



ELSEVIER

Contents lists available at ScienceDirect

INTEGRATION, the VLSI journal

journal homepage: www.elsevier.com/locate/vlsi

BiLink: A high performance NoC router architecture using bi-directional link with double data rate



Jingyang Zhu ^{a,*}, Zhiliang Qian ^b, Chi-Ying Tsui ^a

^a Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Hong Kong

^b Department of Micro- and Nano- Electronics, Shanghai Jiao Tong University, Shanghai, China

ARTICLE INFO

Article history:

Received 8 January 2015

Received in revised form

22 February 2016

Accepted 22 February 2016

Available online 2 March 2016

Keywords:

Network-on-Chip (NoC)

Bi-directional link

Double data rate

ABSTRACT

This paper presents a novel high performance Network-on-Chip (NoC) router architecture design using a bi-directional link with double data rate (BiLink). Ideally, it can provide as high as 2 times speed-up compared with the conventional NoC router. BiLink utilizes an extra link stage between routers and transmits two flits in one link per cycle using phase pipelining if both routers require to use the current link. To further increase the effective bandwidth, the direction of each link can be configured in every clock cycle to cater for different traffic loads from each side. Therefore, the data rate can be as high as 4 times compared with conventional NoC routers under uneven traffic. Centralized mode control scheme is implemented using a finite state machine (FSM) approach. Cycle-accurate simulations are carried out on both synthetic traffic patterns as well as real application benchmarks. Simulation results show that BiLink can provide as high as 90% and 250% speedup compared with conventional NoC routers for even and uneven traffic, respectively. 2X and 3X gains in throughput are obtained under even and uneven traffic, respectively, when compared with the conventional NoC router for the virtual channel flow control. The BiLink router architecture is synthesized using TSMC 65 nm process technology and it is shown that an area overhead of 28% over state-of-the-art bi-directional NoC is introduced while the critical path is about 9% higher than that of the conventional routers. Despite the overhead in critical path and power consumption, a 47.45% improvement of Energy-Delay-Product (EDP) is achieved by BiLink under high injection rate traffic.

© 2016 Elsevier B.V. All rights reserved.

1. Introduction

Network-on-Chip (NoC) has become a promising approach to solve the communication bottleneck in the modern many-core system-on-chip. With the potential deployment of many-core systems on new applications such as big data, artificial intelligence and deep machine learning, the NoC router requires to transfer a larger amount of communication data among processors. For example, the Google Brain project [14,8] uses 1000 machines to train a deep neural network. Each machine contains 16 cores on it and a subset of neural network will be mapped on each of them [8]. The requirement of the data bandwidth is high and uneven due to the interleaving of the feed-forward and back-propagation training phases. To address for the intensive bandwidth requirement of these applications, a higher throughput NoC

router architecture is essential and crucial for the next generation of many-core systems.

As the traffic pattern is usually uneven distributed among the network [13], self-reconfigurable router architectures have been proposed [13,6,20,2] to improve the NoC performance by adapting the direction of the links to the run time traffic conditions. A bi-directional NoC (BiNoC) router architecture was introduced in [13,6] to cater for the uneven traffic patterns. However, most of the emphasis on the existing reconfigurable NoC architecture has been focusing on optimizing the design of the router itself. The optimization of the interconnection between two neighboring routers is rarely touched. On the other hand, in the domain of communications, the introduction of network coding [1] provides an optimized way to use the channel bandwidth and achieves a significant improvement in the system throughput. Borrowing the concept of network coding, in [5], an extra coding unit was inserted between each pair of routers to enable the data transmission from both ends over a single physical channel simultaneously.

* Corresponding author.

E-mail addresses: jzhuak@ust.hk, eetsui@ust.hk (J. Zhu), qianzl@sjtu.edu.cn (Z. Qian).

In this work, to address the high bandwidth requirement for the next generation NoC architecture, we propose *BiLink*, a new NoC router architecture using bidirectional double data rate links. More specifically, in *BiLink*, a customized link stage is designed to transmit two flits over one physical channel in each cycle in a phase pipelined fashion. To further increase the effective bandwidth, the direction of each link can be configured to cater for the uneven distribution of the traffic loads. A centralized controller is implemented using a FSM to dynamically determine the operation mode to support *BiLink* transmission. In this way, data are transmitted in both the clock edges to maximize the potential throughput of the NoC router, leading to a better solution for the future data-intensive applications.

Cycle accurate simulations were executed to verify the performance improvement of the proposed *BiLink* architecture. Simulation results show that the proposed *BiLink* architecture can achieve 90% and 60% improvements in the saturation injection rate compared to Bi-directional (BiNoC) router architectures [13] for even and uneven traffic distributions, respectively. Furthermore, *BiLink* also has a 250% improvement over the conventional NoC router for the uneven traffic distribution. In summary, this work brings the following contributions:

- We combine the idea of self-reconfigurable router structures with a double data rate link for NoC and achieve a significant performance improvement through this joint optimization.
- We implement the proposed *BiLink* structure to verify the performance as well as the hardware overhead.
- We propose three variants of *BiLink* architecture and perform a thorough analysis on the performance and implementation tradeoff of these structures.

The remainder of the paper is organized as follows. In Section 2, we discuss the basic idea of the normal double data rate bidirectional link (*BiLink*) and analyze its timing issue. In Section 3, a self-reconfigurable direction control scheme, namely aggressive bidirectional link (*A-BiLink*) is proposed. In Section 4, the detailed hardware implementation of *BiLink* and *A-BiLink* are addressed. In addition, a new variant of *A-BiLink* which is more suitable for hardware implementation is presented in this section. Simulation and hardware synthesis results are shown in Section 5 and the related work is discussed in Section 6. Finally, Section 7 concludes this work.

2. Bidirectional link stage

To understand the basic principle behind the bidirectional link (*BiLink*), we will first discuss the data flow in *BiLink*. Then the related timing issues will be analyzed to show that *BiLink* can work properly under different timing constraints.

2.1. Motivation for exploring *BiLink*

In both uni-directional and bi-directional NoCs, the data transfer occupies the entire clock cycle. In this work, we propose

to further improve the throughput by allowing the transmission of two flits over the channel in every cycle. More specifically, we use both phases of the clock to transmit two different flits. In the first phase of clock cycle, the routers at both ends of the link send the flits to the link module in the middle of the link (shown in Fig. 1 (a)). Then, in the next phase, the link module sends the two flits to the corresponding destination routers (shown in Fig. 1(b)). Compared to the conventional transfer mode, the transfer data rate is doubled using the proposed *BiLink* scheme and it can transfer up to four flits between routers R1 and R2 in every clock cycle.

The main function of the intermediate link module is to isolate the flits from both routers at the two ends. For the link stage, two D Flip-Flops (DFFs) are required to store the data received from each side during the first half cycle. Moreover, two switches are used to control the direction of the data flow, in order to avoid overwriting the data originally stored in the registers. Fig. 2 (a) shows the hardware implementation of the link stage. When the clock phase is high, the switches S1 and S2 are open and the flits transmitted from both sides will be stored into these two DFFs in the link stage, respectively. Then, at the second phase of the clock cycle, S1 and S2 will be closed. The two DFFs will transmit the stored flits to the corresponding destination. For the router side, the output stage of each router has a similar structure to synchronize with the link stage. It sends flits at the first half clock cycle and receives flits at the second half as shown in Fig. 2(b).

2.2. Analysis of the timing constraints for *BiLink*

With the insertion of the link module, we need to analyze the impact on the timing of the overall system under reasonable clock skew and jitter assumptions.

First we investigate whether the insertion of the link stage will affect the clock frequency performance of the system. The datapath of a router consists of 2 parts, the inner pipeline stage and the link transfer stage. As will be shown in the simulation results in Section 5, the critical path of the inner pipeline stage of the router for the *BiLink* architecture is similar to that of the BiNoC. For the link transfer delay, the insertion of the link module will not cause extra delay. If the long wire delay of the link transfer is the critical path of the design, adding a link module in the middle breaks the long wire into half. Therefore the total delay of driving the long wire will be decreased instead and the overall critical path, which includes the clock to Q delay and the setup time of the DFF inserted in the link module, will be shortened.

We designed and layouted the link stage and the router's output stage in TSMC 65 nm process, and used it to drive different lengths of wires. We simulated the performance of the overall link transfer using HSPICE under a clock skew of 10% of the clock period [19]. The results show that the wire with a link stage is always better in terms of critical path performance than that without a link stage.

The hold time constraint of the link stage has also to be satisfied. The hold time of the DFF in the link module due to the datapath through the wire is easily satisfied because of the large delay of the long wire even under 10% positive clock skew. For the hold time requirement due to the inner loop with the link module

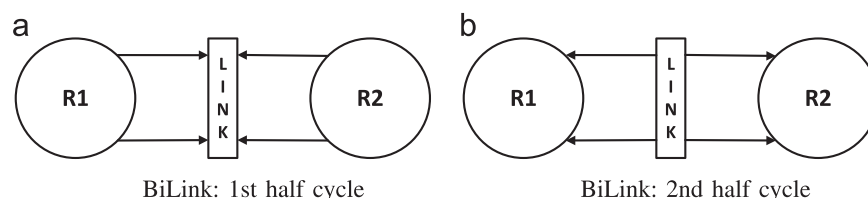


Fig. 1. Data transfer mode for *BiLink*.

Download English Version:

<https://daneshyari.com/en/article/542565>

Download Persian Version:

<https://daneshyari.com/article/542565>

[Daneshyari.com](https://daneshyari.com)