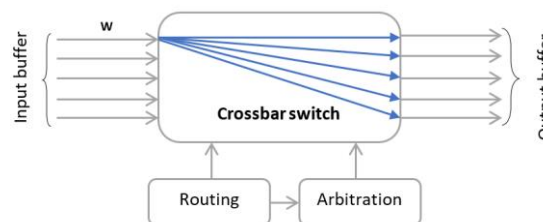


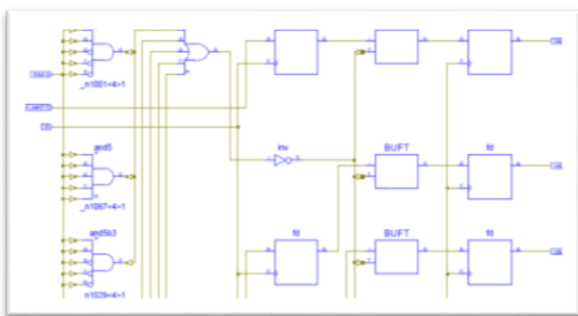**Fig. 3**. Switch and related communication.

Connections based on applications and design requirements can be directional or bidirectional. All possible communication paths for the input port W are specified in Figure 3. One approach to design a switch is to implement a crossbar (matrix) connection, which, based on the generated signals by allocated to theoutput port, connects the output of each input queue to each output queue. There are two common implementations of the crossbar switches, using ass transistors or multiplexers. The implementation of pass transistor results in a decrease of area and power consumption, While the implementation of the multiplexer has the advantage that it can be automatically synthesized from a RTL description and also used in FPGA devices as it is shown in Figure 5. In theory, Crossbar's area grows with the square of the ports number and, in fact, automatic transfer and routing using the EDA[2] tool results in the overheads of more than 32 ports. A t code assignment for 5-port connections is shown in the table1.

---

[1] Processing Element

[2] Electronic design automation

**Table 1.** code assignment for switch ports.

| port | Output<=input | RA |
|---|---|---|
| Output_North | O(0) <= I(0) | 00000 |
| | O(0) <= I(1) | 00001 |
| | O(0) <= I(2) | 00010 |
| | O(0) <= I(3) | 00011 |
| | O(0) <= I(4) | 00100 |



**Fig. 4**. The RTL schema of the switch.

The common approaches are based on using input queue, or output queue, or both of them. However, there are newer approaches, using virtual channels [4], which uses shared queue sharing for all channels. In order to design virtual channels, the multiplexer and demultiplexer are required in the router's input port as it is shown in Figure 1. Multiplexers move the message queue to fit the virtual channel within the packet. This form of buffer implementation, while increasing the number of virtual channels, improves the performance. Increasing the number of buffers for speed results in unacceptable power consumption. Therefore, the buffer sharing between virtual channels is proposed.

It can be said that the buffer size is the number of flites stored in the router. Buffers are costly in area and power consumption, and typically, few flits are used. The number of stored flutes in the buffer is very important, since the higher the number of flit become, the more input or output data can be saved and, this results in lower packet loss and improves the operational capability, including the transmission speed. Although the network on a chip with few buffers will be inefficient in heavy traffic conditions, but in general assigning huge memory for buffer is not a good design since memory enlargement leads to more power and space occupancy on the chip. It should be taken into account that the router is the most common element of every NOCs, and therefore, a small change in the area or power of this over-consuming element has a serious effect on the total area and power consumption of the whole grid on the chip [5, 6]. In order to minimize the implementation cost, the area overhead should be low. Hence, unlike the systems on the chip which have large memory, often DRAM, networks on chip consist of small registers for buffering. Another advantage of using registers instead of large memory is that the latency of the encryption and decryption of addresses and the access delay will be significantly reduced. This is crucial for the time restricted applications, like real time ones, For some heterogeneous NOCs, routers are designed as a re-usable IP block [7]. In such cases, this block is intended to play the role of a specific IP to support the main processor whenever it is needed, beside acting as a router. This reconfigurabilty feature can be also considered in self-buffering which means the ability to change the way of using buffers. The main difference between a flexible buffer and the common ones is using pipelining for flip flops to buffer data instead of adding input buffers to router ports. Usually the word length is the same for every ports. However, in [8, 9], the channel length varies with the router's ability to be reconfigurable.

Buffers is an element with high power consumption, therefore, if a buffer is designed efficiently, it results in a good trade-off between area and power. Buffer size is important since it effects the area occupied by the router on the chip and the maximum power leakage is on the network router on the chip, it consumes about 64% of the total leakage power [11]. In addition, they can significantly increase the dynamic power consumption. Previous studies show that storing a data in a buffer consumes more energy than transmitting it. However, reducing the size of the buffer to decrease the power consumption and the area of silicon is often not a good solution, since, it degrade performance, especially under heavy traffic. Due to effects of buffer on power consumption, distinct researches have been done on it. The arrangement of buffers was limited to the input and output queues to store inbound and outbound traffic. Each input / output port has its own dedicated buffers. In regular topologies, all routers and ports on a router usually have the same buffer size. The optimal buffer size and the number of virtual channels are still needed to be specified. Although in unconventional architectures, each router and port buffer size can be set according to the application's traffic needs.

Communication signals are messages which contain Ack or Nack, exchanged between different routing sections, or the response given to the conditional blocks in a router's flowchart. The drawback of flow control based on acknowledges or nack is its need for additional buffers for support, which needs to have a copy of the transferred flit for the situation that needs to be retransmitted so that it can re-send. Sophisticated and smarter routing algorithms also increase communication

signals. Often, these signals do not require more than one or two bits but the sum of the power used by millions or billions of signals in a system with a large number of cores causes considerable power dissipation. Additionally, the production and transmission of each of these signals, which the performance of the other units depends on their proper arrival, It takes time, so more number of signals means more delay. Although designers are bound to use these signals, especially in decision-making units, it may be possible to use shared signals.

Among the performance parameters for routers, clock frequency and throughput can be mentioned. In general, it is not possible to increase the clock frequency by adding pipeline levels, although this is common in other digital systems, and that's why the low latency in the network is very important on the chip. In some work [10], they regulate the clock frequency locally and globally, which has a significant impact on reducing power consumption.

## 4. POWER REDUCTION SOLUTIONS IN NOC

Considering the importance of power consumption for NOC, several solutions have been proposed to reduce consumption, some of which are: modifying the structure of routers, reducing the number of routers in the networks by modifying their topology, reducing the capacitive capacities of links by replacing shorter links between operating blocks by further exchanging data, reducing the number of routers between operational blocks, modifying routing algorithms to select shorter ways with less delay, and reducing the busy switching activity using appropriate codes. As an example, CVBPM[3] can be referred to as a communication architecture for power management that It tries to address all of above to reduce power consumption. This method uses the DVS[4] technique, It adjusts the clock and voltage of processors on the chip. In this technique, the power of a link at various frequencies can be changed up to ten times. This is done by providing an algorithm that adjusts the feed and frequency of the links based on the bandwidth and the instantaneous frequency. It improves power consumption by up to 3.2 times. Of course, reliability is not considered in this study. If the data in the receiver has reached with error, there are two ways to handle it. The first method is FEC[5] that uses the properties of the code and corrects the error and does not request redundancy at all. This method uses less bandwidth of links and routers. Instead, it needs a complex decoder. The second method is ARQ[6],

requesting retransmission by the receiver. This technique requires a separate channel from the receiver to the sender to request a retransmission. The two methods of FEC and ARQ and their combination in wireless communications are compared; In the same conditions, ARQ uses less power than FEC. In addition, the FEC consumes more energy, while the ARQ transfers more bits. According to these comparisons, the selection of any suitable application technique can be effective on power consumption [11].

There is another technique called DSCDMA[7] that has been recently used in channels. Each sender modulates the information with codes in the way that it becomes orthogonal to the order of the codes of the other senders. These codes are perpendicular to each other with their own, and conversely, they have little relation with the order of other codes. The receiver can even retrieve information even if there are different codes from different senders. This allows to send multiple simultaneous data to a common physical media, in contrast with typical architectures that only one component can send information. Another technique is using WLAN[8] connection. Using CDMA[9] and FDMA[1] , multiple components can send and receive information together. By RF[1] communication, up to 50 data can be sent in parallel, which is due to the high RF bandwidth [11]. OrionI and OrionII are two important papers about modeling the power consumption of networks on the chip that are widely used in noc simulators [12]. In these articles, precision models are presented at the transistor level. The results from the modeling were compared with the results of practical experiments on the Alpha router from IBM, and a high correlation was found between the two(them). In 2012, a model that included static power consumption, in addition to the dynamic power consumption provided by OrionI, was named OrionII. In this paper, the total power of the following equation is calculated.

$$Pt = Pstatic + Pdyn \qquad (1)$$

In the above, Pt is the total power and Pstatic is the static power that is calculated for each gate, such as NAND and NOR in HSPICE, and is added to the model and Pdyn is dynamic power. The overall power that is calculated by the ISIM tool in the FPGA is the sum of the two components, as described by Xilinx corporation. Static power is mainly due to leakage current in transistors of the device. There is also a current leakage from the source to the drain or through the gate oxide,

---

[3] Communication Based Power Management
[4] Dynamic Voltage Scaling
[5] Forward error correction
[6] Automatic Repeat Request
[7] Direct Sequence Code Division Multiple Access

[8] Wireless Local Area Network
[9] Code-Division Multiple Access
[1] Frequency-Division Multiple Access
[1] Radio Frequency

and even when the transistor is logically "off". However, the dynamic power, which consists of the activity of the scheme and the event of switching at the core or the device's input / output, is determined by the capacitive capacitance of the node, the power supply voltage, and the switching frequency (CV2f) [13]. One of the important results of OrionII is the estimation of power consumption by component, as shown in Fig. 5, Where the highest power consumption is 33% at the clocking and 22% in the router's FIFO buffers. This model has been used in various articles in various applications.
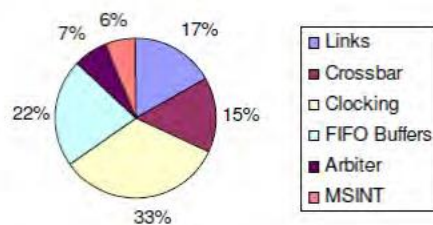


**Fig. 5**. Percentage of the power consumption of each sub-circuit from the router in the grid on the chip relative to the total power consumption of the router [11].

## 5. CONCLUSION

The network on the chip is introduced as a suitable alternative for connection problems of chip systems and the router is one of its most important and most common network elements on the chip, which its power consumption has a significant effect on network performance. For this aim, in this work, router components are described and studied and their power consumption are analyzed. The effect of each individual components on the power consumption of routers and the resulting on the whole network were studied in this work. Using the ability to reconfigure and optimize design for routing buffers, using efficient algorithms and proper mapping were among the most effective factors on power consumption. As noted, a small change on the component design significantly affects the power consumption of the entire router and also NOC. Therefore, it is suggested that special attention should be paid to NOC's routers and also designing its components with high-performance should be considered.

## REFERENCES

[1] R. Afshar Mazayjani, M. Alaee, and F. Yazdanpanah, "**Composition of Wired and Wireless Communications in Networks on Chip**", *International Conference on Computer Engineering and Information Technology,* 1395.

[2] K. Tatas, K. Siozios, D. Soudris, and A. Jantsch, "**Designing 2D and 3D Network-on-Chip Architectures**" *Springer,* 2014.

[3] H. Elmiligi, M. Sallam, and M. W. El-Kharashi, "**A Power-Optimized, Area-Efficient Implementation of Connection-Then-Credit NoC physical layer**" *Microelectronics Journal,* Vol. 68, pp. 55-68, 2017.

[4] N. Kavaldjiev, G. J. M. Smit, and P. G. Jansen, "**A Virtual Channel Router for on-Chip Networks**" in *SOC Conference, 2004. Proceedings. IEEE International*, 2004, pp. 289-293: IEEE.

[5] F.Zogh, M. Alaei, and F. Yazdanpanah "**Architecture of Reconfigurable Routers in NOC**" *the Fourth International Conference on New Findings of Science and Technology,* 1396.

[6] F. Yazdanpanah, M. Alaee, and F.Zogh, "**Dynamic Reconfigurable Network on Chip**" *the Fourth International Conference on New Findings of Science and Technology,* 1396.

[7] P. Mahr and C. Bobda, "**Reconfigurable router for dynamic Networks-on-Chip**" *Rapid System Prototyping (RSP), 2010 21st IEEE International Symposium on*, 2010, pp. 1-6: IEEE.

[8] C. Concatto, D. Matos, L. Carro, F. Kastensmidt, A. Susin, and M. Kreutz, "**Noc Power Optimization using A Reconfigurable Router,**" in *NoC Power Optimization Using a Reconfigurable Router*, 2009: IEEE.

[9] D. Matos, C. Concatto, M. Kreutz, F. Kastensmidt, L. Carro, and A. Susin, "**Reconfigurable Routers for Low Power and High Performance,**" *IEEE Transactions on very large scale integration (VLSI) systems,* Vol. 19, No. 11, pp. 2045-2057, 2011.

[10] R. Kamal and J. M. M. Arostegui, "**A Multi-Synchronous Bi-Directional NoC (MBiNoC) Architecture with Dynamic Self-Reconfigurable Channel for the GALS Infrastructure,**" *Alexandria Engineering Journal,* 2017.

[11] B. Naraqi, Gh. Karimi, "**Study on Power Consumption and Power Consumption Reduction Solutions in Network on Chip and Test Networks**", *the first National Conference on Computer Engineering, Computer Science and Information Technology, Qom, June 1395, Applied Scientific University Qom governorate.*

[12] A. B. Kahng, B. Li, L.-S. Peh, and K. Samadi, "**Orion 2.0: A Power-Area Simulator for Interconnection Networks,**" *IEEE Transactions on Very Large Scale Integration (VLSI) Systems,* Vol. 20, No. 1, pp. 191-196, 2012.

[13] Xilinx, Xilinx Power Tools Tutorial, Vol. UG733 (v1.0), March 15, 2010.