

# A Novel Performance Metric of Routing Algorithm

Minghua Tang

Department of Computer Science and Technology  
 GuangDong University of Finance  
 Guangzhou, 510521, China.  
 fractal218@126.com

**Abstract**—The widely used routing algorithm performance metric of adaptiveness cannot precisely measure performance of routing algorithm. In this paper, we propose a new metric of routing pressure for measuring routing algorithm performance. It has higher precision of measuring routing algorithm performance than degree of adaptiveness. Performance of routing algorithm can be evaluated through routing pressure without simulation. It can explain why congestion takes place in network. In addition, where and when congestion takes place can be pointed out without simulation.

**Keywords**-routing pressure, Network-on-Chip, routing algorithm.

## I. INTRODUCTION

NoC is presented as scalable communication architecture for system-on-chip which will integrate tens or hundreds of processing cores in the near future [1, 2]. The communication efficiency of NoC which is affected by a lot of factors is most important for the whole system.

After NoC topology is set up, routing algorithm plays the most important role in determining NoC performance. To measure performance of a newly designed routing algorithm, researchers often take the degree of adaptiveness as the metric [3, 4, 5]. High degree of adaptiveness means that packets are provided more paths to reach their destinations. Thus they have more chances to avoid congested nodes and arrive at destinations as fast as possible.

Nevertheless, our study shows that the degree of adaptiveness is not a suitable choice as a metric of routing algorithm performance.

Firstly, the degree of adaptiveness cannot accurately measure performance of routing algorithm. A group of routing algorithms may significantly vary in performance although they have the same degree of adaptiveness.

Secondly, it cannot explain why congestion takes place.

Thirdly, it is difficult to compare the degree of adaptiveness of two routing algorithms.

Finally, we cannot know the maximum packet injection rate can be taken under a given degree of adaptiveness of a routing algorithm.

In this paper, we propose routing pressure to be the new metric to measure performance of routing algorithm. Difference in performance of routing algorithms can be explained by routing pressure. Routing pressure can account for why congestion takes place and predict where and when congestion is going to occur. It is then meaningful to compare the channel pressure of two routing algorithms. The

relationship between the maximum packet injection rate and channel pressure is clarified by a formula.

The rest of the paper is organized as follows: In next Section, we analyze the degree of adaptiveness as a metric. Then we propose the new metric of routing pressure in Section 3. In Section 4, we exemplify how to measure routing performance by routing pressure. In Section 5, routing pressure is used to account for congestion. The relationship between routing pressure and packet injection rate is formulated in Section 6. In the last Section, we make some conclusions.

## II. ANALYSIS ON DEGREE OF ADAPTIVENESS

In this paper we focus on two-dimensional mesh topology. Origin locates at the top left corner of mesh. The X axis is horizontal direction, and the Y axis is vertical direction. The positive direction of X axis points to east. The positive direction of Y axis points to south.

In uniform traffic, packets generated at a node are sent to other nodes with the same possibility. For a symmetrical network of size N, packets generated at node (i, j) are all sent to node (N-1-j, N-1-i) in transpose1 traffic scenario. In transpose2 traffic, node (i, j) only sends packets to node (j, i).

In this paper, we study routing algorithms through detail simulations, the configurations are shown in Table 1. We use the *Noxim* [6] simulator which is an open source simulator and based on SystemC. The network payload traffic is regulated by packet injection rate (referred to as PIR).

TABLE I. SIMULATION CONFIGURATIONS

|                             |                             |
|-----------------------------|-----------------------------|
| <b>Simulator</b>            | <i>Noxim</i>                |
| <b>Topology</b>             | Mesh-based                  |
| <b>Network size</b>         | 7×7                         |
| <b>Port buffer</b>          | Four flits                  |
| <b>Switch technique</b>     | Wormhole switching          |
| <b>Arbitration</b>          | Round-Robin                 |
| <b>Selection strategy</b>   | Random                      |
| <b>Traffic scenario</b>     | Transpose1, Transpose2      |
| <b>Packet size</b>          | Eight flits                 |
| <b>Traffic distribution</b> | Poisson                     |
| <b>Routing algorithm</b>    | turn model, OE, APSRA, RABC |

In this paper, we adopt the following definition of degree of adaptiveness.

Definition 1. The adaptiveness for a source-destination pair is the number of shortest paths allowed by a certain routing. For instance, in deterministic routing, there is only one path for any source-destination pair. Therefore, the adaptiveness of every source-destination pair is one.

Definition 2. The adaptiveness of a routing is the summation of adaptiveness for all source-destination pairs.

Using the method proposed in [7], a lot of routing algorithms can be constructed. Consequently, it is possible to completely study routing algorithm due to the abundant research material.

Four routing algorithms are constructed with the method for mesh topology of size  $7 \times 7$ , which are named R1, R2, R3 and R4, respectively. The four routing algorithms have the same adaptiveness.

Figure 1 shows latency variation of the four routing algorithms, under uniform traffic scenario. Although they have the equivalent adaptiveness, there exists huge difference in their performance. For example, average packet latency is 42, 61, 112 and 1153 for routing R1, R2, R3 and R4 respectively, when PIR is 0.014.

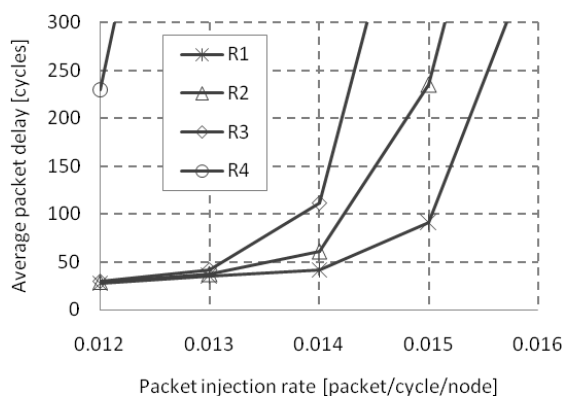


Figure 1. Latency variations for routings R1, R2, R3 and R4, which have the same adaptiveness.

Another three routing algorithms are labeled R5, R6 and R7, respectively. The degree of adaptiveness for R5, R6 and R7 are 16196, 16708 and 17742, respectively. The average packet latency variations for the three routings are depicted in Figure 2. High degree of adaptiveness does not bring about high performance. Although the adaptiveness of R7 is 9.5% higher than R5, performance of R7 is significantly lower than R5.

### III. ROUTING PRESSURE

Given a deadlock free routing algorithm R for a topology TG, and a traffic tr consisting of n communication pairs. In its  $i$ th communication comm <sub>$i$</sub> , the source node and destination node are src <sub>$i$</sub>  and dest <sub>$i$</sub> , respectively. The packet injection rate for the  $i$ th source node is PIR <sub>$i$</sub> .

Suppose there are  $p_i$  paths for the  $i$ th communication and the packets generated at source node src <sub>$i$</sub>  are uniformly distributed across those  $p_i$  paths. Then each channel on a path

will process  $PIR_i/p_i$  of packets coming from the source node. If k paths pass through a channel then that channel will assume  $k*PIR_i/p_i$  of the packets from source node.

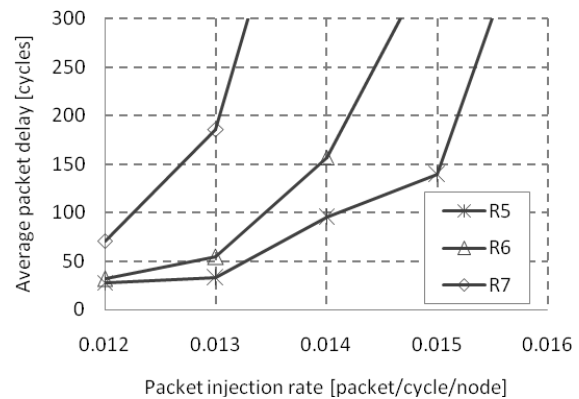


Figure 2. Latency variations for routings R5, R6 and R7, which have different adaptiveness.

For a channel ch, the  $i$ th communication will impose a certain share of packets on it. That share is referred to as shr <sub>$i$</sub> , where  $shr_i = \rho_i * PIR_i$ .

Definition 3. Given a channel ch, the summation of shares of packets imposed by all communication pairs in traffic tr is called its channel pressure, prss<sub>ch</sub>,

$$prss_{ch} = \sum_{i=1}^n shr_i = \sum_{i=1}^n \rho_i * pir_i$$

The channel which has the largest channel pressure will receive the maximum volume of traffic under the given routing algorithm. If the received message is beyond its processing capability congestion occurs at that channel. Otherwise, congestion will not occur at that channel.

If congestion does not take place at the channel with the largest channel pressure, it will not happen at other channels. Consequently, the largest channel pressure can be used to estimate whether congestion will occur in the network or not. Furthermore, it can be used to measure the performance of routing algorithm.

Definition 4. Given a routing algorithm and a traffic pattern, the routing pressure refers to the largest channel pressure among all channels under the given traffic.

### IV. USING ROUTING PRESSURE TO MEASURE ROUTING PERFORMANCE

Four routing algorithms with the same routing pressure for transpose1 traffic are created for  $7 \times 7$  mesh network. They are labeled as R8, R9, R10 and R11. The latency variations are shown in Figure 3. As can be observed in Figure 3, there is no difference in the performance of the four routings which have the same routing pressure.

In the second case, other four routing algorithms with different routing pressure for transpose1 traffic are constructed for  $7 \times 7$  mesh network. They are labeled R12, R13, R14 and R15, which have routing pressure 5.42, 6.31, 7.15 and 8.28, respectively. Figure 4 shows the latency

variations for the four routings. As routing pressure increases, its performance decreases accordingly.

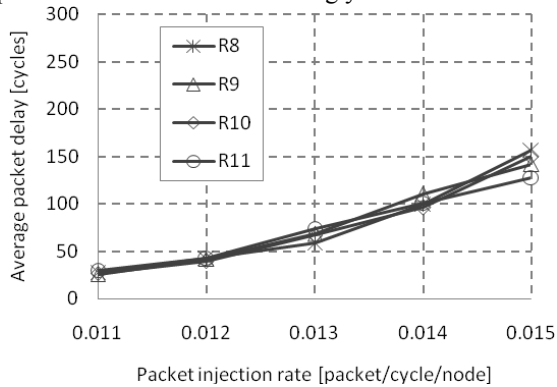


Figure 3. Latency variations for routings R8, R9, R10 and R11, which have the same routing pressure.

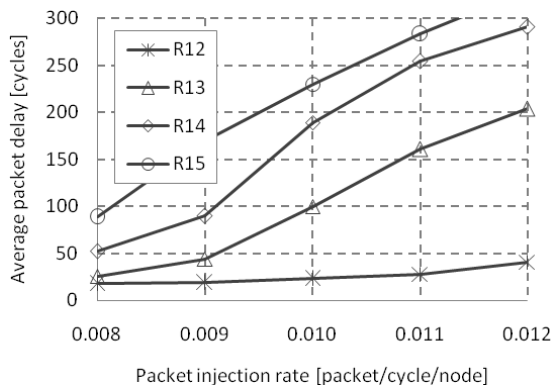


Figure 4. Latency variations for routings R12, R13, R14 and R15, which have different routing pressure.

Several routing algorithms are chosen as examples to study the relationship between routing pressure and routing performance. They are XY, OE, turn model (negative-first) and APSRA.

For transpose1 traffic, routing pressures for the four routings are 6, 4.81, 6, and 8.86, respectively. The simulation results are shown in Figure 5. OE has the best performance because its routing pressure is smallest. XY and negative-first (NF) have the same routing pressure. There is no difference in their performance. APSRA has the worst performance due to its large routing pressure.

For transpose2 traffic, routing pressures for the four routings are 6, 4.81, 2.41, and 5.44, respectively. Figure 6 shows the simulation results for them. As expected, negative-first has the best performance. In this traffic, XY has the worst performance because of its largest routing pressure.

### V. ROUTING PRESSURE AND CONGESTION

To analyze the relationship between congestion and channel pressure, we take routing R12 when PIR is 0.012 as an example. The channel which has the largest pressure of 5.42 is 17-24.

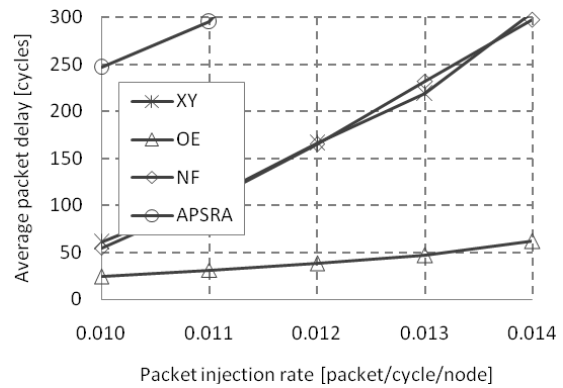


Figure 5. Latency variations of routings XY, OE, NF and APSRA for transpose1 traffic.

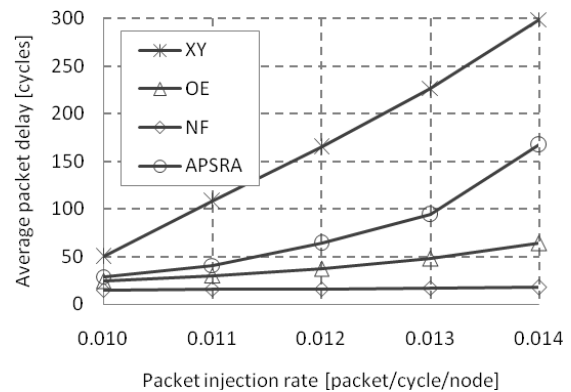


Figure 6. Latency variations of routings XY, OE, NF and APSRA for transpose2 traffic.

The average channel utilization ratio of all the network channels is 24%. However, the utilization ratio of channel 17-24 is as large as 81.4%. It shows that channel 17-24 is highly utilized.

Figure 7 shows channel utilization ratio variations along with channel pressure. As channel pressure increases, utilization of channel rises up. The channel with the largest channel pressure has the highest utilization ratio.

At runtime, the channel utilization ratio may not linearly increase with channel pressure because the high utilization of channel 17-24 brings about congestion. The congestion makes utilization ratios of some channels lower than expected.

For example, there is a low utilization ratio for channel with pressure of 0.037. That channel is 26-27. Under transpose1 traffic, only messages generated at node 0, 1, 2 and 3 will pass through channel 26-27. However, all packets which come from those four source nodes and pass through channel 26-27 have to go through channel 17-24. Unfortunately, there is congestion at channel 17-24, which makes few packets reach channel 26-27. It is why utilization ratio of channel 26-27 is low.

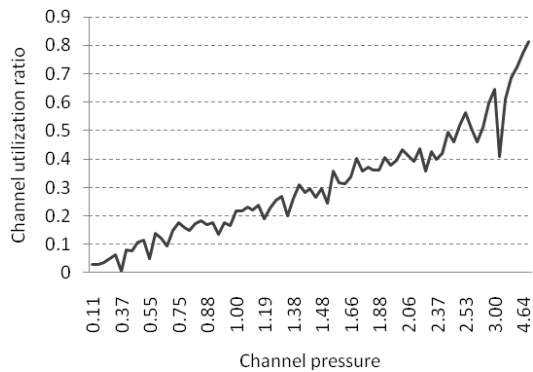


Figure 7. Utilization ratio of channels with different channel pressure.

In transpose1 traffic, one communication is from node 10 to node 26. Packets generated at node 10 reach destination node 26 after only passing through four channels. However, node 10 generated the largest latency packet among all the packets in the traffic. When computing average packet latency for every single communication pair we get average packet latency for each communication pair. The average packet latency for communication between node 10 and 26 is still the largest.

Under routing algorithm R12, there is only one path for communication between node 10 and 26, which is 10-17-24-25-26 labeled by node ID. All packets of the communication have to pass a congested channel. It is not strange packet latency of this communication is large.

Consequently, congestion takes place at the channel which has largest pressure.

#### VI. ROUTING PRESSURE AND PACKET INJECTION RATE

Given a routing pressure, the largest packet injection rate which will not lead to congestion can be computed.

To avoid congestion the amount of packets imposed on a channel cannot be larger than its processing capability.

Suppose all source nodes in the traffic have the same packet injection rate of  $PIR$ , routing algorithm has routing pressure of  $\rho * PIR$ , network channel can process  $frate$  flits per cycle, all packets have the same length of  $len$  flits. The simulation runs  $time$  cycles.

Then a channel can process at most  $frate * time / len$  packets during the simulation. The number of packets imposed on it is  $\rho * PIR * time$ . We have,

$$\rho * PIR * time \leq frate * time / len$$

Then following inequality holds,

$$\rho * PIR \leq frate / len$$

In this paper, the simulator is *Noxim* which channel needs two cycles to forward a flit, that is  $frate = 0.5$ . Each packet has eight flits,  $len = 8$ . Given a routing which has routing pressure of  $\rho * PIR$ , to avoid congestion  $PIR$  has to satisfy,

$$PIR \leq 0.0625 / \rho$$

For transpose1 traffic, APSRA routing pressure is 8.86. The maximum packet injection rate is 0.007 as shown in Figure 8. After that congestion begins to occur in the network and system performance degrades significantly.

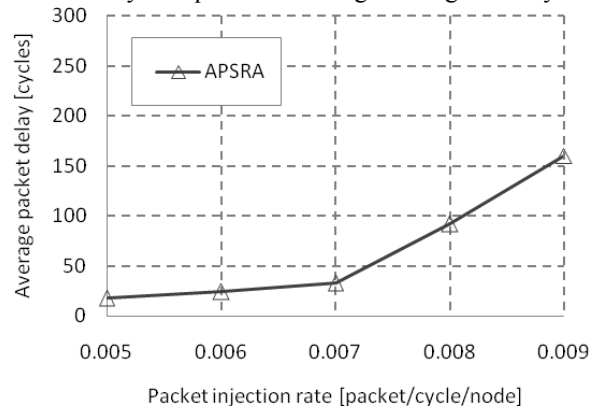


Figure 8. The maximum packet injection rate under APSRA routing for transpose1 traffic.

#### VII. CONCLUSION

In this paper, we propose a new performance metric for routing algorithm. It has higher precision than adaptiveness while measure performance of routing algorithm. In addition, it can account for why congestion takes place in the network. The maximum packet injection rate can be taken by the network can be computed from its routing pressure.

#### ACKNOWLEDGMENT

This work is supported in part by the Scientific Research Foundation for Talent Introduction (NO. 2012RCYJ013).

#### REFERENCES

- [1] W. J. Dally and B. Towles, Route Packets, Not Wires: On-Chip Interconnection Networks, Proc. ACM/IEEE Design Automation Conf., pp. 684-689, 2001.
- [2] L. Benini and G. D. Micheli, Networks on Chips: A New SoC Paradigm, IEEE Computer, vol. 35, no. 1, pp. 70-78, Jan. 2002.
- [3] C. J. Glass and L. M. Ni, The Turn Model for Adaptive Routing, J. Assoc. for Computing Machinery, vol. 41, pp. 874-902, 1994.
- [4] G.-M. Chiu, The Odd-Even Turn Model for Adaptive Routing, IEEE Trans. Parallel and Distributed Systems, vol. 11, no. 7, pp. 729-738, July 2000.
- [5] M. Palesi, R. Holsmark, S. Kumar and V. Catania, Application Specific Routing Algorithms for Networks on Chip, IEEE Trans. Parallel and Distributed Systems, vol. 20, pp. 316-330, 2009.
- [6] Sourceforge.net, Noxim: Network-on-chip simulator, 2008. [Online]. Available: <http://noxim.sourceforge.net>
- [7] M. H. Tang and C. H. Wu, A New Method of Designing NoC Routing Algorithm, 2nd International Conference on Consumer Electronics, Communications and Networks, pp. 3044-3047, 2011.