

A Performance Model for Network-on-Chip Wormhole Routers

Yuhui Zhang

Department of Computer Science and Technology, Tsinghua University, Beijing, 100084, China

Email:zyh02@tsinghua.edu.cn

Xiaoguo Dong and Siqing Gan

Department of Mathematical Science and Computing Technology, Central South University, Changsha, 410075, China

Weimin Zheng

Department of Computer Science and Technology, Tsinghua University, Beijing, 100084, China

Abstract—A generic analytical performance model of single-channel wormhole routers is presented using the M/D/1/B queuing theory. Compared with previous work, the flow-control feedback mechanism is studied in detail, and a computing method bases on Markov chain for the flow-control feedback probability is proposed. Compared with BookSim, a well-known cycle-accurate Network-on-Chip (NoC) simulator, this model presents accurate results on key metrics: the average relative error of flow-control feedback probability is about 7.87%. In addition, based on the model of single-channel routers, the asymmetric multi-channel and symmetric multi-channel structured routers are both modeled respectively.

Index Terms—Network-on-Chip; Markov chain; queuing theory; performance analysis

I. INTRODUCTION

Networks-on-Chip (NoC) [1][2] has been proposed as a solution for addressing the design challenges of high performance nano-scale architectures by separating the on-chip communication from computing and storage. Connecting components through an on-chip network has several advantages over dedicated wiring, potentially delivering high-bandwidth, low-latency, and low-power communication over a flexible, modular medium.

Wormhole-routing [3][4] is a system of simple flow control in NoC based on fixed links, which makes message latency almost independent of the inter-node distance in the absence of blocking. In wormhole routing, network packets are broken into small pieces called flits. The first flit, called the header flit holds information about this packet's route (namely the destination address) and sets up the routing behavior for all subsequent flits associated with the packet. The head flit is followed by zero or more body flits, containing the actual pay load of data. The final flit, called the tail flit, performs some bookkeeping to close the connection between the two

nodes. If the header is blocked, the data flits are blocked behind it occupying all the channels and buffers already taken [3].

Moreover, the routing algorithm defines how to transfer a message through a network path. A key issue for any routing algorithm is deadlock-free. Now deterministic routing for deadlock prevention is widely employed [3].

Currently, many NoC designs are based on the simulation method. However, simulation is a time-consuming procedure, especially within the large design space. Therefore, system designers have to choose limited assessment in the design space and then cannot get the optimized results usually.

Another approach is utilization of an analytical model of the system which is the most cost-effective tool for performance evaluation. Several analytical models of deterministic routing in wormhole-routed mesh-based networks, e.g. hypercube and tori, have been reported.

Ref.[5] introduces a probability model for wormhole network. Ref.[6] is restricted to the k-ary n-cubes topology while [7] is only adapted by the hypercube. Ref.[8] presents a performance model based on some queuing theory, but it only can apply to the switched network. Ref.[9] presents an analytical performance evaluation method for the general wormhole-routing NoC architectures. It analyses different reasons for blocking network packets in the router buffer, but does not consider the flow-control feedback mechanism, which is caused by the fullness of the input queue of the successive router. Ref.[10], based on the same assumptions of [9], provides a novel method: the numerical analysis and iterative computation is used to estimate the NoC performance.

Moreover, several models for adaptive wormhole routing have been introduced for the torus and hypercube [11][12][13].

This paper presents a generic analytical performance evaluation approach of NoC design. Different from the previous work, the flow-control feedback probability between adjacent routers is considered meticulously, which is the important indicator of the buffer utilization

Supported by Chinese National Science Foundation under Grant 60773147, and Chinese National 863 High Technology Programs under Grant 2008AA01A204, Email: zyh02@tsinghua.edu.cn.

and network traffic. And how to extend the single-channel model to the model for multi-channel routers is also presented.

In summary, this paper gives the following contributions:

1) A general analytical model of wormhole routers with single channel is proposed, which supports arbitrary network topologies, the deterministic routing algorithm, arbitrary packet / buffer lengths, and so on.

Based on the router model, the computing method of flow-control feedback probability using the Markov chain is presented in detail.

2) Based on the performance model of the single-channel router, the asymmetric multi-channel and symmetric multi-channel structured routers are both modeled respectively. Multi-channel structures can improve the communication performance significantly.

3) The accuracy of this approach is validated through the comparisons with a well-known cycle-accurate simulator, BookSim [14].

The remaining of this paper is organized as follows. Section 2 introduces the related work on performance modeling. Section 3 gives our modeling assumptions. The model for single channel routers is presented in the next section and the multi-channel versions are presented in Section 5. Experimental results are given in Section 6; the conclusion and future work are introduced at last.

II. RELATED WORK

Ref. [5] develops a model of a single wormhole router, which is based on the probability analysis rather than queuing theory. The model is evaluated through a series of flit-level simulations. Moreover, how to extend the model to networks of routers is also discussed.

Ref. [6] analyzes communication networks of varying dimension under the assumption of constant wire bisection. Models of the latency, average throughput, and hot-spot throughput of k-ary n-cube networks are presented.

Ref. [7] proposes a general analytical model to predict message latency in wormhole-routed k-ary n-cubes with fully adaptive routing. The analysis focuses on a widely-accepted fully adaptive routing algorithm.

A system-level buffer planning algorithm is given in [8]. Using this algorithm, the buffer depth for each input channel in different routers across the chip can be derived to optimize the overall performance, given the traffic characteristics of the target application and the total budget of the available buffering space.

Ref. [9] presents a generalized router model and then utilizes this model for doing NoC performance analysis. The proposed model can be used not only to obtain fast and accurate performance estimates, but also to guide the NoC design process within an optimization loop.

Ref. [10] also presents a generic analytical method to estimate communication latencies and link-buffer utilizations for a given NoC architecture. The accuracy of this method is experimentally compared with the results obtained from Cycle-Accurate SystemC simulations. It is based on the same assumptions of [9], providing a novel

method: the numerical analysis and iterative computation is used to estimate the NoC performance.

III. ASSUMPTIONS OF ROUTER MODELING

A. NoC Router

An illustration of the router structure is presented in Fig.1.

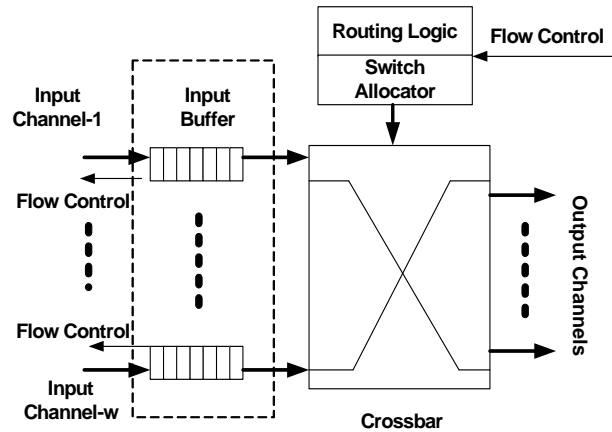


Figure 1. The structure of a router

We assume that a wormhole router contains w ports and adopts a deterministic routing algorithm (like X-Y routing). The local port is regarded as the same as others; each port is associated with a single input buffer (single channel; how to model the multi-channel structure is discussed in Section V). Input buffers are abstracted as groups of FIFO.

The crossbar switch can be configured to connect any input buffer of the router to any output channel, but under the constraints that each input is connected to at most one output, and each output is connected to at most one input. The tasks of resolving all the potential requests to the crossbar and other shared resources of the router fall onto the switch allocator.

Before transfer, data packets are divided into small pieces, called flits. The header flit holds destination information to set up the transfer channel for all subsequent flits of the same packet.

The router also implements flit-level flow control: the

TABLE I.
PARAMETERS OF THE MODEL

| Symbol | Description |
|----------|--|
| w | Number of ports of a router |
| B | Length of the input buffer |
| T | Service time of a packet, which does not include the waiting time in the queue |
| P | Packet size (in flit) |
| H_s | Service time of the header flit, or the time of the header flit going through the router (without the waiting time). It is also the number of pipeline stages of a router because of Assumption 4. |
| (i, j) | Port j ($0 \leq j < w$) at Router i . |

input queue capacity is limited, so if any input queue of a router is full, its upstream router will stop transmitting. In the previous work, flow-control feedback is less involved or even ignored.

For the detailed workflow of a wormhole router, please refer to [15].

B. Modeling Parameters

As in [9][10], the router model introduces the following hypotheses.

1) Network traffic is generated from all nodes uniformly and follows the Poisson process, which implies that packet-arrival events occur continuously and independently of one another.

2) Packet destinations are equally distributed across the network nodes following a uniform traffic pattern.

3) Traffic sinks consume the incoming packets with the constant rate of one flit per cycle.

4) A pipeline stage of a router can deal with one flit per cycle, too.

5) Input buffers of any router have finite capacity.

6) When the system achieves asymptotic stability, the service time of packets is approximately equal to a constant value.

Some parameters of the model are listed in Table 1.

Then, we have (1) because of Assumption 4.

$$T = H_s + P \quad (1)$$

IV. ANALYTICAL MODEL OF SINGLE-CHANNEL ROUTERS

A. Analysis of waiting time in the input queue

This section focuses on modeling a single-channel router as a set of first-come first-serve buffers connected by a crossbar switch to analyze the average waiting time that an incoming packet spends in the queue.

$T_{i,j}$ denotes the average time in the queue of (i, j) . It is composed of the following three parts as described in [9]:

1) *Service time of the packets already waiting in the same buffer;*

2) *The residual service time seen by an incoming packet;*

3) *The packets waiting in other buffers of the same router and served before the incoming packet.*

As stated in [9], the traffic rate at (i, j) , $\lambda_{i,j}$, can be computed by (2).

$$\lambda_{i,j} = \sum_s \sum_{\forall d} x_{s,d} R(s, d, i, j) \quad (2)$$

Here $x_{s,d}$ represents the traffic rate from source router s to the destination d , and $\pi_{s,d}$ indicates the path from s to d . R is the indicator function that returns 1 if the path goes through (i, j) , and returns 0 otherwise.

$$R(s, d, i, j) = \begin{cases} 1, & (i, j) \in \pi_{s,d} \\ 0, & (i, j) \notin \pi_{s,d} \end{cases} \quad (3)$$

I. The first part

Suppose the average number of packets in the input buffers of router i is a vector, N .

$$N = [N_1, N_2, \dots, N_w]^T. \quad (4)$$

Then, the average number of packets waiting in port j is N_j . And the average waiting time for incoming packet is $E(T) \times N_j$, where $E(T)$ indicates the mean of service time.

II. The second part

If incoming packet p_m arrives at the top of the input queue while some other packet p_n is being serviced, then the residual service time R_m for p_m is the time left for the packet p_n to finish its service.

To simplify the analysis, the concept of mean residual service time $R(\lambda_{i,j})$ is used. It is the average value of service time for all packets, as well as a function of the traffic rate and service time. When the system reaches the asymptotic steady-state, the following equation represents its value:

$$R(\lambda_{i,j}) = \frac{1}{2} \lambda_{i,j} \times E(T^2). \quad (5)$$

where $E(T^2)$ is the second moment of service time.

III. The third part

In [9][10], both Part 1 and 2 have been analyzed completely. But for Part 3, they do not consider the flow-control feedback. Therefore, we focus on this issue.

In detail, when a head flit intends to go to the specific output port, it has to compete with all other flits applying for the same direction. Moreover, another necessary condition for any winner to continue is that the input queue of the downstream router is not full, which is called the flow-control feedback.

Then, a packet transmitted from (i, j) to $(i+1, k)$ consists of two processes: *competition* and *flow-control*.

Suppose $F_{i,j,k}$ is the probability of the header flit transmitted from (i, j) to $(i+1, k)$, and $p_{i+1,k}$ is the flow-control feedback probability produced from $(i+1, k)$ and $f_{i,j,k}$ is the competition probability of the header flit. Then we have:

$$F_{i,j,k} = f_{i,j,k} \times (1 - p_{i+1,k}). \quad (6)$$

$\lambda_{i,j,k}$ is the traffic rate from (i, j) to $(i+1, k)$ and we get

$$f_{i,j,k} = \frac{\lambda_{i,j,k}}{\sum_{l=1}^w \lambda_{i,l,k}}. \quad (7)$$

$c_{i,j,q}$ denotes the competition probability of the header flits in (i, j) and (i, q) transmitting to the same input port of Router $(i+1)$.

We have $c_{i,j,q} = 1$ if $j=q$.

If $1 \leq j, q \leq p$ and $j \neq q$, we can get

$$c_{i,j,q} = \sum_{k=1}^w F_{i,j,k} F_{i,q,k} = \sum_{k=1}^w f_{i,j,k} f_{i,q,k} (1 - p_{i+1,k})^2 \quad (8)$$

Therefore, the blocking delay caused by packet competitions and flow controls can be denoted by (9).

$$E(T) \sum_{q=1, q \neq j}^w c_{i,j,q} N_q = E(T) \sum_{q=1, q \neq j}^w \sum_{k=1}^w f_{i,j,k} f_{i,q,k} (1 - p_{i+1,k})^2 N_q \quad (9)$$

Summing up the three parts, we obtain the average waiting time of an incoming packet buffered in (i, j) .

$$T_{i,j} = E(T) N_j + \frac{1}{2} \lambda_{i,j} E(T^2) + E(T) \sum_{q=1, q \neq j}^w \sum_{k=1}^w f_{i,j,k} f_{i,q,k} (1 - p_{i+1,k})^2 N_q \quad (10)$$

To calculate $T_{i,j}$, it is necessary to computer the flow-control feedback probability, $p_{i+1,k}$.

B. A computing method of flow-control feedback probability

When the system arrives at the asymptotic steady-state, the service time of packets can be regarded as the mean value, $E(T)$. In general, we consider the flow-control feedback probability $p_{i+1,k}$ of the input queue at $(i+1, k)$

In this figure, $p_{i+1,k}$ is produced by $(i+1, k)$ with the

arrival rate $\sum_{j=1}^w f_{i,j,k} \lambda_{i,j,k}$ and the service rate $\frac{1}{E(T)}$.

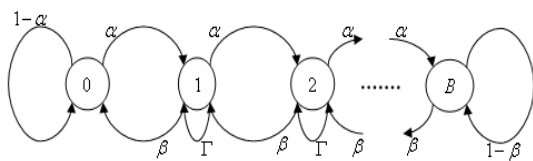


Figure 2. State transition diagram for M/D/1/B queue

Using Markov chain to analyze the changes of flits in the input queue, the state transition diagram for the queue is shown in Fig.2 and the state transition matrix can be written as follows.

$$M = \begin{bmatrix} 1-\alpha & \alpha & 0 & \cdots & 0 & 0 & 0 \\ \beta & \Gamma & \alpha & \cdots & 0 & 0 & 0 \\ 0 & \beta & \Gamma & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \Gamma & \alpha & 0 \\ 0 & 0 & 0 & \cdots & \beta & \Gamma & \alpha \\ 0 & 0 & 0 & \cdots & 0 & \beta & 1-\beta \end{bmatrix} \quad (11)$$

where

$$\alpha = \left(\sum_{j=1}^w f_{i,j,k} \lambda_{i,j,k} \right) \times \left(1 - \frac{1}{E(T)} \right); \quad (12)$$

$$\beta = \left(1 - \sum_{j=1}^w f_{i,j,k} \lambda_{i,j,k} \right) \times \frac{1}{E(T)}; \quad (13)$$

$$\tau = \left(\sum_{j=1}^w f_{i,j,k} \lambda_{i,j,k} \right) \times \frac{1}{E(T)} + \left(1 - \sum_{j=1}^w f_{i,j,k} \lambda_{i,j,k} \right) \times \left(1 - \frac{1}{E(T)} \right); \quad (14)$$

According to the state transition diagram, we get the equilibrium distribution vector,

$$S_{i+1,k} = [S_{i+1,k,0}, S_{i+1,k,1}, \dots, S_{i+1,k,B}]^T, \quad (15)$$

$$\text{where } \sum_{n=1}^B S_{i+1,k,n} = 1 \quad (16)$$

In (15), $S_{i+1,k,n}$ is the probability of the state having n flits filled in the input queue of $(i+1, k)$ and $S_{i+1,k,0}$ is the probability of an empty queue; $S_{i+1,k,B}$ is the probability of a full queue, which can be called the probability generating the flow-control feedback from $(i+1, k)$.

The difference equations for the state transition distribution vector can be written as follows.

$$\alpha S_{i+1,k,0} - \beta S_{i+1,k,1} = 0 \quad (17)$$

$$\alpha S_{i+1,k,n-1} - (\alpha + \beta) S_{i+1,k,n} + \beta S_{i+1,k,n+1} = 0 \quad (0 < n < B) \quad (18)$$

Then, the solution of the above difference equations can be gotten as

$$S_{i+1,k,n} = \left(\frac{\alpha}{\beta} \right)^n S_{i+1,k,0} \quad (0 \leq n \leq B) \quad (19)$$

We define the duty factor of the system as

$$\rho = \frac{\alpha}{\beta} = \frac{\left(\sum_{j=1}^w f_{i,j,k} \lambda_{i,j,k} \right) \left(1 - \frac{1}{E(T)} \right)}{\left(1 - \sum_{j=1}^w f_{i,j,k} \lambda_{i,j,k} \right) \frac{1}{E(T)}} \quad (20)$$

where

$$\sum_{n=0}^B S_{i+1,k,n} = S_{i+1,k,0} \sum_{n=0}^B \rho^n = 1 \quad (21)$$

Now, we get

$$S_{i+1,k,0} = \frac{1 - \rho}{1 - \rho^{B+1}} \quad (22)$$

And then we have

$$p_{i+1,k} = S_{i+1,k,B} = \rho^B \frac{1 - \rho}{1 - \rho^{B+1}} \quad (1 \leq k \leq p) \quad (23)$$

V. ANALYTICAL MODEL FOR MULTI-CHANNEL STRUCTURES

As the increase of network traffic, the waiting time of packets in the input queue of single-channel routers increases, which results in longer transmission delay. To solve this problem, there are two ways: reduce the packet arrival rate and improve the service rate. At present, NoC routers usually adapt multi-channel architecture to achieve the two functions.

Ref. [16] has classified multi-channel routers into two categories: the asymmetric multi-channel and symmetric multi-channel structures. Both are modeled based on our above-mentioned work. In [16], the flow-control feedback possibility is computed based on the M/G/1 queue theory, while we calculate the possibility based on the M/D/1/B queue theory and the Markov chain.

A. Two typical multi-channel NoC routers

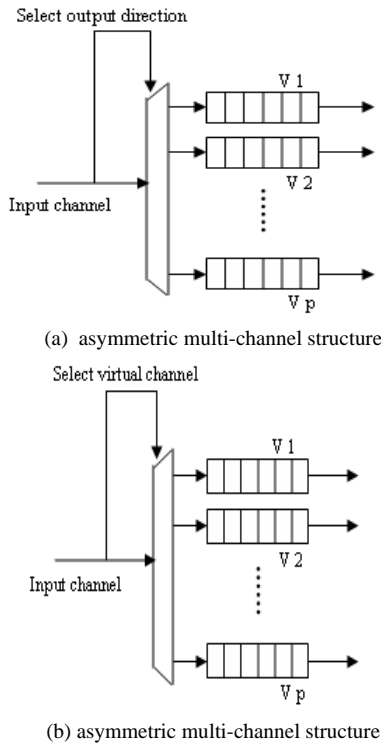


Figure 3. Multi-channel structure NoC routers

Fig.3(a) shows the asymmetric multi-channel structure: each input port is associated with p channels which match different output ports respectively; in another word, each channel matches only one output port and different channels match different output ports respectively. Flits select output direction according to their own routing information. Then the traffic is divided into different channel queues for transmission, thus the arrival rate for each buffer is reduced and the service rate is also improved.

Fig.3(b) shows the symmetric multi-channel structure: each input port is also associated with p channels, and each has the same function. The incoming packet select virtual channel with the same probability $1/p$. Different from the asymmetric, here each channel can match any one of output ports.

In the symmetric structure, the traffic is divided into several groups with same probability, and then each channel send the transmission request to the corresponding output port. If the current packet fails in the competition and is being blocked, the subsequent packets to the same output port can enter the other channels for transmission. The structure can also reduce the arrive rate and improve the service rate.

B. Analytical model for multi-channel structure

I. Analytical model for asymmetric multi-channel structure

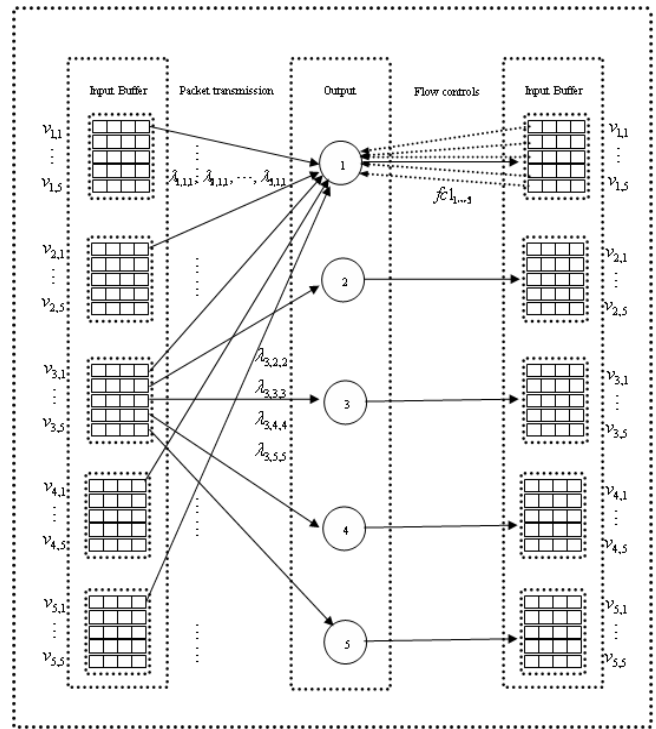


Figure 4. Arbitration model of asymmetric multi-channel structure

Assuming each input port of an asymmetric multi-channel router is associated with p channels. For example, for the 2D-mesh topology, each input port has five different output directions; the traffic transmission is illustrated in Fig.4. In this structure, the flits in the queue only compete with the flits with the same transfer direction from other input ports. For example, the flits in channel $V_{2,1}$ only compete with the flits from channel $V_{1,1}$, $V_{3,1}$, $V_{4,1}$, and $V_{5,1}$. Different from the single-channel structure, the flow-control occurs only if all five channels of the corresponding downstream router are fulfilled.

$\lambda_{i,j,k,k}$ represents the traffic rate from the channel k ($1 \leq k \leq p$) of (i, j) , which can be computed by (24).

$$\lambda_{i,j,k,k} = \sum_{vs} \sum_{vd} x_{s,d} \times R(s,d,i,j,k) \times S(s,d,k,k) \quad (24)$$

Similar with (2), here R is the indicator function that returns 1 if the path goes through the k channel of (i, j) , and returns 0 otherwise. Function $S(s,d,k,k)$ returns 1

if and only if the flits in the channel k will be sent to output port k , otherwise it returns 0. Because the flits in each input channel k can only be sent to the output port k , we have $S(s, d, k, k) = 1$.

The analysis of waiting time in the input queue k of (i, j) also includes three parts. The first and the second parts are similar with the single-channel structure. We focus on the third.

Suppose $F_{i,j,k,k}$ is the probability of the header flit transmitted from the channel k in (i, j) to $(i+1, k)$, $p_{i+1,k}$ is the flow-control feedback probability produced from $(i+1, k)$ and $f_{i,j,k,k}$ is the competition probability of the header flit. We have

$$F_{i,j,k,k} = f_{i,j,k,k} \times (1 - p_{i+1,k}^p) \quad (25)$$

$$f_{i,j,k,k} = \frac{\lambda_{i,j,k,k}}{\sum_{l=1}^p \lambda_{i,j,l,k}} \quad (26)$$

$c_{i,j,q,k}$ denotes the competition probability of the header flits in the channel k of (i, j) and (i, q) transmitting to the same input port of router $(i+1)$.

We have $c_{i,j,q,k} = 1$ if $j=q$.

If $1 \leq j, q \leq p$ and $j \neq q$, we can get

$$c_{i,j,q,k} = F_{i,j,k,k} \cdot F_{i,q,k,k} = \begin{cases} 1 & (j=q) \\ f_{i,j,k,k} f_{i,q,k,k} (1 - p_{i+1,k}^p)^2 & (1 \leq j, q \leq p; j \neq q) \end{cases} \quad (27)$$

Therefore, the blocking delay caused by packet competitions and flow controls happened in the two channels of (i, j) and (i, q) can be denoted by (28).

$$\begin{aligned} E(T) & \sum_{q=1, q \neq j}^p c_{i,j,q,k} N_{q,k} \\ & = E(T) \sum_{q=1, q \neq j}^p f_{i,j,k,k} f_{i,q,k,k} (1 - p_{i+1,k}^p)^2 N_{q,k} \end{aligned} \quad (28)$$

Summing up the three parts, the average waiting time of an incoming packet buffered in channel k of (i, j) is:

$$\begin{aligned} \tau_{i,j,k} & = E(T) N_{j,k} + \frac{1}{2} \lambda_{i,j,k} E(T^2) \\ & + E(T) \sum_{q=1, q \neq j}^p f_{i,j,k,k} f_{i,q,k,k} (1 - p_{i+1,k}^p)^2 N_{q,k} \end{aligned} \quad (29)$$

where $N_{q,k}$ ($1 \leq q, k \leq p$) stands for the average number of flits in the channel k of (i, q) .

II. Analytical model for symmetric multi-channel structure

Routers can also adopt symmetric multi-channel structure. Each input port is associated with p channels. The incoming packet selects any virtual channel with the same probability, $1/p$. Still taking the 2D-mesh topology as the example, each input port has five channels, the flits select a virtual channel with the probability of 0.2, and each channel can send flits to arbitrary output port. It means that any output port can receive the flits from all 25 channels.

The traffic transmission is shown in Fig.5. Compared with the asymmetric structure, the transfer principle is the same, but the flow-control probability is different.

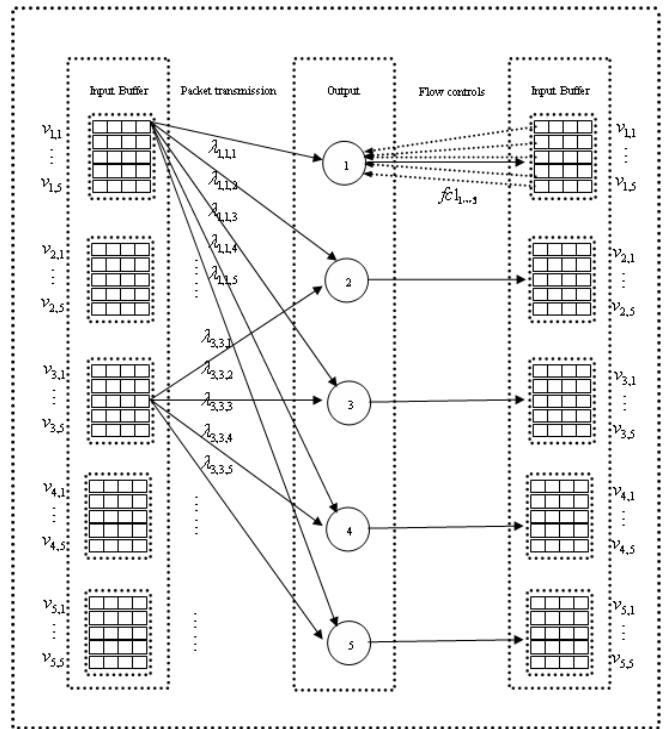


Figure 5. Arbitration model of symmetric multi-channel structure

$\lambda_{i,j,k,h}$ presents the traffic rate from channel k of (i, j) to output port h ($1 \leq h \leq p$), then we have

$$\lambda_{i,j,k,h} = \sum_{\forall s} \sum_{\forall d} x_{s,d} \times \frac{R(s, d, i, j, k) S(s, d, k, h)}{p} \quad (30)$$

Here, R is the indicator function owing the same meaning of the asymmetric version. Function S returns 1 if and only if the flits in the channel k will be sent to output port h , otherwise it returns 0.

The analysis of waiting time in the input queue k of (i, j) still includes three parts. The first and second parts are the same with the asymmetric multi-channel structure. We also focus on the third part and we have

$$F_{i,j,k,h} = f_{i,j,k,h} \times (1 - p_{i+1,h}^p) \quad (31)$$

$$f_{i,j,k,h} = \frac{\lambda_{i,j,k,h}}{\sum_{j=1}^p \sum_{l=1}^p \lambda_{i,j,l,h}} \quad (32)$$

Here F , f , and p have the same meanings of the asymmetric model.

Similarly, $c_{i,j,q,k}$ denotes the competition probability of the header flits in the channel k of (i, j) and (i, q) transmitting to the same input port of Router $(i+1)$.

If $j = q$, we have

$$c_{i,j,q,k,n} = \sum_{h=1}^p F_{i,j,k,h} F_{i,q,n,h} = \sum_{h=1}^p f_{i,j,k,h} f_{i,q,n,h} (1-p_{i+1,h}^p)^2 \quad (33)$$

If $j \neq q$, we have

$$c_{i,j,q,k,m} = \sum_{h=1}^p F_{i,j,k,h} F_{i,q,m,h} = \sum_{h=1}^p f_{i,j,k,h} f_{i,q,m,h} (1-p_{i+1,h}^p)^2 \quad (34)$$

Then, the blocking delay caused by packet competitions and flow controls can be denoted by (35):

$$\begin{aligned} E(T) & \left(\sum_{q=1, q \neq j}^p \sum_{m=1}^p c_{i,j,q,k,m} N_{q,m} + \sum_{n=1, n \neq k}^p c_{i,j,j,k,n} N_{j,n} \right) \\ & = E(T) \sum_{q=1, q \neq j}^p \sum_{m=1}^p \sum_{h=1}^p f_{i,j,k,h} f_{i,q,m,h} (1-p_{i+1,h}^p)^2 N_{q,m} \\ & \quad + E(T) \sum_{n=1, n \neq k}^p \sum_{h=1}^p f_{i,j,k,h} f_{i,j,n,h} (1-p_{i+1,h}^p)^2 N_{j,n} \end{aligned} \quad (35)$$

And the average waiting time is,

$$\begin{aligned} \tau_{i,j,k} & = E(T) \sum_{q=1, q \neq j}^p \sum_{m=1}^p \sum_{h=1}^p f_{i,j,k,h} f_{i,q,m,h} (1-p_{i+1,h}^p)^2 N_{q,m} \\ & \quad + E(T) \sum_{n=1, n \neq k}^p \sum_{h=1}^p f_{i,j,k,h} f_{i,j,n,h} (1-p_{i+1,h}^p)^2 N_{j,n} \\ & \quad + E(T) N_{j,h} + \frac{1}{2} \lambda_{i,j,k} E(T^2) \end{aligned} \quad (36)$$

VI. EXPERIMENT RESULTS

A. Test method

A well-known, third-party NoC simulator, BookSim, is used to validate the accuracy of the proposed single-channel model. Originally developed for and introduced by the *Principles and Practices of Interconnection Networks* book [15]. BookSim's functionality has been continuously extended. We use BookSim 2.0, which supports a wide range of topologies, provides diverse routing algorithms and includes numerous options for customizing the router's micro-architecture.

BookSim cannot provide the flow-control feedback probability as a result directly; therefore we employ the following method:

At each simulation cycle, we record the flit number of any input queue; when an input queue is full, the flow control signal is regarded as issued till there is some free space. Then, the probability can be computed.

In this experiment, we adopt the XY deterministic routing and a 5x5 2D-mesh network. The observed results are obtained by simulating 2×10^7 cycles after a warm-up phase of 2×10^7 cycles, and then compared with analysis results. The performance of our model is analyzed under uniform traffic patterns, where any node transfers packets towards the destinations with equal probability.

The injection rate is specified in packets per cycle. For example, the injection rate is 0.25, which means each node injects a new packet every four cycles.

The error between the analytical results V_{test} and simulation results V_{sim} is calculated by the following equation.

$$Err = \frac{|V_{sim} - V_{test}|}{V_{sim}} \times 100\%$$

B. Accuracy validation for the computing method of flow-control feedback probability

In this section, we focus on the influence of input buffer size (B), number of pipeline stages (H_s) and packet size (P) on the flow-control feedback probability.

Simulation results in Fig.6 reveal that for the different design parameters, the trend of flow-control feedback probability is roughly similar: increasing slowly, increasing rapidly and tends to balance with the increased injection rates.

Our proposed analysis model tallies closely with the simulation results. When the injection rate increases to a certain extent, the traffic rates are much higher and buffers become full frequently. Therefore, the flow-control feedback probability increases significantly. If the rate increases continually, the network will be gradually saturated and lead to the equilibrium of flow-control feedback probability.

In detail, as Fig.6(a) shows, when the injection rate lies between 0.008 and 0.032, the probability changes sharply; as the injection rate continues to increase, the probability remains steady. The analysis results track the simulated closely and the mean error is 8.65%.

Comparing Fig.6(b) with 6(a), we can see that, when the rate is 0.016, the probability reaches steady. The mean error is 7.83%.

With the different packet size, the probability of flow-control feedback slightly changes, but the saturation point and the overall trend are similar, as shown in Fig.6(a) and 3(c). In Fig.6(c), the mean error is 6.22%.

As Fig.6(d) shows, when the input buffer size is 8, the probability is smaller than the other situations with low injection rates. It increases sharply to be steady as the injection rate lies between 0.016 and 0.032. We conclude that input buffer size is one of the most important factors which impact the probability as expected. The mean error is of 8.78%.

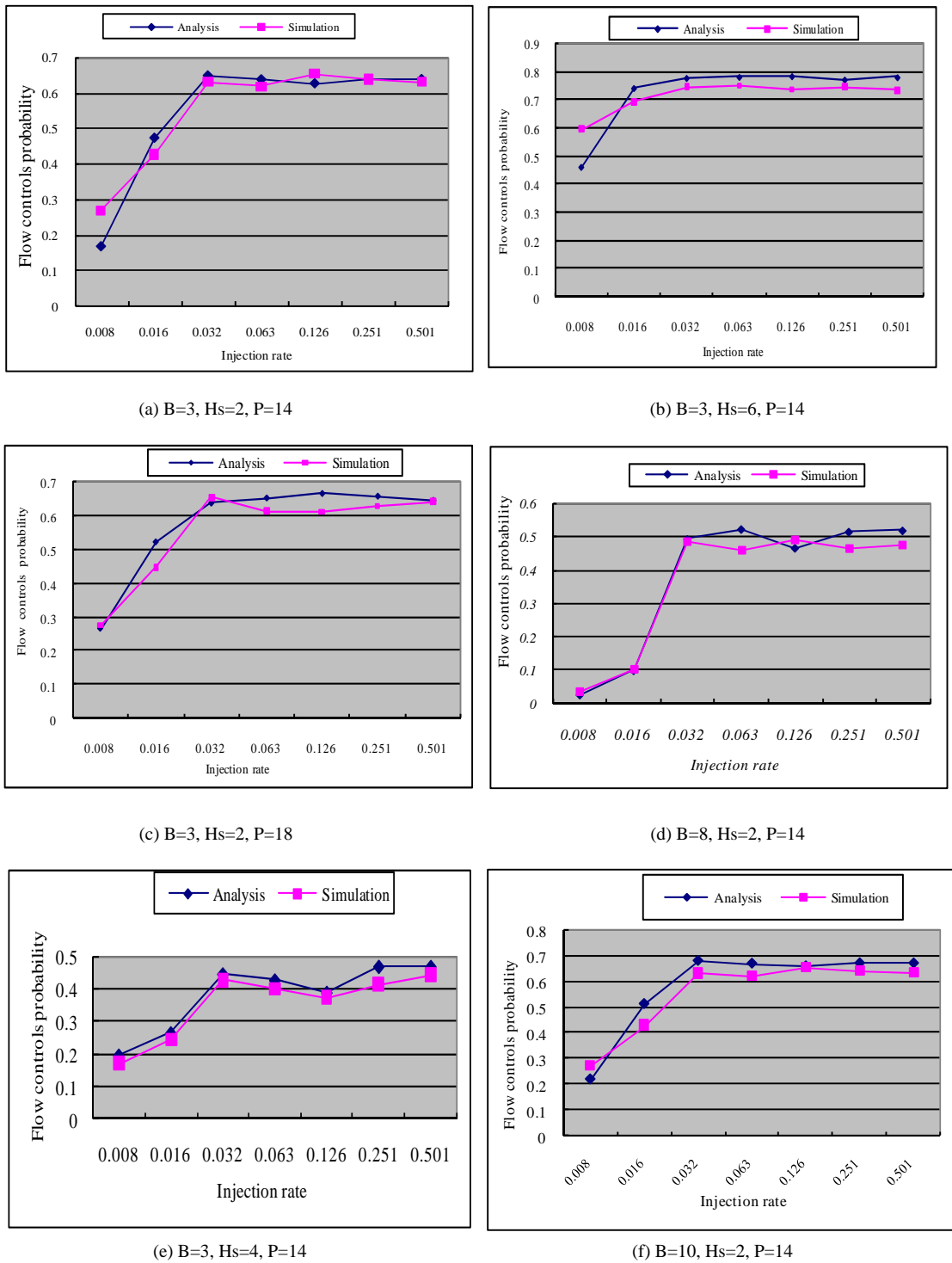


Figure 6. Flow-control feedback probability with different parameters

In summary, the computing method achieves a mean error of 7.87%, compared with the cycle-based simulation.

VII. CONCLUSION

A router model for NoC performance analysis is presented, which uses the M/D/1/B queuing theory to analyze various packet blocking-conditions. In addition, considering the effect of flow control to increase the blocking possibilities, a computing method of the flow control probability is proposed.

Experimental results show that the average error of the computing method for flow-control feedback probability is 7.87%.

In addition, we give the extended models for routers with multi-channel structure. Future work will refine the models and give their validation.

ACKNOWLEDGMENT

This research is supported by Chinese National Science Foundation under Grant 60773147, and Chinese National 863 High Technology Programs under Grant 2008AA01A204.

REFERENCES

- [1] W. J. Dally and B. Towles, "Route packets, not wires: on-chip interconnection network", DAC, 2001, pp. 684–689.
- [2] M. Sgroi, M. Sheets, A. Mihal, K. Keutzer, S. Malik, J. Rabaey and A. Sangiovanni-Vincentelli, "Addressing the System-on-a-Chip Interconnect Woes Through Communication-Based Design", DAC 2001, pp. 667–672.
- [3] J. Duato, S. Yalamanchili, L. Ni, "Interconnection Networks: An Engineering Approach", Morgan Kaufmann, 2002.
- [4] J. Kim, C.R. Das, "Hypercube communication delay with wormhole routing", IEEE Transactions on Computers C-43 (7) (1994) 806–814.
- [5] O. Lysne, "Towards a generic analytical model of wormhole routing networks", Microprocessors and Microsystems 21(7-8), 1998, 491–498.
- [6] W. J. Dally, "Performance analysis of k-ary n-cube interconnection networks", IEEE Transactions on Computers, vol. 39, pp. 775-785, 1990.
- [7] A. Khonsari, M. Ould-Khaoua and J. Ferguson, "A General Analytical Model of Adaptive Wormhole Routing in k-Ary n-Cube Interconnection Networks", SIMULATION SERIES, vol. 35, pp. 547-554, 2003.
- [8] J. Hu, U.Y. Ogras and R. Marculescu, "System-level buffer allocation for application-specific Networks-on-Chip router design", IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems 25(12), 2006, 2919-2933.
- [9] Y.U. Ogras and R. Marculescu, "Analytical router modeling for networks-on-chip performance analysis", In: Proceedings of Design, Automation and Test in Europe Conference (DATE'07), Acropolis, 2007, 1096-1101.
- [10] S. Foroutan, Y. Thonnart, R. Hersemeule and A. Jerraya, "An analytic method for evaluating network-on-chip performance", In: Proceedings of Design, Automation and Test in Europe Conference (DATE'10), Dresden, 2010, 1629-1632.
- [11] H. Sarbazi-Azad, M. Ould-Khaoua, L.M. Mackenzie, "An accurate analytical model of adaptive wormhole routing in k-ary n-cube interconnection networks", Performance Evaluation 43 (2,3) (2001) 165–179.
- [12] H. Sarbazi-Azad, M. Ould-Khaoua, L.M. Mackenzie, "On the performance of adaptive wormhole routing in the bi-directional torus network: A hotspot analysis", Microprocessors and Microsystems 25 (6) (2001) 277–285.
- [13] H. Sarbazi-Azad, M. Ould-Khaoua, L.M. Mackenzie, "Analytical modeling of wormhole-routed k-ary n-cubes in the presence of matrix-transpose traffic", Journal of Parallel and Distributed Computing 63 (2003) 396–409.
- [14] <http://nocs.stanford.edu/cgi-bin/trac.cgi/wiki/Resources/BlookSim>.
- [15] W.J. Dally and B. Towles, "Principles and Practices of Interconnection Networks", Morgan Kaufmann, San Francisco, 2004.
- [16] Ming-che Lai, Lei Gao, Nong Xiao, Zhiying Wang. "An accurate and efficient performance analysis approach based on queuing model for network on chip". Proceedings of International Conference on Computer Aided Design, pp. 563-570, 2009.

Youhui Zhang received the BSc and PhD degrees in computer science from Tsinghua University, China, in 1998 and 2002 respectively. He is currently an associate professor in the Department of Computer Science at Tsinghua University. His research interests include computer architecture, storage systems and high-performance computing. He is a member of the IEEE and the IEEE Computer Society.

Xiaoguo Dong is a graduate student in the Department of Mathematical Science and Computing Technology, Central South University, China. His research interests include mathematical modeling and analysis.

Siqing Gan is a professor and Chair in the Department of Mathematical Science and Computing Technology, Central South University, China. His research interests include numerical solution of differential equations, and scientific and engineering computing.

Weimin Zheng received the BSc and MSc degrees in computer science from Tsinghua University, China, in 1970 and 1982 respectively. Now he is a Professor in the Department of Computer Science at the University of Tsinghua, China. His research interests include high performance computing, network storage and parallel compiler. He is a member of the IEEE and the IEEE Computer Society.