

# A case of area- and energy-efficient heterogeneous mesh network-on-chip

Jili Yan, Xiaola Lin

School of Information Science and Technology  
Sun Yat-sen University  
Guangzhou, 510006, China  
e-mail: [linxl@mail.sysu.edu.cn](mailto:linxl@mail.sysu.edu.cn)  
[yanjili@mail2.sysu.edu.cn](mailto:yanjili@mail2.sysu.edu.cn)

Guoming Lai

College of Computer Engineering and Technology  
Guangdong Institute of Science and Technology  
Zhuhai, 519090, China  
e-mail: [laigm@mail3.sysu.edu.cn](mailto:laigm@mail3.sysu.edu.cn)

**Abstract**—In this paper, based on observation on performance analysis of mesh network, we proposed a case of area- and energy-efficient heterogeneous mesh network by redistributing and reconfiguring scarce network resources, buffers, links and ports, to enhance network performance. Experimental results show that proposed network can achieve maximum saturation improvement by up to 16.7% and improve network latency by up to 35% while reduce about 31.7% router area. Experimental results also show that diagonal link is efficient design for mesh network topology.

**Keywords**- *network-on-chip; interconnection; system-on-chip; multiprocessor*

## I. INTRODUCTION

As Very Large Scale Integration (VLSI) technology advances, more cores can be integrated into a single chip to form a multiprocessor system-on-chip (MPSoC) [1]. Network-on-chip (NoC) architectures rather than traditional shared buses are becoming the de facto communication fabric for these SoCs [2] due to lower message latency, better scalability and more reliable performance predictability. Compared to off-chip networks or bus architectures, the design of network-on-chip faces with different design constraints, for example, power/energy consumption, area overhead and network latency, so designing NoC with area- and energy-efficient, low-latency and better scalability is therefore desirable [2,3].

Topology of network-on-chip, which is one of the research hotspots, has significant influence on performance of NoCs as whole [2]. Past decades, researchers have proposed numerous topologies for network-on-chip, such as, mesh, torus, and flattened butterfly etc. In these topologies, the mesh network is the most widely studied in research and employed in prototype systems due to its better scalability, regularity and ease of implementation in silicon. Recently, many variants of mesh had been proposed [5,6]. In these mesh variant, scarce resources like buffers and bandwidth are uniformly distributed in all nodes. While some researches show that network resource utilization in non-edge symmetric network like mesh network is not non-uniform under deterministic and adaptive routing [7]. The main ineffective of mesh interconnect fabric is the large network diameter and the artifact of non-edge symmetric network employing deterministic XY-routing. How to efficiently and fairly use the scarce resources while meeting area and energy constraints are design challenges. Figure 1

shows the buffer utilizations of  $4 \times 4$  and  $8 \times 8$  homogeneous mesh networks employing XY-routing under uniform random (UR) traffic with packet injection rate (*pir*) 0.4 and 0.3, respectively. The non-uniform of buffer utilization in mesh network is obvious.

In a mesh NoC, buffer, port and bandwidth of link are the scarce resources and have significantly influence on performance of system [14]. Some previous work explored different designs of mesh network by redistributing buffers and links or reconfiguring different number ports or bandwidth of links [5,6,7,8]. In addition, diagonal link design in chip manufacturing for the 2D mesh network has been proposed [15]. The diagonal link not only reduces inter-node distance but alleviates traffic congestion in the network so that network performance is enhanced dramatically. In this paper, we explored the design space of heterogeneous mesh network and proposed a case of area- and energy-efficient heterogeneous mesh, which achieves higher throughput, lower latency and consumes lower energy by redistributing and reconfiguring buffers, bandwidth of link and ports of router in mesh network.

The remainder of this paper is organized as follows: in Section II, we discuss related work; the design of area- and energy-efficient heterogeneous mesh network and routing algorithm are discussed in detail in Section III; Section IV shows the performance evaluation and benefits; conclusion is given in Section V.

## II. RELATED WORK

Topology defines how nodes are placed and connected, which affects the bisection bandwidth and the latency of a network, and many different topologies have been proposed, such as mesh, torus, fat-tree and so on. Mesh network topology in all network topologies draws special attention and is adopted in research and practical systems because of the scalability, regularity and ease of implementation. In recent years, researchers have also proposed many modified mesh topologies [5,6]. In the work [5], the X-network is employed to implement area-efficient router blocks, and one processing element (PE) is connected to four routers, and then one router also is connected to four PEs. In the work [6], diagonal links are added between switches. However, these mesh networks are homogeneous and non-edge symmetrical mesh networks, thus they are inherited shortcomings of mesh network.

Recent years, some heterogeneous mesh network topologies have been proposed in literatures [4, 7, 8]. Authors in [4] explore the heterogeneous mesh design for specific-application in order to improve the system

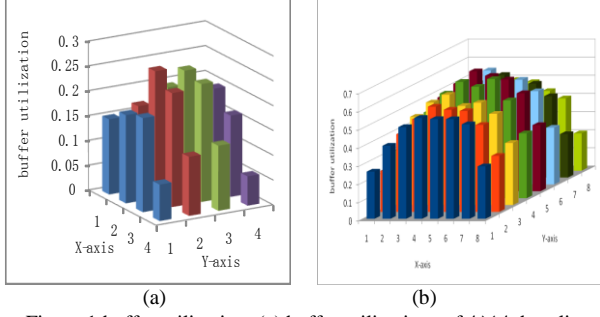


Figure 1 buffer utilization. (a) buffer utilizations of  $4 \times 4$  baseline mesh and (b) buffer utilizations of  $8 \times 8$  baseline mesh

performance by inserting long-range links. Work [7] proposed a case of heterogeneous mesh network for general CMPs. Authors use two kinds of routers with different performance and different bandwidth to improve the system performance. However, this heterogeneous mesh is different to scale and that has same network diameter as the traditional mesh networks. Kilo-NOC [8] uses also two kinds of routers, QoS-enabled and not QoS-enabled, to provide low cost, scalable and energy-efficient QoS guarantees in a network. Prior works have also investigated co-designing the NoC with caches [9] and memory controllers [10]. In particular, work in [9] examined heterogeneous wires with varying width, latency and energy, and proposed mapping coherence messages with differing latency and bandwidth characteristics onto the different wires. Work in [11] proposed two asymmetric networks, one customized for coherence and short messages and the other for cache bank reply packets. Most of these past works have investigated heterogeneity or customization in the network based on micro-architectural techniques or hardware characteristics. Our approach is different these work because our design provides a case of general area- and energy-efficient heterogeneous mesh network for a variant of traffic patterns.

### III. HETEROGENEOUS ARCHITECTURE DESIGN

#### A. design of basically heterogeneous mesh unit

Figure 2 (a) is a  $4 \times 4$  baseline mesh NoC, every router has 5 input and 5 output ports and every input port has a 8-flit buffer with 192 bit bandwidth. We start to design heterogeneous mesh network with the  $4 \times 4$  baseline mesh. For non-edge symmetric network, the aggregated properties of a variant of traffics in network is apparent. We design the heterogeneous mesh network by using two kinds of routers. One kind of multi-port router (MPR), has 7 input and 7 output ports and then the router has wider router bandwidth than that of baseline mesh network, every input port has a 6-flit buffer. Other kind, conventional-port router (CPR), has 5 input and 5 output ports like the router in the baseline mesh network with 8-flit buffer in every input port.

**Redistribution of links:** In our design, we keep the bisection bandwidth constraint with the baseline homogeneous mesh. With a baseline homogeneous mesh with 192b links, we have 128b links in heterogeneous mesh. As discussed in Section II, the diagonal link design is efficient improving network performance [15]. Figure 2 (b) shows the view of heterogeneous mesh with diagonal links.

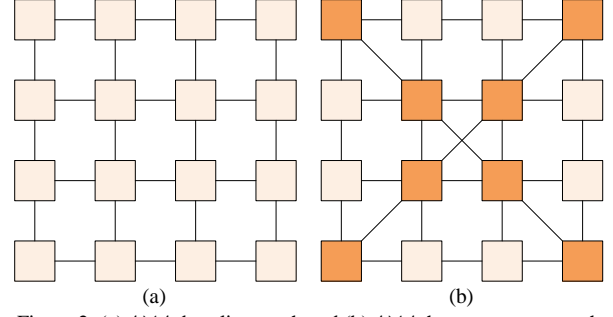


Figure 2 (a)  $4 \times 4$  baseline mesh and (b)  $4 \times 4$  heterogeneous mesh

This design, taking full advantage of abundant wiring tracks, is widely adopted in literatures, such as work [4, 5, 6]. In case of fig 2 (b), we can observe that there is higher bandwidth and buffer density in the center of heterogeneous mesh network than the peripheral area of network, which can better map the aggregated traffic of mesh network and shorten the network diameter.

**Redistribution of buffers:** since buffers consume about 35% of router power [7], having more buffers in a router increases the router power consumption. However, more buffers in a router improve performance. In our design, we configure 8-flit buffer for each input port in a baseline mesh router. In heterogeneous mesh network, a MPR has 7 input ports and each input port has 6-flit buffer; each input port has 8-flit buffer in a CPR. This configuration can reduce 31.7% percent buffer area while can achieve better performance compared to the baseline mesh network.

**Number of CPR and MPR:** our design goal is to build a area- and energy-efficient heterogeneous NoC, which has to meet constrains of area overhead and power consumption with respect to that of the baseline mesh NoC. These constrains determine the number of MPR in heterogeneous mesh network. For the case  $4 \times 4$  heterogeneous mesh network, to calculate the number of MPR and CPR in the network, we use the inequality:

$$0.017 \times 16 \geq 0.012 \times n + 0.021 \times (16 - n) \quad (1)$$

where 0.17 is the power consumption (in watts) of the baseline router; 0.21 and 0.12 are the power consumption (in watts) of the MPR and CPR, respectively. These power consumption values were achieved from ORION 3.0 [13]. The number 16 is the total number of routers in the baseline mesh network and  $n$  is the number of CPRs. Simplifying the inequality gives us  $n \geq 7.1$ . So we can have a minimum of 8 CPRs in the heterogeneous mesh, which is power and area efficient configuration than the same size homogeneous network. For diagonal redistribution in the heterogeneous network, we select 8 MPRs.

**Redistribution of CPR and MPR:** for heterogeneous mesh network, the number of buffer, bandwidth of link and placement of CPR and MPR will influence the performance of network. It is time to decide the placement of two types of routers after redistributing buffers and bandwidth of links. In our design, we select diagonal link, so the MPRs are placed on the diagonal. Diagonal placement of multi-port router is efficient [7] for enhancing network performance, which is been employed in our design. This placement helps center

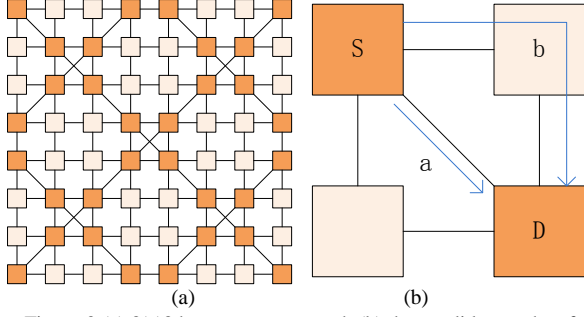


Figure 3 (a)  $8 \times 8$  heterogeneous mesh (b) the candidate paths of different source-destination configurations in MPR

traffic get through the congest region, and then alleviate the network congestion.

### B. Construction of larger heterogeneous mesh network

Section III (A) shows a case of heterogeneous mesh architecture. When we build the larger heterogeneous mesh network, the built heterogeneous mesh network should meet the basic constrains of area and power with respect to same size baseline mesh network. The extended heterogeneous mesh network also has lower diameter and scalability. Using the case of heterogeneous mesh architecture to construct the larger heterogeneous mesh is ease. Figure 3 (a) shows a case of built heterogeneous mesh network. This heterogeneous mesh network meets constrains on area and power consumption due to containing four heterogeneous mesh architectures. The larger heterogeneous mesh network is of properties of heterogeneous mesh architecture. Detailed performance evaluation is seen in Section IV.

### C. Routing design

There are two kinds of router in our proposed heterogeneous mesh network. One kind is MPR on main and second diagonals, which has 7 input and 7 output ports: East, South, West, North, West-North, East-South and Local. The other kind is CPR, which has 5 input and 5 output ports as the router of baseline mesh: East, South, West, North and Local. We adopt XY routing algorithm in our heterogeneous network and the baseline mesh network due to its simplicity. However, the main and second diagonal links in heterogeneous mesh network can provide more path selection than the baseline mesh network, how to leverage efficiently these diagonal links to improve network performance needs to carefully design the routing strategy.

We have developed a Quasi-Dimension-Ordered Routing (QDOR) for MPR due to heterogeneous mesh network with diagonal links, which the routing rule is routing packet taking x-dimension or diagonal link first and then taking y-dimension,. The other routers, CPRs, use XY routing. When the routing packet arrives at the diagonal router, the packet may have more than one output selection. For example, fig.3 (b) shows that source node S sends packet to destination D. Obviously the shortest path is the path *a* rather than XY dimension-ordered path *b*. However, if we always choose the shortest path *a*, there will be severe congestion on diagonal

links and low utilization on vertical or horizontal links. In order to avoid congestion and balance traffic in network, QDOR relaxes the output port selection. In this example, it allows packets to take path *b* depending on the network state. QDOR is free from deadlock and livelock. We employ selection strategy based on the number of free buffer on the next hop node, which select always the next hop node that has more free buffer.

## IV. PERFORMANCE EVALUATION

### A. Evaluation Methodology

We evaluate the effectiveness of our proposed heterogeneous mesh network by using the extended open source simulator, noxim [12], which is a flit-accurate simulator based on systemC. In simulation, the number of the bits included in a flit in the baseline mesh and our heterogeneous mesh is 192 and 128, respectively. A data packet consists of 1000 bits and is decomposed into 6 flits in the baseline mesh and 8 flits in our heterogeneous mesh. After a packet is generated, it is stored in an infinite queue at the source node and waits for being injected into the network. This mechanism referred to as the open-loop measurement configuration isolates the packet generation from the network condition. Each simulation executes 1,000 clock cycles for worm-up and then continues for 10,000 cycles during which performance measurements are conducted.

### B. Performance Evaluation

In this section, we show the simulation results and analysis of it on  $4 \times 4$  and  $8 \times 8$  heterogeneous mesh networks.

Figures 4 illustrate network latency, average throughput and energy consumption as a function of injection rate under uniform random (UR) synthetic traffic pattern for network sizes of  $4 \times 4$  and  $8 \times 8$ , respectively. We have the following observations based on these simulation results:

- Compared to the baseline mesh network configuration, our heterogeneous mesh reduces network latency by up to 31% and 35%, improves throughput up to 16.7% and 12.5% and lowers energy consumption up to 7.3% and 9.6% over  $4 \times 4$  and  $8 \times 8$  mesh networks under UR traffic pattern, respectively. The reason is that, on one hand, the diagonal links provide more bandwidth in centre of our proposed network than that of the baseline mesh to transfer packets stayed in network; on the other hand, diagonal links reduce average routing hops, which can alleviate quickly the center congestion in network.
- Though maximum reduction 31.7% of buffer area, heterogeneous mesh also provides greater than or equal to average latency and throughput performance provided by the baseline mesh. It all most enjoys obvious higher saturation throughput than the baseline mesh network. This is mainly due to the non-uniform distribution of the traffic which can take advantage of diagonal links. Packets experience lower hop count in heterogeneous mesh even high heavy traffic. However, as traffic increases, more contention occurs hence more waiting time in buffers which leads to increasing latency.

## V. CONCLUSION

In this paper, we leverage the non-uniform resource utilization in homogeneous mesh network to build area- and energy heterogeneous mesh network, composed of MPR and CPR, under the constraints of area, power and bandwidth. We explore the design space analysis in choosing the size, number and placement of these routers placed along the mesh network diagonals performs significantly better than the traditional homogeneous network under a variety of traffic patterns. Extensive experimentation shows that for synthetic traffic patterns, our proposed heterogeneous mesh design with MPR and CPR along the diagonals provides the maximum benefits.

## ACKNOWLEDGMENT

This work is partly supported by National Science Foundation of China under Grants No. NSFC 60773199, U0735001, NSFC 61073055, NSFC 4103470, 985-III fund and the Project of Department of Education of Guangdong Province No.2013KJCX0128.

## REFERENCES

- [1] Dally W J, Towles B. Route packets, not wires: On-chip interconnection networks[C]//Design Automation Conference, 2001. Proceedings. IEEE, 2001: pp. 684-689.
- [2] Benini L, De Micheli G. Networks on chips: A new SoC paradigm [J]. Computer, 2002, 35(1): pp. 70-78.
- [3] Marculescu R, Ogras U Y, Peh L S. Outstanding research problems in NoC design: system, microarchitecture, and circuit perspectives[J]. Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on, 2009, 28(1): pp. 3-21.
- [4] Ogras U Y, Marculescu R, Lee H G, et al. Communication architecture optimization: making the shortest path shorter in regular networks-on-chip[C]//Design, Automation and Test in Europe, 2006. DATE'06. Proceedings. IEEE, 2006, pp. 1-6.
- [5] Wang X M, Bandi L. X-Network: An area-efficient and high-performance on-chip wormhole interconnect network [J]. Microprocessors and Microsystems, 2013, 37(8): pp. 1208-1218.
- [6] Wang C, Hu W H, Lee S E, et al. Area and power-efficient innovative congestion-aware Network-on-Chip architecture[J]. Journal of Systems Architecture, 2011, 57(1): pp. 24-38.
- [7] Mishra A K, Vijaykrishnan N, Das C R. A case for heterogeneous on-chip interconnects for CMPs[C]//Computer Architecture (ISCA), 2011 38th Annual International Symposium on. IEEE, 2011: pp. 389-399.
- [8] Grot B, Hestness J, Keckler S W, et al. Kilo-NOC: a heterogeneous network-on-chip architecture for scalability and service guarantees[J]. ACM SIGARCH Computer Architecture News, 2011, 39(3): pp. 401-412.
- [9] Cheng L, Muralimanohar N, Ramani K, et al. Interconnect-aware coherence protocols for chip multiprocessors[J]. ACM SIGARCH Computer Architecture News, 2006, 34(2): pp. 339-351.
- [10] Abts D, Enright Jerger N D, Kim J, et al. Achieving predictable performance through better memory controller placement in many-core CMPs[C]//ACM SIGARCH Computer Architecture News. ACM, 2009, 37(3): pp. 451-461.
- [11] Volos S, Seiculescu C, Grot B, et al. CCNoC: Specializing on-chip interconnects for energy efficiency in cache-coherent servers[C]//Networks on Chip (NoCS), 2012 Sixth IEEE/ACM International Symposium on. IEEE, 2012: pp. 67-74.
- [12] Noxim. <http://sourceforge.net/projects/noxim>.
- [13] ORION3.0. <http://vlsicad.ucsd.edu/ORION3/>.
- [14] Duato J. Interconnection networks: an engineering approach [M]. Morgan Kaufmann, 2003.
- [15] Teig S L. The X. architecture: not your father's diagonal wiring [C]//Proceedings of the 2002 international workshop on System-level interconnect prediction. ACM, 2002: pp. 33-37.

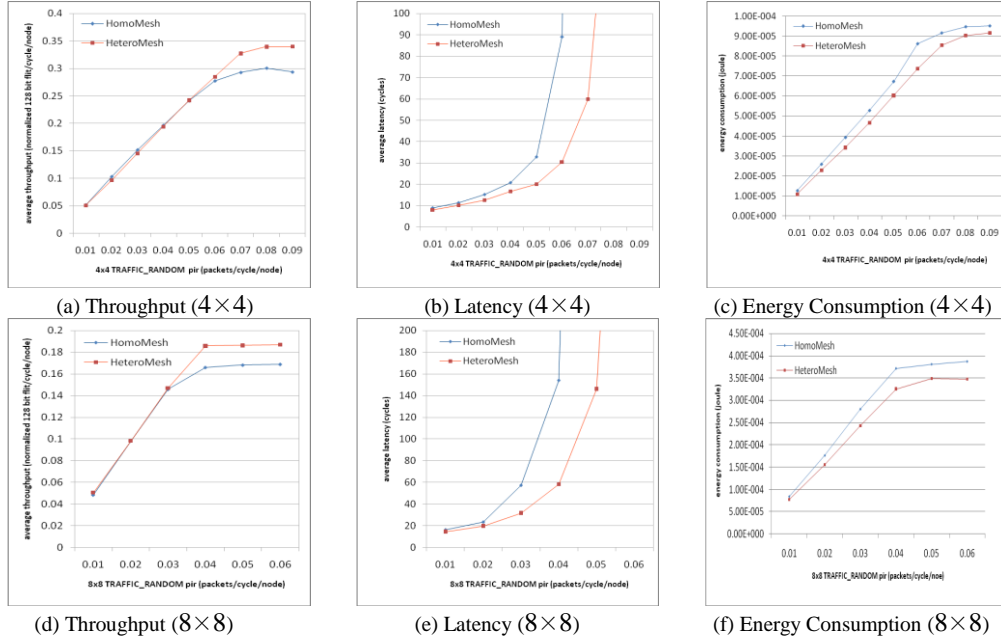


Figure 4 performance evaluations. HomoMesh standing for homogeneous mesh and HeteroMesh standing for heterogeneous mesh