THREE DIMENSIONAL NETWORKS-ON-CHIP:

A PERFORMANCE EVALUATION

By

BRETT STANLEY FEERO

A thesis submitted in partial fulfillment of
the requirements for the degree of

MASTER OF SCIENCE IN COMPUTER ENGINEERING

WASHINGTON STATE UNIVERSITY
School of Electrical Engineering and Computer Science

MAY 2008

To the Faculty of Washington State University:

The members of the Committee appointed to examine the thesis of BRETT STANLEY FEERO find it satisfactory and recommend that it be accepted.

_____

Chair

_____

_____

# ACKNOWLEDGEMENT

I would first like to thank Katie, my beautiful wife, as well as my family. They have provided strong support throughout my education and have pushed me to achieve more than I thought possible. I would also like to thank Dr. Partha Pande, my advisor, for his passion, criticism, and guidance, and I also thank Washington State University, a great place to be. Go Cougs!

Lastly, I extend thanks to the referees of my conference and journal papers. Their suggestions and criticisms significantly increased the quality of my papers, and, as a result, this thesis.

THREE DIMENSIONAL NETWORKS-ON-CHIP:

A PERFORMANCE EVALUATION

Abstract

by Brett Stanley Feero, M.S.
Washington State University
May 2008

Chair: Partha P. Pande

The Network-on-Chip (NoC) paradigm has emerged as a revolutionary methodology for integrating a very high number of intellectual property (IP) blocks in a single die. The achievable performance benefit arising out of adopting NoCs is constrained by the performance limitation imposed by the metal wire, which is the physical realization of communication channels. With technology scaling, only depending on the material innovation will extend the lifetime of conventional interconnect systems a few technology generations. According to International Technology Roadmap for Semiconductors (ITRS) for the longer term, new interconnect paradigms are in need. The conventional two dimensional (2D) integrated circuit (IC) has limited floor-planning choices, and consequently it limits the performance enhancements arising out of NoC architectures. Three dimensional (3D) ICs are capable of achieving better performance, functionality, and packaging density compared to more traditional planar ICs. On the other hand, NoC is an enabling solution for integrating large numbers of embedded cores in a single die. 3D NoC architectures combine the benefits of these

two new domains to offer an unprecedented performance gain.

This thesis quantifies the performance of 3D NoC architectures. It demonstrates functionality in terms of throughput, latency, energy dissipation, and wiring area overhead. It also addresses the temperature concerns that are apparent in 3D integrated circuits in general as well as many emerging 2D applications, showing that the characteristics of 3D NoCs limit what would otherwise be a dramatic increase in temperature, and in a certain case, even reduce temperature.

TABLE OF CONTENTS

# LIST OF FIGURES

Page

**CHAPTER ONE**

**INTRODUCTION**


The current trend in System-on-Chip (SoC) design in the ultra deep sub-micron (UDSM) regime and beyond is to integrate a huge number of functional and storage blocks in a single die [1]. The possibility of this enormous degree of integration gives rise to new challenges in designing the interconnection infrastructure for these big SoCs. Extrapolating from the existing CMOS scaling trends, traditional on-chip interconnect systems have been projected to be limited in their ability to meet the performance needs of SoCs at the UDSM technology nodes and beyond [2]. This limit stems primarily from global interconnect delay significantly exceeding that of gate delays. While copper and low-$k$ dielectrics have been introduced to decrease the global interconnect delay, they only extend the lifetime of conventional interconnect systems a few technology generations. According to the International Technology Roadmap for Semiconductors (ITRS) [2], for the longer term, material innovation with traditional scaling will no longer satisfy the performance requirements. New interconnect paradigms are in need. Continued progress of interconnect performance will require employing approaches that introduce materials and structures beyond the conventional metal/dielectric system, and one of the promising approaches is 3D integration. Shown in Figure 1, three-dimensional (3D) ICs, which contain multiple layers of active devices, have the potential for enhancing system performance [3] [4] [5] [6]. According to [3], three-dimensional ICs allow for performance enhancements even in the absence of scaling. A clear way to

reduce the burden of high frequency signal propagation across monolithic ICs is to reduce the line length needed, and this can be done by employing stacking of active devices using 3D interconnects. Here, the multiple layers of active devices are separated by a few tens of micrometers. Consequently, 3D interconnects allow communication among these active devices with smaller distances required for signal propagation.

Three-dimensional ICs will have a significant impact on the design of multi-core SoCs. Recently, Networks-on-Chip (NoCs) have emerged as an effective methodology for designing big multi-core SoCs [7] [8]. However, the conventional two dimensional (2D) IC has limited floor-planning choices and, consequently, limits the performance enhancements arising out of NoC architectures. The performance improvement arising from the architectural advantages of NoCs will be significantly enhanced if 3D ICs are adopted as the basic fabrication methodology. The amalgamation of two emerging paradigms, namely NoCs in a 3D IC environment, allows for the creation of new structures that enable significant performance enhancements over more traditional solutions. With freedom in the third dimension, on-chip network architectures that were



Figure 1. 3DIC from a SOI Process

impossible or prohibitive due to wiring constraints in planar ICs are now possible [9] [10].

However, 3D ICs are not without limitations. Thermal effects are already impacting interconnect and device reliability in 2D circuits [11]. Due to the reduction of chip size in a 3D implementation, 3D integrated circuits exhibit a profound increase in power density. Consequently, increases in heat dissipation will give rise to circuit degradation and chip cracking, among other side-effects [12]. As a result, there is a real need to keep the temperature low for reliable circuit operation. Furthermore, in ICs implementing NoCs, the interconnect structure dissipates a large percentage of energy. In certain applications [13], this percentage has been shown to approach 50%. As a result, the interconnection network has a significant contribution to the thermal performance of 3D NoCs.

This thesis characterizes the performance of multiple 3D NoC architectures in the presence of realistic traffic patterns through cycle-accurate simulation and establishes the performance benchmark and related design trade-offs. The metrics of throughput, latency, energy dissipation, area, and temperature are each evaluated.

This thesis is organized as follows. Chapter 2 covers previous research and developments pertinent to this study. It covers other works which evaluate three-dimensional integrated circuits, various types of three-dimensional NoCs, and various methods of controlling the thermal performance of integrated circuits. Chapter 3 is an introduction to the different 3D NoC architectures evaluated in this thesis. Both mesh-based and tree-based topologies are considered in this study. Chapter 4 reviews the

performance metrics used to evaluate the network architectures and details the experimental analysis. In this chapter, the various NoCs are evaluated in the presence of real traffic patterns for throughput characteristics, latency, energy dissipation, and silicon area. The effects of temperature are discussed in detail in chapter 5. Finally, chapter 6 concludes this thesis with an overview and future research directions.

# CHAPTER TWO

# RELATED WORK

Current SoCs are implemented predominantly following 2D architectures. However, the emergence of 3D ICs will present a fundamental change. Topol et al., in [3], describe, in detail, the challenges of manufacturing in a 3D IC process. They show that 3D ICs are capable of improvements in power, noise, logical span, density, performance, and functionality. One major advantage of the 3D IC paradigm is that it allows for the integration of "dissimilar technologies", e.g. memory, analog, MEMS, etc. in a single die. The paper describes the benefits and drawbacks of different fabrication methods including face-to-face bonding and face-to-back bonding. With their most sophisticated SOI face-to-back process, the via pitch is minimized, at 0.4 µm with a separation of 2 µm between layers of SOI devices [3].

Jacob, et al. [14] propose using 3D ICs to improve the performance of microprocessors by forming a processor-memory stack. They show that the integration of processor and memory in a stack enables a large increase in performance. In particular, 3D integration enables the use of very wide buses (>1024 bits) for vertical communication. In addition to ultra wide buses, a stack provides a very short distance between processor and memory, decreasing memory access times considerably.

In [9], 3D ICs were proposed to improve performance of chip multi-processors. Drawing upon 3D IC research, they chose a hybridization of busses and networks to provide the interconnect fabric between CPUs and L2 caches. The performance of this

fusion of NoC and bus architectures was evaluated using standard CPU benchmarks. However, this analysis pertains only to chip multiprocessors and does not consider the use of three-dimensional network structures for application-specific SoCs. Three dimensional NoCs are analyzed in terms of temperature in [15]. Pavlidis, et al., in [10], compared 2D MESH structures with their 3D counterparts by analyzing the zero-load latency and power consumption of each network. This is an evaluation that shows some of the advantages of 3D NoCs, but it neither applies any real traffic pattern, nor does it measure other relevant performance metrics. This thesis aims to address these concerns by applying real traffic patterns in a cycle-accurate simulation, and measuring performance through established metrics.

The thermal consequences of conventional 2D NoCs were discussed in [16]. This research introduces a runtime solution to temperature mitigation by throttling traffic through the routers. The efficient algorithm presented offers reactive and proactive routing in addition to the distributed throttling of traffic. In the specific application shown in this paper, their algorithm reduces temperature by 10 °C, while incurring throughput degradation of less than 1% and latency degradation of less than 1.2%. Although this work does not pertain to 3D NoCs, the algorithm can be adapted for this scenario.

Another approach to reduce temperature in an integrated circuit is the inclusion of thermal vias, thus improving its thermal conductivity. This is of particular importance in 3D ICs. Thermal vias can improve the thermal conductivity of any 3D IC, network-on-chip applications inclusive. In [17], a thermal-via allocation algorithm is presented. This work improves on previous thermal via algorithms by considering the spatial and

temporal variations of energy dissipation. In addition to showing significant speedup compared to previous algorithms, it reduces the area overhead from these thermal vias by one half.

A third method of mitigating temperature concerns is thermally-intelligent placement of processing elements. In [15], the author extends previous work covering the optimal placement and mapping of on-chip components of 2D NoCs to a three-dimensional environment. This work proposes the use of genetic algorithms to enable optimum placement of processing elements in three-dimensional networks-on-chip.

Each of [16], [17], and [15] provide unique approaches to mitigating the temperature problem. One of the principal characteristics of 3D NoCs is that they dissipate lower communication energy compared to 2D implementations, and eventually this reduction has a positive effect on the temperature issues arising from increased power density in nascent 3D NoCs. However, these effects have not been physically quantified on an architectural level. This thesis addresses that void by characterizing the thermal performance of an array of 3D NoC topologies.

**2.1 Conclusions**

This chapter has introduced previous works on three-dimensional integrated circuits, networks-on-chip, and thermal concerns. The preceding works have been invaluable to this thesis, and they have provided a foundation for what are presented. This thesis aims to address the limitations of previous three-dimensional NoC research by expanding upon the limitations introducing more architectures, using real traffic patterns,

and applying the concepts of other works, namely the inclusions of ultra-wide buses and

the introduction of thermal analysis.

# CHAPTER THREE

# 3D NOC ARCHITECTURES

Enabling design in the vertical dimension permits a large degree of freedom in choosing an on-chip network topology. Due to wire-length constraints and layout complications, the more conventional two-dimensional integrated circuits have placed limitations on the types of network structures that are possible. With the advent of 3D ICs, a wide range of on-chip network structures that were not explored earlier are being considered [9] [10]. This paper investigates five different topologies in 3D space and compares them with three well-known NoC architectures from 2D implementations. This thesis considers a SoC with a 400mm$^2$ floor plan and 64 functional IP blocks. This system size was selected to reflect the state of the art of emerging SoCs. At ISSCC 2007, design of an 80-core processor arranged in an 8x10 regular grid built on fundamental NoC concepts was demonstrated [18]. Therefore, the system size assumed in this work is representative of the current trends. IP blocks for a 3D SoC are mapped onto four 10mm×10mm layers, in order to occupy the same total area as a single-layer, 20mm×20mm layout.

## 3.1. Mesh-Based Networks

One of the well-known 2D NoC architectures is the 2D Mesh as shown in Figure 2a. This architecture consists of an *m×n* mesh of switches interconnecting IP blocks placed along with them. It is known for its regular structure and short interswitch wires.
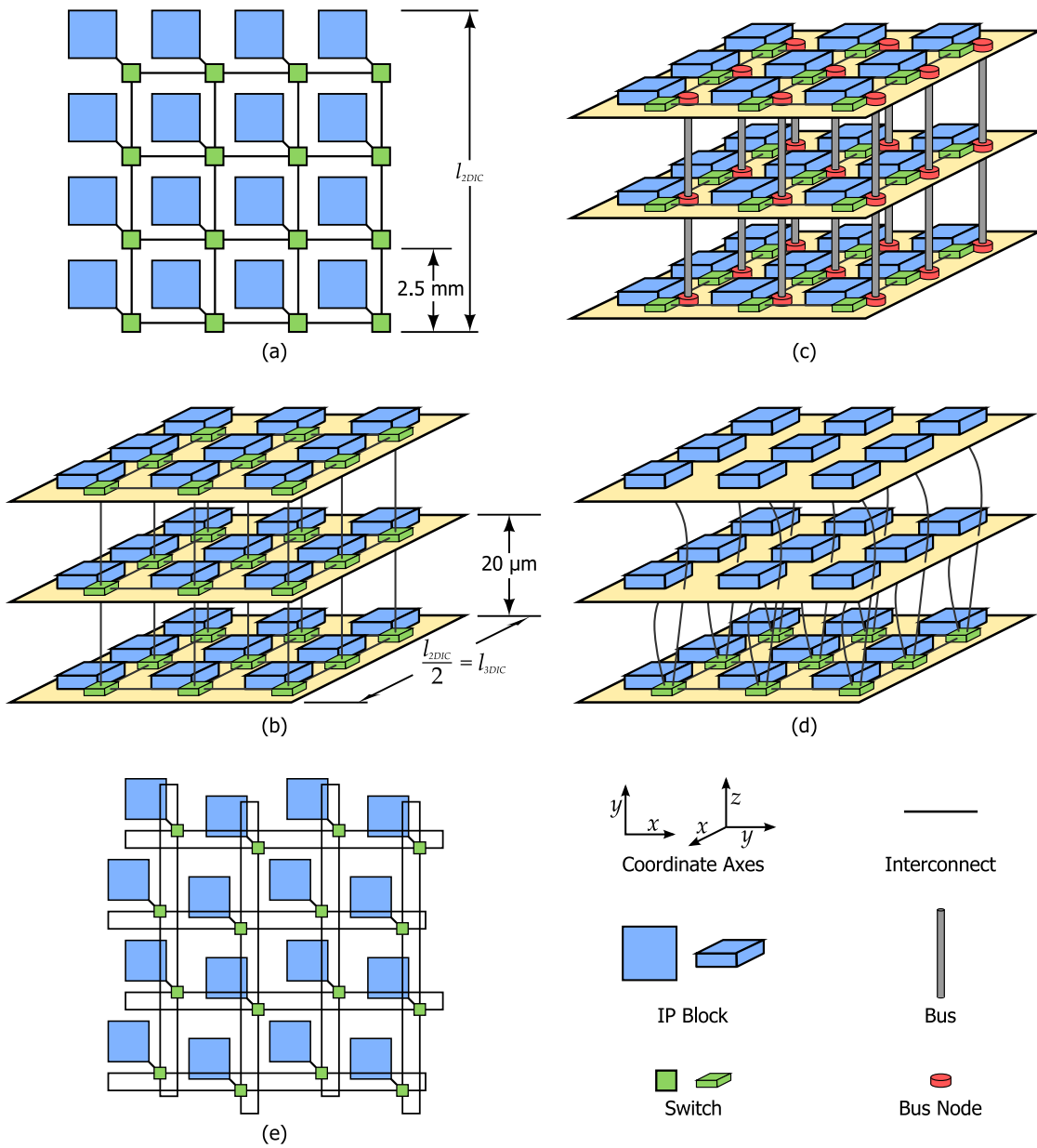
Figure 2. Mesh-based NoC architectures: (a) 2D Mesh, (b) 3D Mesh, (c) Stacked Mesh, and (d) Ciliated 3D Mesh

From this structure, a variety of three-dimensional topologies can be derived. The straightforward extension of this popular planar structure is the 3D Mesh. Figure 2b shows an example of 3D Mesh NoC. It employs 7-port switches: one port to the IP block, one each to switches above and below, and one in each cardinal direction (North, South, East, and West), as shown in Figure 3a. A second derivation, 3D Stacked Mesh (Figure 2c), takes advantage of the short inter-layer distances that are characteristics of a 3D IC, which can be around 20 μm [3]. The 3D Stacked Mesh architecture is a hybrid between a packet-switched network and a bus. It integrates multiple layers of 2D Mesh networks by connecting them with a bus spanning the entire vertical distance of the chip. As the distance between the individual 2D layers in 3D IC is extremely small, the overall length of the bus is also small, making it a suitable choice for communicating in the $z$-dimension [9]. Furthermore, each bus has only a small number of nodes (i.e. equal to the number of layers of silicon), keeping the overall capacitance on the bus small and greatly simplifying bus arbitration. For consistency with [9], this analysis considers the use of a dynamic, time-division multiple-access (dTDMA) bus, although any other type of bus may be used as well. A switch in a 3D Stacked Mesh network has, at most, 6 ports: one to the IP, one to the bus, and four for the cardinal directions (Figure 3b). Additionally, it is possible to utilize ultra wide buses similar to the approach introduced in [14] to implement cost-effective, high-bandwidth communication between layers.

A third method of constructing a 3D NoC is by adding layers of functional IP blocks and restricting the switches to one layer or a small number of layers. With this in mind, this thesis introduces a new architecture, 3D Ciliated Mesh. This structure is

11

Figure 3. Switches for mesh-based NoCs: (a) 3D Mesh, (b) Stacked Mesh, and (c) Ciliated 3D Mesh

essentially a 3D Mesh network with multiple IP blocks per switch. The 3D Ciliated Mesh is a 4×4×2 3D mesh-based network with 2 IPs per switch, where the two functional IP blocks occupy, more or less, the same footprint, but reside at different layers. This is shown in Figure 2d. In a Ciliated 3D Mesh network, each switch contains seven ports (one for each cardinal direction, one either up or down, and one to each of two IP blocks) as shown in Figure 3c. This architecture will clearly exhibit lower overall bandwidth than a complete 3D Mesh due to multiple IP blocks per switch and reduced connectivity; however, chapter 4 will show that this type of network offers an advantage in terms of energy dissipation, especially in the presence of specific traffic patterns.

It is important to note that each mesh-based network introduced in this section can be easily translated into a toroidal structure. Toroidal structures differ in that the switches at the edges wrap around to the opposite side. This ensures that all switches have equal numbers of ports. However, this leads to long wrap-around wires, so the Folded Torus architecture (Figure 2e), in which all wires are the same length, was designed as a

mitigating solution [19]. Furthermore, toroidal networks they are advantageous for many parallel processing algorithms, like Fast Fourier Transforms for example, since they are implementations of hypercubes. For instance, a 4×4 2D torus is equivalent to a 2×2×2×2 4D hypercube, and a 4×4×4 3D torus is equivalent to a 2×2×2×2×2×2 6D hypercube. Three-dimensional folded tori are easily extensible from the 3D Mesh, 3D Stacked Mesh, and 3D Ciliated Mesh structures in Figure 2.

### 3.2. Tree-Based Networks

Two types of tree-based interconnection networks that have been considered for network-on-chip applications are Butterfly Fat Tree (BFT) [20], [21] and the generic Fat Tree, or SPIN [22]. This paper endeavors to quantify the enhancements achieved when these networks are instantiated in a 3D IC environment. Unlike the work with mesh-based NoCs, this thesis does not propose any new topologies for tree-based systems. Instead, it investigates the achievable performance benefits by instantiating already-existing tree-based NoC topologies in a 3D environment.

The considered BFT topology is shown in Figure 4a. For a 64-IP SoC, a BFT network will contain 28 switches. Each switch (Figure 5a) in a Butterfly Fat Tree network consists of 6 ports, one to each of four child nodes and two to parent nodes, with the exception of the switches at the topmost layer. When mapped to a 2D structure the longest inter-switch wire length for a BFT-based NoC is $l_{2DIC}/2$, where $l_{2DIC}$ is the die length on one side [21] [23]. If the NoC is spread over a 20mm×20mm die, then the longest inter-switch wire is 10 mm [23], as shown in Figure 4c. Yet, when the same BFT

13

Figure 4. Tree architectures: (a) Butterfly Fat Tree, (b) SPIN, (c) 2D BFT Floorplan, (d) 3D BFT Floorplan for the first two layers, and (e) the first three layers of a 3D BFT Floorplan as seen in elevation view

14

Figure 5. Switches for tree networks: (a) BFT and (b) SPIN

network is mapped onto a four-layer 3D SoC, wire routing becomes simpler, and the longest inter-switch wire length is reduced by at least a factor of two, as can be seen in Figure 4d. This will lead to reduced energy dissipation as well as less area overhead. The fat tree topology of Figure 4b will have the same advantages when mapped on to a 3D IC as the BFT.

### 3.3 Conclusions

This chapter has introduced the various 3D NoC architectures that are evaluated in this thesis. Three-dimensional integrated circuits allow the creation of new topologies, and they allow the utilization old topologies in such a way that they take advantage of some of the inherent benefits of 3D ICs. This chapter has shown how the traditional 2D Mesh can be expanded in the third dimension with the creation of 3D Mesh, 3D Stacked Mesh, and 3D Ciliated Mesh topologies. It has introduced how the concept of ultra-wide buses can apply to the bus-NoC hybrid 3D Stacked Mesh architecture as well. These

three-dimensional structures take advantage of more efficient topology to produce savings in terms of energy and improvements in terms of throughput. Lastly, two traditional tree-based topologies were introduced, and it was shown that without creating any new tree-based topologies, both the Fat Tree and Butterfly Fat Tree architectures can be adapted in three-dimensions in such a way that energy is saved and area is reduced.

# CHAPTER FOUR

# PERFORMANCE EVALUATION

## 4.1 Performance Metrics

In order to properly analyze the various 3D network-on-chip topologies, a standard set of metrics must be used [24]. Wormhole routing [25] is assumed as the data transport mechanism where the packet is divided into fixed length flow control units or flits. The header flit holds the routing and control information. It establishes a path, and subsequent payload or body flits follow that path. This comparative analysis focuses on the four established benchmarks [24] of throughput, latency, energy, and area overhead.

Throughput is a metric that quantifies the rate in which message traffic can be sent across a communication fabric. It is defined as the average number of flits arriving per IP block per clock cycle, so the maximum throughput of a system is directly related to the peak data rate that a system can sustain. For purposes of a message-passing system, throughput $T$ is given by the equation

$$T = \frac{(Total\ Messages\ Completed) \times (Message\ Length)}{(Number\ of\ IP\ Blocks) \times (Time)}. \tag{1}$$

*Total Messages Completed* are the number of messages which successfully traverse the network from source to destination. *Message Length* refers to the number of flits a message consists of, and *Number of IP Blocks* signifies the number of intellectual property units that send data over the network. *Time* is length of time in clock cycles between the generation of the first packet and the reception of the last. It can be seen that

throughput is measured in flits/IP block/cycle, where a throughput of 1 signifies that every IP block is accepting a flit in each clock cycle. Accordingly, throughput is a measure of the maximum amount of sustainable traffic. Throughput will be dependent on a number of parameters including the number of links in the architecture, the average hop count, the number of ports per switch, and injection load. Injection load is measured by the number of flits injected in to the network per IP block per cycle. Consequently, it has the same unit as the throughput, and an injection load of 1 signifies that every IP block is injecting a flit in each clock cycle.

Next, latency refers to the length of time elapsed between the injection of a message header at the source node and the reception of the tail flit at the destination. Latency is defined as the time in clock cycles elapsed from the transfer of the header flit by the source IP to the acceptance of the tail flit by the destination IP block. Latency is characterized by three delays: sender overhead, transport latency, and receiver overhead.

$$L_i = L_{sender} + L_{transport} + L_{receiver} \tag{2}$$

Flits must traverse a network while traveling from source to destination. With different routing algorithms and switch architectures, each packet will experience a unique latency. As a result, network topologies will be compared by average latency. Let $P$ be the number of packets received in a given time period, and let $L_i$ be the latency of the $i$th packet. Average latency is therefore given by the equation:

$$L_{avg} = \frac{\sum_{i=1}^{P} L_i}{P} . \tag{3}$$

Additionally, the transport of messages across a network leads to a quantifiable

amount of energy dissipation. Activity in the logic gates of the network switches as well as the charging and discharging of interconnection wires lead to the consumption of energy. This thesis examines two types of energy: energy per cycle and packet energy. Cycle energy is defined as the amount (in Joules) of energy dissipated by the entire network in one clock cycle. On the other hand, packet energy is defined as the amount of energy incurred by a single packet as it traverses the network from source to destination over many clock cycles. It will be shown that each of these types of energy reveals unique information about the behavior of the varying network architectures.

Lastly, the amount of silicon area used by an interconnection network is a necessary consideration. As the network switches form an integral part of the infrastructure, it is important to determine the amount of relative silicon area they consume. Additionally, area overhead arising from layer-to-layer vias, inter-switch wires, and buffers incurred by relatively longer wires need to be considered. The evaluation of area in this thesis includes each form of area overhead.

## 4.1 Performance Analysis of 3D Mesh-Based NoCs

In this section, the performance of the 3D mesh-based NoC architectures is analyzed in terms of the parameters mentioned above: throughput, latency, energy dissipation, and area overhead.

Throughput is given in the number of accepted flits per IP per cycle. This metric, therefore, is closely related to the maximum amount of sustainable traffic in a certain network type. Any improvements in throughput in 3D networks are principally related to

19

two factors: the number of physical links and the average number of hops.

In general, for a mesh-based NoC, the number of links is given as follows:

$$links = N_1 N_2 (N_3 - 1) + N_1 N_3 (N_2 - 1) + N_2 N_3 (N_1 - 1), \tag{4}$$

where $N_i$ represents the number of switches in the $i^{th}$ dimension. For instance, in an 8×8 2D Mesh NoC, this yield 112 links. In a 4×4×4 3D Mesh NoC, the number of links turns out to be 144. With a greater number of links, a 3D Mesh network, for example, is able to contain a greater number of flits and therefore transmit a greater number of messages.

However, only considering the number of links will not characterize the overall throughput of a network. The average hop count also has a definitive effect on throughput. Following [10], the average number of hops in a mesh-based NoC is given by

$$hops_{Mesh} = \frac{n_1 n_2 n_3 (n_1 + n_2 + n_3) - n_3 (n_1 + n_2) - n_1 n_2}{3(n_1 n_2 n_3 - 1)}, \tag{5}$$

where $n_i$ is the number of nodes in the $i^{th}$ dimension. This equation applies both to the 4×4×4 3D Mesh and 4×4×2 3D Ciliated Mesh networks. The number of hops for the 3D Stacked Mesh is equal to

$$hops_{Stacked} = \frac{n_1 + n_2}{3} + \frac{n_3 - 1}{n_3}. \tag{6}$$

For the 4×4×4 3D Mesh and 8×8 2D Mesh, average hop counts are 3.81 and 5.33, respectively. There are 40% more hops in the 2D Mesh compared to that in 3D Mesh. Consequently, flits in the 3D Mesh needs to traverse fewer stages between a pair of source and destination than the 2D counterpart. As a result of this, a corresponding increase in throughput is expected. A lower average hop count will also allow more flits to be transmitted through the network. With a lower hop count, a wormhole-routed packet

will utilize fewer links, thus leaving more room to increase the maximum sustainable traffic.

Transport latency, like throughput, is also affected by average hop count. It is also affected heavily by the number of links and the injection load. In 3D architectures, a decrease in latency is expected due to a lower hop count and an increased number of links.

In the System-on-Chip realm, energy dissipation characteristics of the interconnect structures are crucial, as the interconnect fabric can consume a significant portion of the overall energy budget [15]. The energy dissipation in a NoC depends on the energy dissipated by the switch blocks and the inter-switch wire segments. Both of these factors depend on the network architecture. Additionally, the injection load has a significant contribution as it is the cause for any activity in the switches and inter-switch wires. Intuitively, it is clear that with more packets traversing the network, power will increase. This is why packet energy is an important attribute for characterizing NoC structures. The energy dissipated per flit per hop is given by

$$E_{hop} = E_{switch} + E_{wire} , \tag{7}$$

where $E_{switch}$ and $E_{wire}$ are the energy dissipated by each switch and inter-switch wire segments respectively. The energy of a packet of length $n$ flits that completes $h$ hops is given by

$$E_{packet} = n \sum_{j=1}^{h} E_{hop,j} . \tag{8}$$

From this, a formula for packet energy can be realized. If $P$ packets are transmitted then

the average energy dissipated per packet is given as

$$\overline{E_{packet}} = \frac{\sum_{i=1}^{P} E_{packet,i}}{P} = \frac{\sum_{i=1}^{P} \left( n_i \sum_{j=1}^{h_i} E_{hop,j} \right)}{P} \, . \tag{9}$$

Now, it is clear that a strong correlation exists between packet energy and the number of hops from source to destination. Consequently, a network topology that exhibits smaller hop counts will also exhibit correspondingly lower packet energy. As all 3D mesh-based NoC architectures exhibit a lower hop count they should also dissipate less energy per packet.

Lastly, the area overhead must be analyzed for mesh-based NoCs. Area overhead for a NoC includes switch overhead and wiring overhead. Switch area is affected by the overall number of switches and the area per switch, which is highly correlated to the number of ports. Since all 3D mesh-based NoCs have more ports, the area per switch will increase. However, the ciliated structure has a reduced number of switches, which should significantly reduce the overall switch area. For 3D NoCs in general, wiring overhead includes the interlayer via footprint in addition to the area incurred by horizontal and vertical wiring. The addition of interlayer vias and their corresponding area overhead is a characteristic that is unique to three-dimensional ICs, and it is included in the area overhead calculations presented later.

Wire overhead is reduced when moving to a 3DIC. However, this is not due to reductions in the length of most interswitch wires in the case of mesh-based NoCs. Horizontal wire length is given by $l_{IC}/n_{side}$, where $n_{side}$ represents the number of IPs in one dimension of the IC and $l_{IC}$ is the die length on one side as shown earlier in Figure 2a and

Figure 2b. For the 8×8 2D Mesh, this evaluates to 20mm/8 or 2.5mm, and for all 3D mesh-based architectures, the expression evaluates to 10mm/4, also 2.5mm. With this in mind, reductions in wire overhead come from the interlayer wires. The 3D structures have a reduced number of horizontal links due to the presence of interlayer wires. These interlayer wires are very small and hence, they are the source of wire overhead savings in mesh-based 3D NoCs.

**4.2 Performance Analysis of 3D Tree-Based NoCs**

Unlike the previous discussion pertaining to mesh-based NoCs, the tree-based networks considered for 3D implementations have identical topologies to their 2D counterparts. The only variable is the inter-switch wire length. As a result, there are significant improvements both in terms of energy and area overhead.

In 2D space, the longest inter-switch wire length in a BFT or SPIN network is equal to $l_{2DIC}/2$ [21] [23], where $l_{2DIC}$ is the die length on one side. This inter-switch wire length corresponds to the top-most level of the tree. In a 3D IC, however, this changes significantly. For instance, as shown in Figure 4d and Figure 4e, the longest wire length for 3D, tree-based NoC is equal to the length of horizontal travel in addition to the length of the vertical via. Considering a 20mm×20mm 2D die, the longest inter-switch wire length is equal to 10mm, whereas with a 10mm×10mm stack of four layers, the maximum wire length is equal to the sum of $l_{3DIC}/4$, or 2.5mm, and the span of two layers, 40μm. This is almost a factor-of-4 reduction compared to 2D implementations. Similarly, mid-level wire lengths are reduced by a factor of 2. As a result, this reduction

23

Table 1. Inter-Switch Wire Lengths in 3D tree-based NoCs

|  | 2D NoC | 4-layer 3D NoC |
|---|---|---|
| 1$^{st}$ Level | $\leq l/8 = 2.5$ mm | $\leq l/4 = 2.5$ mm |
| 2$^{nd}$ Level | $l/4 = 5$ mm | $l/4 = 2.5$ mm |
| 3$^{rd}$ Level | $l/2 = 10$ mm | $l/4 = 2.5$ mm |

in wire length, shown in Table 1, causes a significant reduction in energy.

In addition to benefits in terms of energy, 3D ICs effect area improvements for tree-based NoCs. Again, as with energy, area gains pertain only to the inter-switch wire segments; there is neither a change in the number of switches nor in the design of the switch.

As with the 3D mesh-based NoCs, wire overhead in a 3D tree-based NoC consists of the horizontal wiring in addition to the area incurred by the vertical wires and vias. Also, the longer inter-switch wires, which are characteristics of 2D tree-based NoCs, require repeaters, and this is taken into account. For a Butterfly Fat Tree, the number of wires in an arbitrary tree level $l$ as defined in [20] is

$$wires_{\text{layer } l} = w_{link} \cdot \left\lceil \frac{N}{2^{l-1}} \right\rceil,$$
(10)

where $N$ is the number of IP blocks and $w_{link}$ is the link width in bits. For a generic Fat Tree, the number of wires in a tree level $l$ is given by

$$wires_{\text{layer } l} = w_{link} \cdot N .$$
(11)

For instance, in a 64-IP BFT network with 32-bit wide bi-directional interswitch links, there are 2048 wires in the first level, 1024 wires in the second level, and 512 wires in the

third. Similarly, a 64-IP Fat Tree will have 2048 wires in every level.

## 4.3 Simulation Methodology

To model performance of different NoC structures, a cycle-accurate network simulator is employed that can also simulate dTDMA buses. The simulator is flit-driven and uses wormhole routing. In this work, a self-similar injection process [24] [26] [27] [28] is assumed. This type of traffic has been observed in the bursty traffic typical of on-chip modules in MPEG-2 video applications [28], as well as various other networking applications [27]. It has been shown to closely model real traffic [28]. In terms of spatial distribution it is capable of producing both uniform and localized traffic patterns for injected packets. In order to acquire energy and area characteristics, the network

Table 2. Wire Delays

| Wire Type | Wire Length | Delay (ps) | Architectures Used |
|---|---|---|---|
| Interlayer | 20 µm | 16 | all 3D mesh-based |
| Vertical Bus | 60 µm | 110/450** | 3D Stacked Mesh |
| Horizontal | 2.5 mm | 219 | mesh-based, 2D tree-based |
| Horizontal + Interlayer | 2.54 mm | 231 | all 3D tree-based |
| Horizontal | 5 mm | 436* | Mid-level in all 2D tree-based |
| Horizontal | 10 mm | 550* | Top-level in all 2D tree-based |
| | | *Repeaters Necessary | **Bus Arbitration Included |

Table 3. Architectural Parameters

| Topology | Port Count | Switch Area (mm²) | Switch Static Energy (pJ) | Longest Wire Delay (ps) |
|---|---|---|---|---|
| 2D Mesh | 5 | 0.0924 | 65.3 | 219 |
| 3D Mesh | 7 | 0.1385 | 91.4 | 219 |
| 3D Stacked Mesh | 6 (+ bus arbitration) | 0.1225 | 81.3 | 219 |
| Ciliated 3D Mesh | 7 | 0.1346 | 91.2 | 219 |
| 2D BFT | 6 | 0.1155 | 78.3 | 550 |
| 3D BFT | 6 | 0.1155 | 78.3 | 231 |
| 2D Fat Tree | 8 | 0.1616 | 104.5 | 550 |
| 3D Fat Tree | 8 | 0.1616 | 104.5 | 231 |

switches, dTDMA arbiter, and FIFO buffers were modeled in VHDL. The network switches were designed in such a way that their delay can be constrained within the limit of one clock cycle. The clock cycle is assumed to be equal to 15FO4 (fan-out-of 4) delay units. With the 90nm standard cell library from CMP [29], this corresponds to a clock frequency of 1.67 GHz. As the switches were designed with differing numbers of ports, their delays vary with one another. However, it was important to ensure that all the delay numbers were kept within the 15FO4 timing constraint. Consistent with [23], the longest delays were in the 2D/3D Fat Tree switches as they had the highest number of ports. Yet, even it can be run with a clock frequency of 11FO4, well within the 15FO4 limit. To have a consistent comparison, all the switches were run with a 15FO4 clock.

Similarly, all interswitch wire delays must hold within the same constraints. As shown in Table 2, wire RC delays remain within the clock period of 600ps [29]. For

Stacked Mesh, even considering the bus arbitration, the delay is constrained within one clock cycle. For the vertical wires, the via resistance and capacitance are included in the analysis. Thus, all network architectures are able to run at the same clock frequency of 1.67 GHz. Additional architectural parameters for each topology are shown in Table 3.

Although the simulator is capable of running with an arbitrary specification, each switch was designed with 4 virtual channels per port and 2-flit-deep virtual channel buffers as discussed in [24]. Synopsys Design Vision was used to synthesize the hardware description using a 90nm standard cell library from CMP [29], and Synopsys PrimePower was used to gather energy dissipation statistics. To calculate $E_{switch}$ and $E_{wire}$ from (7), the methodology discussed in [24] is followed. The energy dissipated by each switch, $E_{switch}$, is determined by running its gate-level netlist through Synopsys PrimePower using large sets of input data patterns. In order to determine the interconnect energy, $E_{interconnect}$, the interconnects' capacitance is estimated, taking into account each inter-switch wire's specific layout, by the following expression [24]:

$$C_{interconnect} = C_{wire} \cdot w_{a+1;a} + n \cdot m \cdot (C_G + C_J), \tag{12}$$

where $C_{wire}$ represents the capacitance per unit length of the wire, $w_{a+1;a}$ is the wire length between two consecutive switches, $n$ is the number of repeaters, $m$ represents the size of those repeaters with respect to minimum-size devices, and lastly, $C_G$ and $C_J$ represent the gate and junction capacitance, respectively, of a minimum size inverter. While determining $C_{wire}$, the worst-case scenario is considered, where adjacent wires switch in

opposite directions [30].

The simulation is initially run for 10,000 cycles to allow the 64-IP network to stabilize, and it is subsequently run for 100,000 more cycles. The simulator provides statistics for energy, throughput, and latency.

## 4.4 Experimental Results for Mesh-Based Networks

This thesis first considers the performance of 3D mesh-based NoC architectures. Figure 6a shows the variation of throughput as a function of the injection load. A network cannot accept more traffic than is supplied, and limitations in routing and collisions cause saturation before throughput reaches unity. From Figure 6a, it is clear that both the 3D Mesh and Stacked Mesh topologies exhibit throughput improvements over their two-dimensional counterparts. It is also clear that the ciliated 3D Mesh network shows only a small throughput improvement. However, this is not where a ciliated structure exhibits the best performance. It will be shown later that this network topology has significant benefits both in terms of energy dissipation and silicon area.

These results coincide with the analysis of 3D mesh-based NoC provided in chapter 4, section 1. Equation (4) shows that a 3D mesh will have 29% more interconnection links than a 2D version; hop count calculations have shown that a flit in a 2D mesh network will, on average, traverse 40% more hops than a flit navigating a 3D mesh (according to Table 4); and 3D mesh switches have higher connectivity with the increased number of ports. These all account for throughput improvements. In general, the lower hop count allows a wormhole-routed packet to occupy fewer resources, freeing
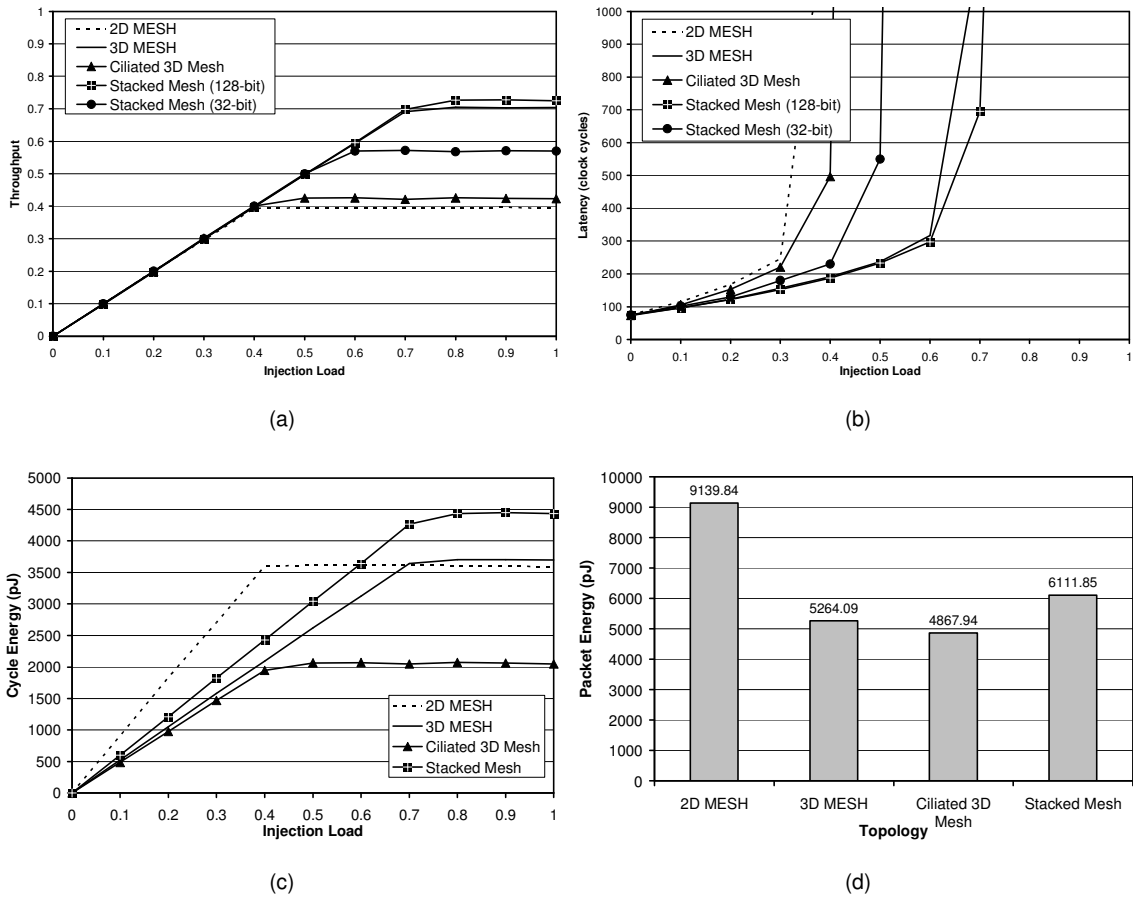
Figure 6. Experimental results for mesh-based NoCs: (a) Throughput vs. injection load, (b) Latency vs. injection load, (c) Cycle energy vs. injection load, and (d) Packet energy
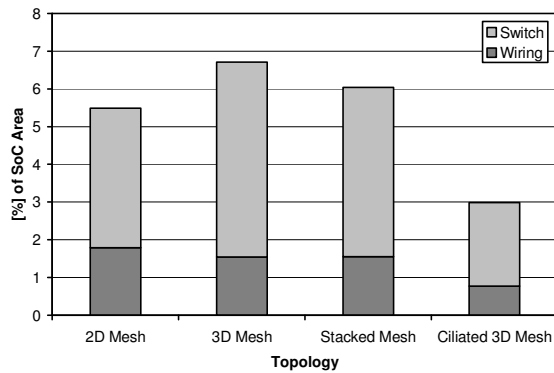


Figure 7. Area Overhead for mesh-based NoCs

Table 4. Average hop count in mesh-based NoCs

| | |
|---|---|
| 2D Mesh | 5.33 |
| 3D Mesh | 3.81 |
| Stacked Mesh | 3.42 |
| Ciliated 3D Mesh | 3.10 |

up links for additional packets. Consequently, there is a corresponding increase in throughput.

Next, the 3D Stacked Mesh architecture is considered. An increase in throughput is evident, as shown in Figure 6a. However, with a 32-bit bus (corresponding to the flit width) connecting the layers of the NoC, throughput improvements are not as substantial as with the 3D Mesh. Contention issues in the bus limit the attainable performance gains. Yet, since communication between layers is bus-based, one may increase the size of the bus without modifying the switch architectures. As a result, the bus width is increased to 128 bits. Any further increase did not have any significant impact on throughput, except to increase the total capacitance on the bus. With this improvement, 3D Stacked Mesh saturates at a slightly higher injection load than a 3D Mesh network. The 3D Stacked Mesh topology also offers a lower hop count in comparison to a strict 3D Mesh. From (6), the average hop count is equal to 3.42. With the lower hop count in addition to the wide, 128-bit bus for vertical transmission, this architecture offers the highest throughput among all the 3D mesh-based networks.

Throughput characteristics of the ciliated 3D Mesh topology differ significantly from the other 3D networks. This network has a saturating throughput that is slightly

higher than a 2D Mesh network and considerably less than both 3D Mesh and Stacked Mesh networks. This is true despite having the lowest hop count at an average of 3.10 hops. However, with only 64 interswitch links, compared to 144 in the 3D Mesh and 112 in the 2D Mesh, throughput improvements due to hop count are negated by the reduced number of links. The fact that there are multiple functional IP blocks for every switch is also responsible for considerable lower throughput due to contention issues in the switches.

Figure 6b depicts the latencies for the architectures under consideration. Here, it is seen that 3D mesh-based NoCs have superior latency characteristics over the 2D versions. This is a product of the reduced hop count characteristic of 3D mesh-based topologies.

Energy dissipation characteristics for three-dimensional mesh-based NoCs reveal a substantial improvement over planar NoCs. The energy dissipation profiles of the mesh-based NoC architectures under consideration are shown in Figure 6c. Energy dissipation is largely dependent on two factors: architecture and injection load. These two parameters are considered the independent factors in this analysis. As shown in (7), the energy dissipation in a NoC depends on the energy dissipated by the switch blocks and the inter-switch wire segments. Both these factors depend on the architectures. The design of the switch varies with the architecture and inter-switch wire length is also architecture dependent [24]. Besides the network architecture, injection load has a clear effect on the total energy dissipation of a NoC, in accordance with Figure 6c. Intuitively, it is clear that with more packets traversing the network, power will increase. This is why packet

energy, in Figure 6d, is an important attribute for characterizing NoC structures. Notice that, at saturation, a 2D Mesh network dissipates less power than both 3D Stacked Mesh and 3D Mesh networks. This is the result of the lower 2D Mesh throughput, and the 3D networks consume more energy because they transmit more flits at saturation. Packet energy is a more accurate representation of the cost of data transmission. With packet energy in mind, it can be seen that every 3D topology provides a very substantial improvement over a 2D Mesh. Also, the energy dissipation of the ciliated mesh topology is less, still, than that of 3D Mesh network. These results follow closely the hop count calculations summarized in Table 4, with the exception of the packet energy for a 3D Stacked Mesh network. Energy is heavily dependant on interconnect energy, and this is where the 3D Stacked Mesh suffers. Since vertical communication takes place through wide busses, the capacitive loading on those busses results in a significant amount of energy. As a result, though 3D Stacked Mesh has a lower hop count compared to 3D Mesh, it dissipates more packet energy on average. Regardless, the profound energy savings possible in these 3D architectures provides serious motivation for a SoC designer to consider a three dimensional integrated circuit.

The final performance metric considered in this study is the overall area overhead incurred with the instantiation of the various networks. Figure 7 shows the area penalty from each NoC design, both in terms of switch area and interconnects area. It shows that while the 3D Mesh and 3D Stacked Mesh NoCs reduce the amount of wiring area, switch overhead is increased. For both 3D Mesh and 3D Stacked Mesh NoCs, the number of longer inter-switch links in *x-y* plane is reduced. There are 96 *x-y* links for both

32

topologies, for 3D Stacked Mesh, 16 buses are present, and for the 3D Mesh, 48 vertical links are present. In comparison, the conventional 2D mesh-based NoC has 112 links in the horizontal plane. As the 3D NoCs have fewer long horizontal links they incur less wiring area overhead. Although there are a large number of vertical links, the amount of area incurred by them is very small due to the $2\mu m \times 2\mu m$ interlayer vias. However, an increased number of ports per switch results in larger switch overhead for both of these NoC architectures, ultimately causing the 3D Mesh and 3D Stacked Mesh topologies to incur more silicon area in spite of wiring improvements. On the other hand, 3D Ciliated Mesh shows a significant improvement in terms of area. The $4 \times 4 \times 2$ 3D Ciliated Mesh structure involves half the number of switches as the other mesh-based architectures in addition to only 64 links. As a result, the area overhead is accordingly smaller.

### 4.5 Experimental Results for Tree-Based Networks

In this section, the performance of the three-dimensional tree-based NoCs is evaluated. It has already been established that 2D and 3D versions of the tree topologies should have identical throughput and latency characteristics, and Figure 8a and Figure 8b support this. Consistent with the analysis of mesh-based NoCs, Figure 8a shows the variation of throughput as a function of injection load, and Figure 8b shows the effect of injection load on latency. The assumption here was that the switches and the inter-switch wire segments are driven by the same clock as explained earlier. Consequently under this assumption, in terms of throughput and latency there is no advantage to choosing a 3D IC over a traditional planar IC for a tree-based NoC. However, this is eclipsed by the
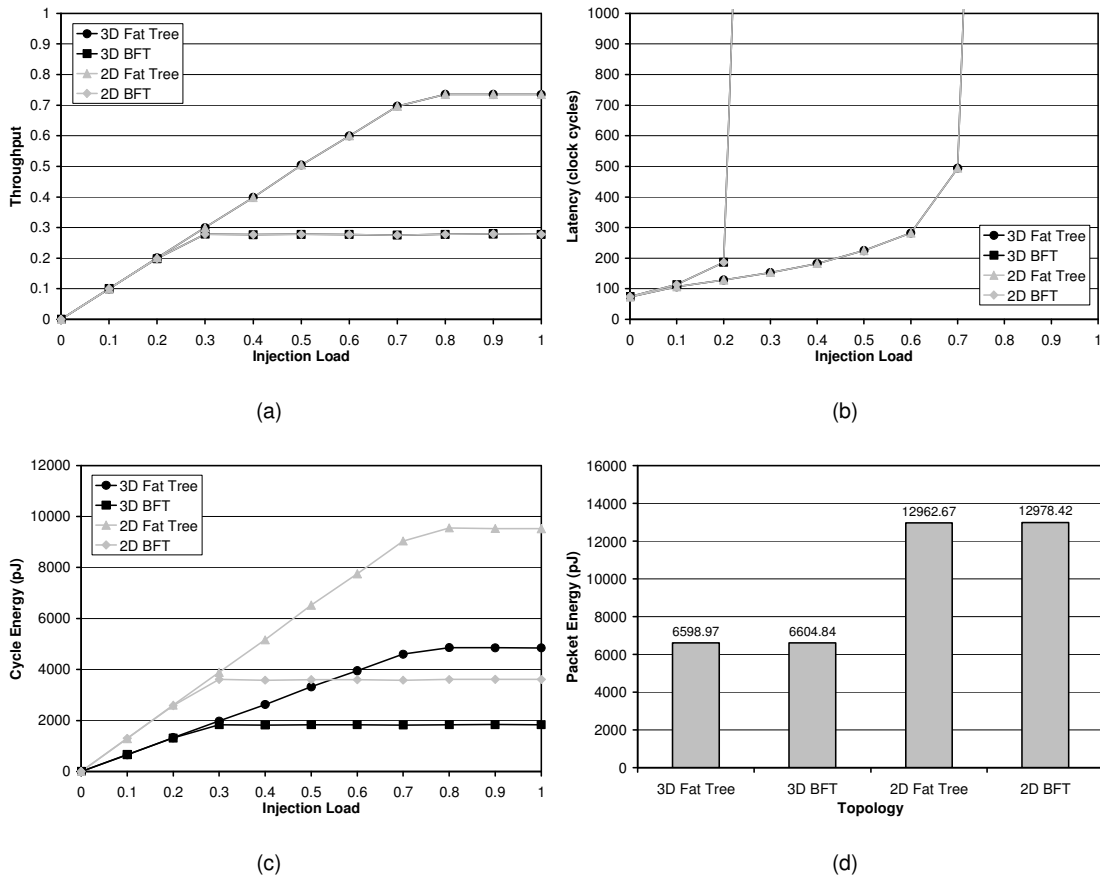
33

Figure 8. Experimental results for tree-based NoCs: (a) Throughput vs. injection load, (b) Latency vs. injection load, (c) Cycle energy vs. injection load, and (d) Packet energy
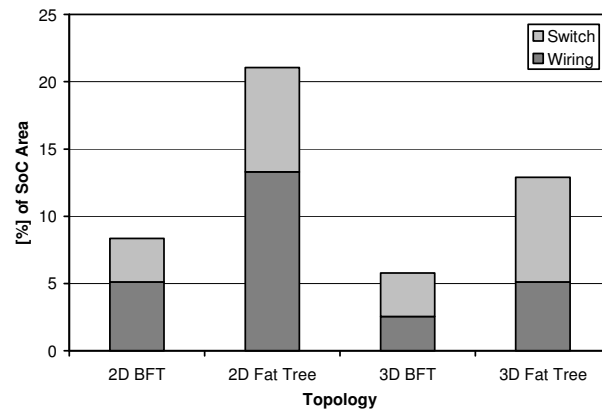


Figure 9. Area Overhead for tree-based NoCs

superior performance achieved in terms of energy and area overhead. If the NoC switches can be designed to operate as fast as the inter-switch wires, then it will show later that the 3D tree-based architectures will benefit significantly in terms of latency and bandwidth.

The energy profiles for 3D tree-based NoCs (Figure 8c) reveal significant improvements over 2D implementations. Both BFT and Fat Tree (SPIN) networks show a very large reduction in energy when 3D ICs are used. Once again, energy dissipation is largely dependant both on architecture and injection load. Each NoC shows that energy dissipation increases with injection load until the network becomes saturated, similar to the throughput curve shown in Figure 8a. The energy profiles show that the Fat Tree networks cause higher energy dissipation than the Butterfly Fat Tree instantiations, but this is universally true only at high injection load.

Again, this is the motivation to consider packet energy of the networks as a relevant metric for comparison, shown in Figure 8d. Energy savings in excess of 45% are achievable by adopting 3D ICs as a manufacturing methodology, and both BFT and Fat Tree networks show similar improvements. In case of tree-based NoCs, where the basic network topology remains unchanged in 3D implementations, all improvements in energy dissipation are caused by the shorter wires. As showed earlier in Table 1, a three-dimensional structure greatly reduces the inter-switch wire length. The overall energy dissipation in a NoC is heavily dependant on the interconnect energy, and this reduction in inter-switch wire length effects very large savings.

Besides advantages in terms of energy, three-dimensional ICs enable tree-based NoCs to reduce silicon area overhead by a sizable margin. Figure 9 shows the overall

area overhead of tree-based NoCs. Although no improvements are made in terms of switch area, the reductions in inter-switch wire lengths and amount of repeaters are responsible for substantial reductions in wiring overhead. This is especially true of the Fat Tree network, which has more interconnects in the higher levels of the tree; wiring overhead is reduced more than 60% by instantiating the network into a 3D IC.

## 4.6 Effects of Traffic Localization

Until this point, a uniform spatial distribution of traffic has been assumed. In a SoC environment, different functions would map to different parts of the chip and the traffic patterns would be expected to be localized to different degrees [31]. This thesis therefore considers the effect of traffic localization on the performance of the 3D NoCs, and in particular it considers the illustrative case of spatial localization where local messages travel from a source to the set of the nearest destinations. In the case of BFT and Fat Tree, localized traffic is constrained to within a cluster consisting of a single sub-tree while, in the case of 3D Mesh, it is constrained to within the destinations placed at the shortest Manhattan distance [24].

On the other hand, the 3D Stacked Mesh architecture is created simply to take advantage of the inexpensive vertical communication. The research pursued by Li et al. in [9] suggested that in a 3D multi-processor SoC, much of the communication should take place vertically, taking advantage of the short inter-layer wire segments. This is a result of a large proportion of network traffic occurring between processor and the closest cache memories, which are often placed along the $z$-dimension. Consequently, in these

36

situations, the traffic will be highly localized, and this study therefore considers localized traffic to be constrained to within a pillar for 3D Stacked Mesh. Figure 10 summarizes these effects, revealing the benefits of traffic localization. More packets can be injected into the network, improving the throughput characteristics of each topology as shown in Figure 10a and Figure 10c, which also shows the throughput profile of the 2D topologies for reference. Analytically, increasing localization reduces the average number of hops that a flit must travel from source to destination.

Figure 10a reveals that the 3D Stacked Mesh network provides best performance in terms of throughput in the presence of localized traffic. However, this is achieved by using a wide bus for vertical communication. Let us consider what occurs when the bus size is equal to the flit width of 32 bits. With low localization, the achieved throughput is higher than that in a 2D Mesh network. However, when the fraction of localized traffic in the vertical pillars is increased, a huge performance degradation is seen. This is due to the contention in the bus. When the bus width is increased to 128 bits, throughput increases significantly with increase in localized traffic. This happens due to less contention in a wider communication channel.

Figure 10b and Figure 10d depict the effects of localization on packet energy, and, unsurprisingly, there is a highly linear relationship between these two parameters. Packet energy is highly correlated with the number of hops from source to destination, and the resultant reduction of packet energy with localization supports this correlation. For the mesh-based networks, 3D Ciliated Mesh exhibits the lowest packet energy due to its low hop count and very short vertical wires. In fact, at highest localization, the packet energy
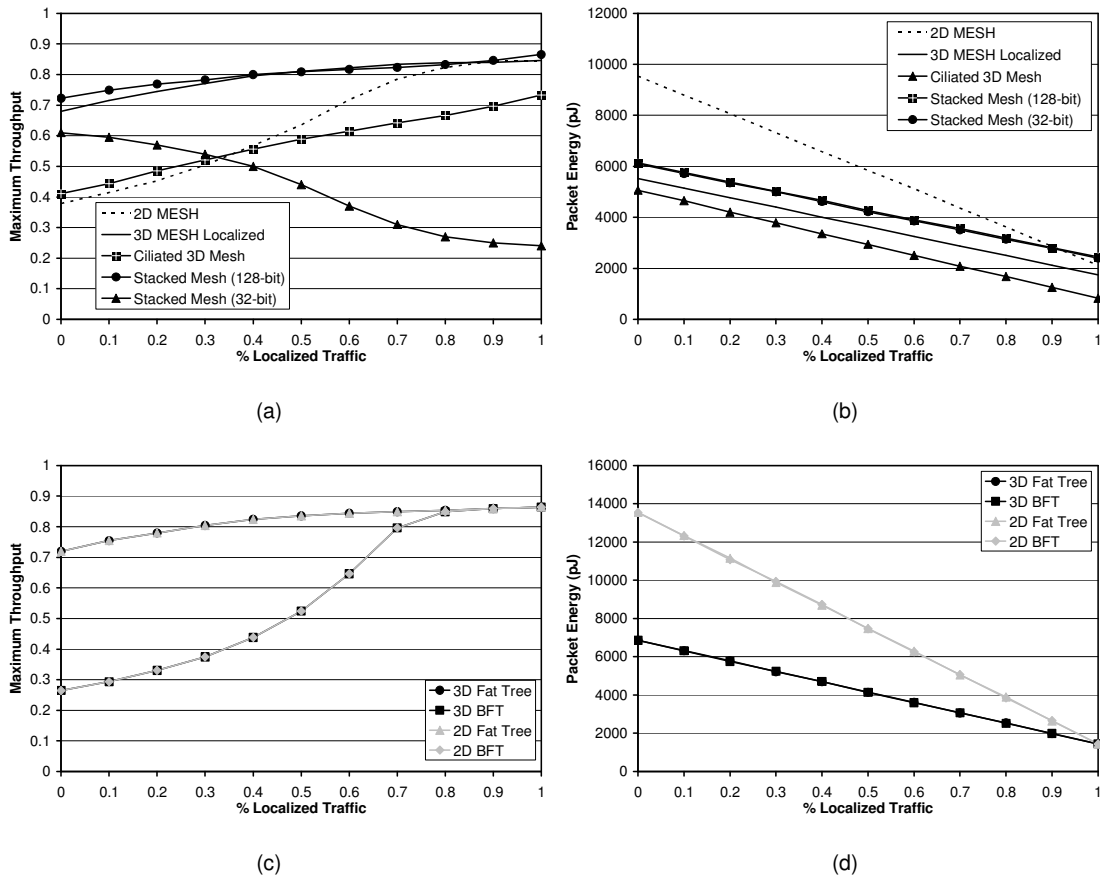
37

Figure 10. Localization effects on mesh-based NoCs in terms of: (a) throughput and (b) packet energy; and on tree-based NoCs in terms of (c) throughput and (d) packet energy

for a 3D Ciliated Mesh topology is less than 50% of that of the next-best-performing topology: 3D Mesh. For the tree-based NoCs, both 3D networks have much-improved packet energy with traffic localization.

As can be seen from Figure 10, there are tradeoffs between packet energy and throughput. For instance, the best-performing topology in terms of energy, ciliated Mesh, operates at the lowest throughput even when traffic is highly localized. On the other hand, although a 3D Stacked Mesh network with wider bus width achieves superior

38

throughput without necessitating a highly local traffic distribution, it incurs more energy

dissipation than other structures under local traffic due to the capacitive loading on the

interlayer busses. However, the other topologies lie in some middle ground between these

two extremes, and in general, it is clear that 3D ICs continue to effect improvements on

NoCs under localized traffic.

**4.7 Effects of Wire Delay on Latency and Bandwidth**

In NoC architectures, the inter-switch wire segments, along with the switch

blocks, constitute a pipelined communication medium as shown in Figure 11. The overall

latency (in nanoseconds) will be governed by the slowest pipelined stage. Table 2

showed earlier that the maximum wire delays for the network architectures are different.

Though the vertical wire delays are very small, still the overall latency will be depended

on the delay of the switch blocks. Though the delays of the switch blocks were

constrained within the 15FO4 limit, they were still the limiting stages in the pipeline,

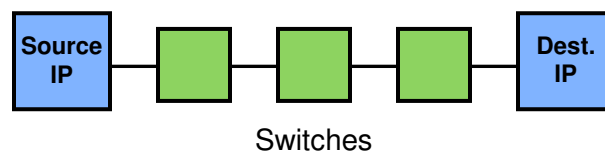specifically when compared to the fast vertical links. Yet, considering a hypothetical case,

Figure 11. The pipelined nature of NoCs

Table 5. Bandwidth of Network Architectures at Simulated and Hypothetical Frequencies (Terabits/s)

| | $f$ = 1.67 GHz | $f$ = 1/(max wire delay) | % increase |
|---|---|---|---|
| 2D Mesh | 1.357 | 3.711 | 173.5 |
| 3D Mesh | 2.412 | 6.596 | 173.5 |
| Ciliated 3D Mesh | 1.457 | 3.983 | 173.5 |
| 3D Stacked Mesh | 2.488 | 6.804 | 173.5 |
| 2D BFT | 0.9543 | 1.039 | 8.9 |
| 2D Fat Tree | 2.515 | 2.738 | 8.9 |
| 3D BFT | 0.9543 | 2.474 | 159.2 |
| 3D Fat Tree | 2.515 | 6.520 | 159.2 |

which ignores the implications of switch design, where the clock period of the network is equal to the inter-switch wire delay, then the clock frequency can be increased, and, resultantly, the latency can be reduced significantly. With this in mind, latency in nanoseconds (instead of latency in clock cycles) and bandwidth (instead of throughput) are calculated. All other network parameters are kept consistent with the previous analysis.

A plot of latency for all network topologies is shown in Figure 12, and Table 5 depicts the network bandwidth in units of Terabits per second. To calculate bandwidth, the following expression is followed:

$$BW = TP_{max} \cdot \frac{1}{f} \cdot w_{flit} \cdot N,$$

(13)

where $TP_{max}$ represents the throughput at saturation, $f$ represents the clock frequency, $w_{flit}$ is the flit width, and $N$ is the number of IP blocks. Table 5 shosw the performance difference achieved by running the NoC with a clock as fast as the inter-switch wire,
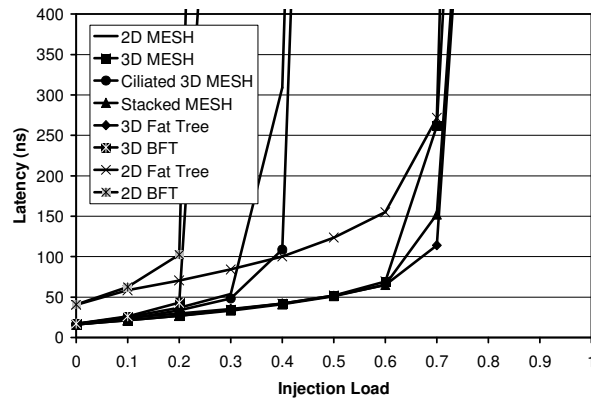
Figure 12. Latency in ns at hypothetical clock frequencies

disregarding the switch design constraints. It is evident that the tree-based architectures show the greatest performance improvement in this scenario going from 2D to 3D implementations, as the horizontal wire lengths are also reduced.

## 4.8 Network Aspect Ratio

The ability to stack layers of silicon is not without nuances. Upcoming 3D processes have a finite number of layers due to manufacturing difficulties and yield issues [3]. Furthermore, it is speculated [3] that the number of layers in a chip stack are not likely to scale with transistor geometries. This has a nontrivial effect on the performance of 3D NoCs. Consequently, future NoCs may have a greater number of intellectual property blocks in the horizontal dimensions than vertically. The effect of this changing aspect ratio must be characterized.

For a more in-depth illustration of these effects, the overall performance of a mesh-based NoC in a 2-layer IC will be evaluated in comparison to the previously-analyzed 3D 4×4×4 Mesh and 2D 8×8 Mesh. Here, a 64-IP 8×4×2 Mesh is considered to
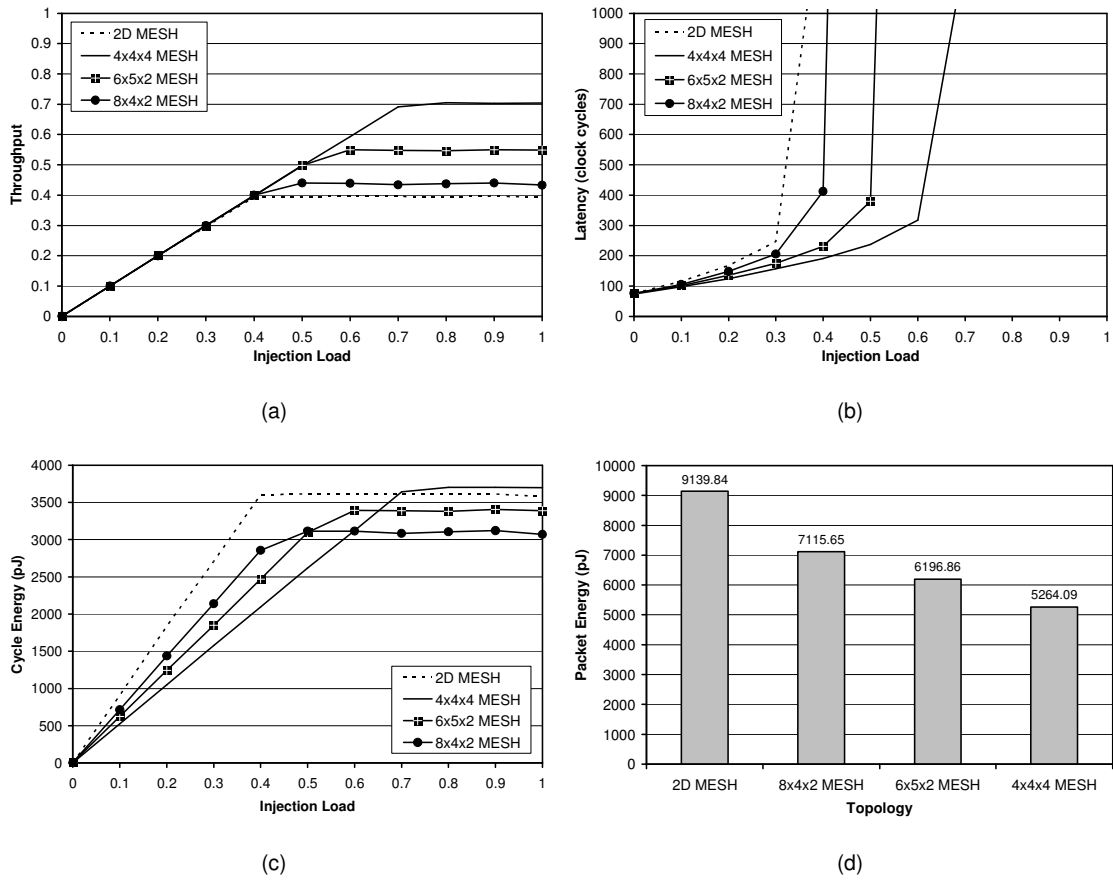
41

(a)

(b)

(c)

(d)

Figure 13. Comparing two 2-layer NoCs: (a) Throughput vs. injection load, (b) Latency vs. injection load, (c) Cycle energy vs. injection load, and (d) Packet energy
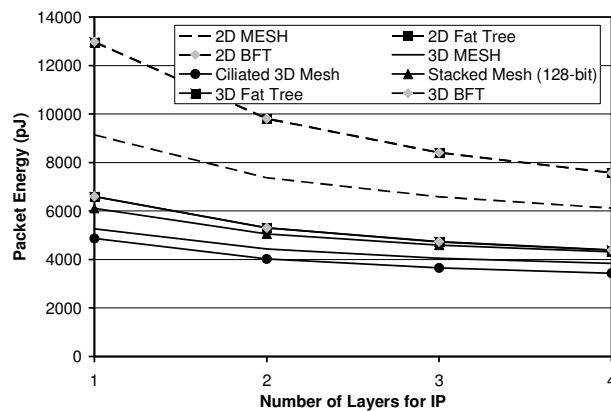


Figure 14. The effect of multi-layer IPs

match the 64-IP network size, in order to make the comparison of latency and energy as fair as possible, along with a 60-IP 6×5×2 Mesh to show a network which is similar in size and that results in a more square overall footprint than the 8×4×2 Mesh. Figure 13 summarizes the analysis of these 2-layer ICs. Throughput characteristics are seen in Figure 13a. It shows clearly that the 6×5×2 Mesh achieves a significantly higher throughput than the 2D 8×8 Mesh and the 8×4×2 Mesh, which suffers from a high average hop count (4.44 vs 4.11 for the 6×5×2 Mesh), while achieving a lower maximum throughput than the 4-layer mesh. Likewise, the 2-layer mesh NoCs outperform the 2D Mesh in terms of latency, shown in Figure 13b, without exceeding the performance of the 4-layer 3D instantiation. This trend continues when considering cycle energy (Figure 13c) and packet energy (Figure 13d). These results are as expected. With the first layer added, significant improvements are apparent in terms of each performance metric over the 2D case. Though the multi-layer NoC exhibits superior performance characteristics compared to a 2D implementation, it will have to circumvent significant manufacturing challenges. Yet, even if implementations are limited to two-layer 3D realizations, they will still significantly outperform the planar NoCs.

## 4.9 Multi-Layer IPs

Throughout this thesis, each IP block has been assumed to be instantiated in one layer of silicon. However, as discussed in [10], it is certainly possible for the IP blocks to be designed using multiple layers. So, each network architecture is analyzed with mult-layer IPs. The pipelined communication shown in Figure 11 is assumed; i.e. the NoCs are

43

constrained by the switch delay and it cannot be driven as fast as the inter-switch wire. Considering this, multi-layer IPs have no effect on either throughput or latency (assuming the same clock frequency for all networks), but there are nontrivial effects on the energy dissipation profile. This effect on packet energy is depicted in Figure 14. The energy savings come from reduced horizontal wire lengths. For instance, if a 2.5mm×2.5mm IP block is instantiated in 2 layers, the IP's circuitry is spread over 2 layers, and the footprint reduces by a factor of 1.414. Similarly, if instantiated in 3 layers, the footprint reduces by a factor of 1.732, and with 4 layers, the factor is 2. Although the vertical wire lengths are increased 2, 3, and 4 times, respectively, in order to span the entire multi-layer IP, the negative effects on energy incurred by this are eclipsed by the significant reductions in horizontal wire lengths. However, multi-layer IPs increase the number of layers in a 3D IC, placing an increased burden on manufacturability.

## 4.10 Conclusions

This chapter covers a very significant thrust of this thesis. Through the application of real traffic patterns, the mesh-based and tree-based architectures of chapter 3 were evaluated in terms of throughput, latency, energy, and area overhead. In terms of throughput and latency, the new mesh-based topologies effect significant improvements over the 2D Mesh, while due to an exact topology match, no improvements are made for the tree-based architectures. However, when energy is considered, all 3D architectures show very impressive gains compared to the 2D architectures. Additionally, area is reduced for all tree-based architectures and some mesh-based architectures.

Also, when traffic localization is considered, 3D NoCs continue to show improvements in terms of throughput and energy. At high localizations, the 3D Ciliated Mesh architecture shows its greatest advantage, low packet energy; this may be reason to use this architecture despite minimal throughput improvements over the traditional 2D Mesh. Next, the implications of wire delay are investigated. If a network switch can be designed to match the maximum wire delay, the architectures with the shortest interswitch wire lengths, such as all mesh-based architectures and the 3D tree-based architectures, can exhibit dramatic improvements in terms of bandwidth and latency. Finally, two-level 3D NoCs are evaluated, and the implications of multi-layer IPs are established. Two-level NoCs show a less-significant improvement over 2D versions as expected, and if many-layer 3D ICs become more manufacturable in time, multi-layer IPs will allow further reductions in energy to be realized. Overall, the advantages of 3D NoCs introduced in this chapter provide very compelling reasons to adopt this methodology as 3D ICs become available in the ensuing months and years.

# CHAPTER FIVE

# HEAT DISSIPATION PROFILE OF 3D NOCS

Heat dissipation is an extremely important concern in 3D ICs. Already, thermal effects have been known to have significant implications on device reliability and interconnects in traditional 2D circuits [32]. With the reduced footprint inherent to 3D ICs, this problem is exacerbated as the energy dissipated throughout the entire chip is now constrained to a smaller area, therefore increasing the energy density of these circuits. As a result, it is imperative that thermal issues are addressed in any system involving 3D integration.

Accordingly, an analysis of three-dimensional networks-on-chip is incomplete without an examination of temperature. It is especially important since the interconnect structure of a NoC can consume close to 50% of the overall power budget [13]. As temperature is closely related to the energy dissipation of the IC, this analysis will draw heavily upon the discussion of energy from chapter 4, sections 1 and 2. This chapter considers the 2D and 3D NoC architectures introduced in chapter 3 and evaluates them in the presence of real traffic patterns. Furthermore, chapter 4 has shown that the energy dissipated by the interconnection infrastructure, i.e. the communication energy, can be reduced compared to a 2D implementation by virtue of the inherent nature of the network architecture. Consequently, it will have a positive effect on heat dissipation.

## 5.2 Temperature Analysis

Temperature in a 3D IC is related to a variety of factors including power dissipation and power density. In an integrated circuit, according to [33], the steady state temperature distribution is given by the following Poisson equation:

$$\nabla^2 T(\mathbf{r}) = \frac{-g(\mathbf{r})}{k_l(\mathbf{r})} \quad . \tag{14}$$

Here, $\mathbf{r}$ is the three-dimensional coordinate $(x,y,z)$. $T(\mathbf{r})$ is the temperature inside the chip at point $\mathbf{r}$, $g(\mathbf{r})$ is the volume power density at that point, and $k_l(\mathbf{r})$ is the thermal conductivity. An important fact to note is that $k_l(\mathbf{r})$, the thermal conductivity, is constrained by the manufacturing process and a designer has little or no control over it. Therefore, the volume power density, $g(\mathbf{r})$, is the parameter over which a designer has the most control. The challenge facing all designers of 3D ICs is to exercise control over this parameter.

In a 3D integrated circuit, the volume power density of the chip is increased. The lateral dimensions are significantly smaller, and as a result, the total power of the circuit is dissipated in a much smaller area. For instance, in a four-layer 3D IC, the floor area is reduced by a factor of 4, for an eight-layer 3D IC, that area is reduced by a factor of 8, etc. Clearly, the energy of the entire chip is now constrained to a much smaller footprint, and the volume power density increases with respect to this.

## 5.4 The Relationship between Temperature and Energy

According to (14), it is clear that lower energy corresponds to lower heat. With an increase in volume power density, there is a corresponding increase in temperature.
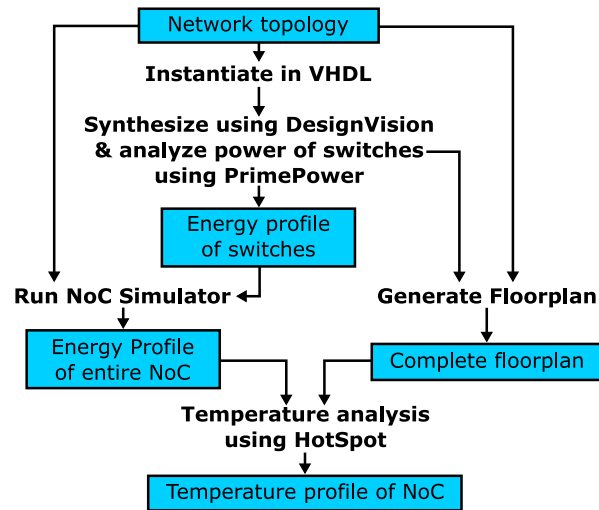
Figure 15. Design Flow

However, with 3D integration, the density of the chip is increased. As a result, there is a factor leading to higher heat in 3D NoCs. On the other hand, in 3D NoCs the communication energy can be reduced compared to a 2D implementation due to the various factors explained above. Consequently, it will lead to lesser heat dissipation in 3D NoCs. To quantify the overall effects of these two opposing factors, the heat dissipation profile for the aforementioned 3D NoC architectures is evaluated in presence of realistic traffic patterns.

## 5.5 Simulation Methodology

The temperature profiles of the 3D NoCs were obtained through simulations following the methodology shown in Figure 15. First, the network architecture is chosen. Subsequently, the network switches are instantiated in VHDL and synthesized using Synopsys DesignVision and a 90-nm standard cell library from CMP [29]. Here,

Synopsys PrimePower is run to generate the energy profiles of the network switches. Next, the overall floorplan of the NoC is created. In order to generate the energy profile of the entire NoC it is necessary to incorporate the energy dissipated by each inter-switch stage. This is calculated taking into account the specific layout of each topology following the method elaborated in [23].

Following this, the NoC simulator is run. It is flit-driven and utilizes wormhole routing; it is cycle-accurate and can model various network structures. A cycle-accurate simulation is run for 10,000 cycles in order for the network to stabilize, and it is subsequently run for 100,000 cycles to gather data. The simulator generates the overall energy profile of the NoC as well as traffic statistics, such as throughput and latency. Network topology, switch energy profiles, and injection load (which is the average number of flits injected into the network each cycle per processing element) are inputs to the simulator.
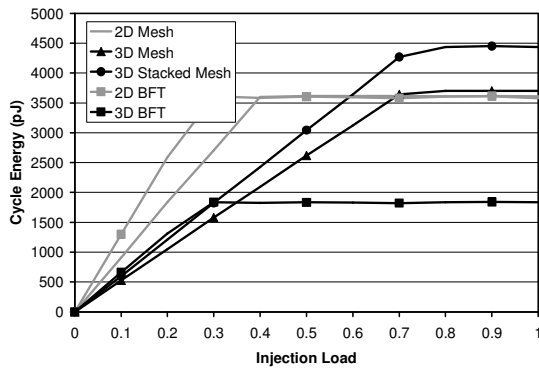
Finally, with a complete floorplan and power profile, the Hotspot tool, developed by a research team at the University of Virginia [34], is used to generate the temperature profile. Hotspot takes the floorplan and applies the power profile to it, and with this information, it calculates the volume power density. From this, the temperature profile is generated. This design flow affords sufficient accuracy with the use of Synopsys PrimePower and Hotspot, but it also affords great flexibility in NoC design with the inclusion of the network simulator, as it is flexible in supporting many different structures. The simulator is capable of handling different types of traffic patterns, as well. For this analysis, a self-similar injection process [26] is followed, and traffic is

considered to be distributed randomly and evenly across the network, although future work will examine the effects of highly-localized traffic and other injection processes on temperature.
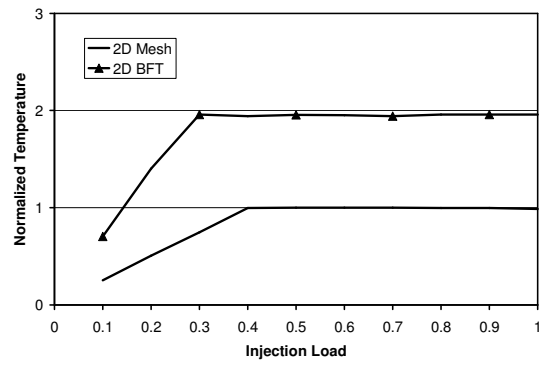
## 5.6 Experimental Results

In accordance with the prescribed methods, 64-IP instantiations of each 3D NoC architecture were analyzed for thermal performance, with temperature taken as a function of injection load. As explained in section 4, temperature is closely related to power density, so, likewise, these temperature profiles are very similar in form to the energy profiles, shown in Figure 16a. The analysis begins with the 2D topologies. A plot of the temperature characteristics of the two architectures is shown in Figure 16b with the temperature normalized to the maximum temperature of the 2D Mesh, considered as the baseline case. Figure 16a shows temperature saturating at different values and at different injection loads for each topology, like the communication energy dissipation profiles.
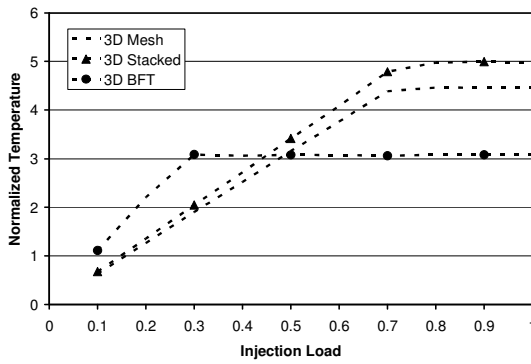
With a 3D network implementation, this thesis has shown significant improvements in terms of energy dissipation, particularly packet energy, which is shown again in Figure 16e. Packet energy is a more accurate representation of the cost of data transmission than cycle energy. Intuitively, it is clear that with more packets traversing the network, power will increase. This is why the packet energy in Figure 16e is an important attribute for characterizing NoC structures. The 3D Stacked Mesh and 3D Mesh structures show 33% and 42% improvements over 2D Mesh, and 3D BFT shows 49% improvement over 2D BFT. These improvements in energy have substantial effects
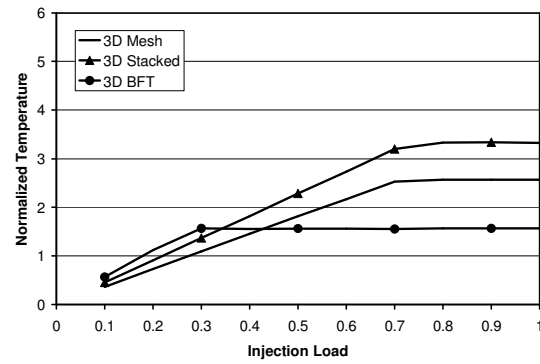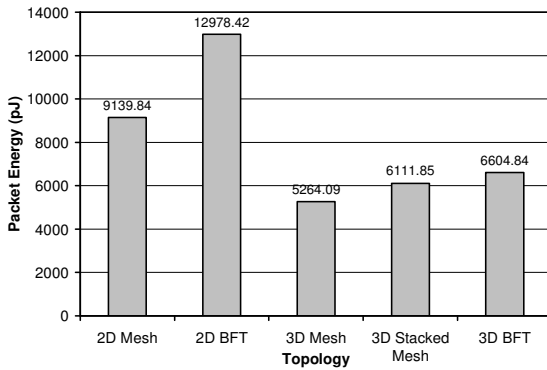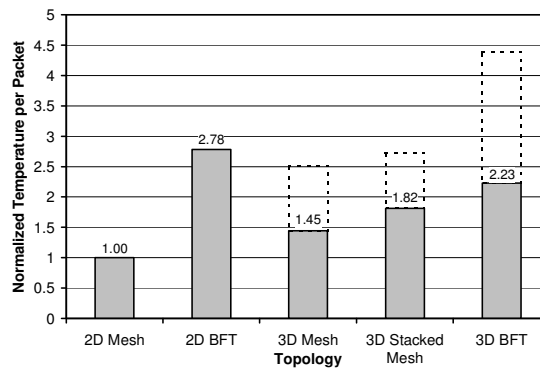
Figure 16. Experimental results for 3D NoCs: (a) Cycle energy, (b) Maximum temperature in 2D architectures, (c) Hypothetical temperature for 3D architectures, (d) Maximum temperature for 3D architectures, (e) Packet energy, and (f) the normalized contribution to temperature per packet.

51

on the temperature characteristics of these 3D networks. Let us first consider the hypothetical case where 3D implementations of these topologies dissipate the same communication energy per packet as the 2D versions. This case is shown by the dotted lines in Figure 16c. It is very clear that in the absence of any packet energy gains, the result is a much hotter network. This is due to the inherent nature of 3D ICs. As discussed in chapter 2, when moving to a 3D NoC, the overall chip area remains constant while the footprint is reduced. In these 10mm×10mm 4-layer, 3D implementations, the entire energy dissipation of the chip is constrained to an area one quarter the size of the 20mm×20mm 2D implementations. As a result, the power density should be significantly increased.

However, the actual temperature profiles of the 3D networks, depicted by the solid lines in Figure 16d, show a marked difference. This highlights a very important characteristic of NoC in a 3D environment: the savings in communication energy incurred by choosing a 3D NoC implementation partially mitigate what would otherwise be a drastic increase in temperature.

To help describe this effect, Figure 16f presents the normalized temperature contribution per packet, using the 2D mesh architecture again as the baseline case. The dotted bars represent the hypothetical case discussed earlier. The contribution to temperature per packet metric follows a similar idea to that of packet energy. Each packet sent through the network is responsible for a certain amount of energy dissipation. This, in turn, causes a rise in temperature. Therefore, as packet energy quantifies the energy efficiency of a NoC, the temperature contribution per packet thus quantifies the

temperature efficiency of a NoC. All topologies show real improvements over the hypothetical case, and, in fact, the 3D version of the BFT network have lower temperature than its 2D counterpart. This can be attributed, in part, to the very high (49%) decrease in packet energy that is characteristic of a 3D BFT implementation over a 2D BFT instantiation.

## 5.7 Conclusions

The Network on Chip (NoC) paradigm has emerged as an effective methodology for designing big multi-core SoCs. On the other hand, 3D integration is emerging as a new interconnect paradigm to overcome the performance limitation of the conventional two-dimensional IC. 3D NoCs incorporate the advantages of these two new paradigms. When instantiating in a 3D environment, NoCs can reduce the temperature of the whole chip, which is one of the principal limiting factors of any 3D process. This happens due to the reduction of communication energy in 3D network structures. Depending on the specific network structure, some 3D NoCs are capable of even reducing the heat dissipation compared to a 2D implementation. Thus, 3D NoCs are efficient in addressing the heat dissipation concerns of SoCs implemented via 3D integration processes.

**CHAPTER SIX**

**CONCLUSION**


This thesis has demonstrated that besides reducing the footprint in a fabricated design, three-dimensional network structures provide a better performance compared to traditional, 2D NoC architectures. It has demonstrated that both mesh- and tree-based NoCs are capable of achieving better performance when instantiated in a 3D IC environment compared to more traditional 2D implementations. The mesh-based architectures show significant performance gains in terms of throughput, latency and energy dissipation with a small area overhead. On the other hand, the 3D tree-based NoCs achieve significant gain in energy dissipation and area overhead without any change in throughput and latency. However, if the NoC switches are designed to be as fast as the interconnect, even the 3D tree-based NoCs will exhibit performance benefits in terms of latency and bandwidth. Furthermore, 3D NoCs are efficient in addressing the temperature issues characteristic of 3D integrated circuits.

The Network-on-Chip (NoC) paradigm continues to attract significant research attention in both academia and industry. With the advent of 3D ICs, the achievable performance benefits from NoC methodology will be more pronounced as shown in this paper. Consequently this will facilitate adoption of the NoC model as a mainstream design solution for larger multi-core system chips.


**6.1 Future Directions**

*Wireless NoC*

In a relatively radical departure from traditional copper wires, wireless signals in the RF and microwave spectra may be used to transmit data on-chip, and transmission can be either through waveguides or through packaging and IC structures. Research groups such as Chang, et al. [35], propose using wave-guiding structures for on-chip communication instead of truly wireless transmission due to the prohibitive space concerns with antenna aperture. However, Pande, et al. [36], believe that antennas based on carbon nanotubes can be utilized to create a NoC interconnection network using optical (as opposed to RF or microwave) frequencies.

This wireless approach may be of particular use to 3D NoCs. Principal limiting factors to the cost-effective fabrication of 3D integrated circuits are yield limitations [3]. In particular, this is related to the processes of vertical via insertion and bonding. If the vertical vias can be eliminated by the integration of wireless vertical communication, yield improvements may be possible. On the other hand, the fabrication of these nanotubes may be yield-inhibiting in their own right. These issues remain to be seen, and the implications of short vertical wireless transmission should be investigated. Furthermore, the possibility of instantiating the entirety of the 3D NoC using wireless communication should be examined as well. Lastly, the throughput, latency, energy, area, and thermal tradeoffs of wireless communication must be established.

*Further Temperature Analysis*

A principal limitation of the temperature analysis in chapter 5 is that only the

energy contribution of the interconnection network is considered and the energy dissipation of the IP blocks is ignored. This was done to free the analysis from the constraints of any specific application. Furthermore, as comparisons were made only between the different network structures, the conclusions will hold for most if not all applications. That said, however, a thorough analysis of temperature is not truly complete until the temperature contributions of the entire IC are exhausted. For proper thoroughness, future analysis should include multiple IC designs, each for a different application, and each evaluated using each of the NoC architectures introduced in chapter 3. The applications may include a multi-core processor, a parallel processing ASIC such as an LDPC decoder or FFT computer, a video processing MPSoC, and a low-power ASIC.

*LDPC Decoder*

Networks-on-chip are well-suited for parallel processing algorithms due to the regular structure both of the networks and the algorithms. One algorithm that has attracted much research attention recently is low-density parity check (LDPC) decoding. LDPC codes [37] allow very high levels of error correction; these codes exhibit error performance that is closer to the Shannon Limit of any error correction code to date [38]. However, these codes require very intensive processing, and this can map well to a NoC, as shown by Theocharides, et al., in [39]. This application should prove to be an excellent test bench for demonstrating the network structures introduced in this thesis and their respective tradeoffs and advantages.

# BIBLIOGRAPHY

[1]     P. Magarshack and P.G. Paulin, "System-on-Chip beyond the Nanometer Wall," *Proceedings of 40th Design Automation Conf. (DAC 03)*, ACM Press, 2003, pp. 419-424.

[2]     International Technology Roadmap for Semiconductors 2005: Interconnect, [online] http://www.itrs.net/

[3]     A. W. Topol et al., "Three-dimensional integrated circuits," *IBM Journal of Research & Development*. Vol. 50 No. 4/5 July/September 2006.

[4]     W. R. Davis et al. "Demystifying 3D ICs: The pros and cons of going vertical." *IEEE Design and Test of Computers*, 22(6), Nov. 2005.

[5]     Y. Deng et al. "2.5D System Integration: A Design Driven System Implementation Schema" *Proc. of the Asia South Pacific Design Automation Conference*, 2004.

[6]     M. Ieong et al. "Three Dimensional CMOS Devices and Integrated Circuits." *Proc. of IEEE Custom Integrated Circuits Conference*, 2003.

[7]     L. Benini and G. De Micheli, "Networks on Chips: A New SoC Paradigm," *IEEE Computer*, Jan. 2002, pp. 70-78.

[8]     W. J. Dally, and B. Towles, "Route Packets, Not Wires: On-Chip Interconnection Networks", *Proceedings of the 2001 DAC*, June 18-22, 2001, pp. 683-689.

[9]     F. Li et al., "Design and Management of 3D Chip Multiprocessors Using Network-in-Memory", *Proceedings of the 33rd International Symposium on Computer Architecture (ISCA'06)*, pp. 130-141

[10]    V. F. Pavlidis and E. G. Friedman, "3-D Topologies for Networks-on-Chip", *IEEE Transactions on Very Large Scale Integration (VLSI)*, pp. 1081-1090, October 2007.

[11]    J. Srinivasan et al., "Exploiting Structural Duplication for Lifetime Reliability Enhancement," *Proc. Int'l Symp. Computer Architecture* (ISCA 05), IEEE CS Press, 2005, pp. 520-531.

[12]     J. Tsai, C. C. Chen, G. Chen, B. Goplen, H. Qian, Y. Zhan, S. Kang, M. D. F. Wong, and S. S. Sapatnekar, "Temperature-Aware Placement for SoCs" *Proceedings of the IEEE*, vol. 94, no. 8, pp. 1502-1518, Aug. 2006.

[13]     T. Theocharides, G. Link, N. Vijaykrishnan, and M. Irwin, "Implementing LDPC Decoding on Network-On-Chip", *Proceedings of the International Conference on VLSI Design, 2005* (VLSID 2005), pp. 134-137.

[14]     Jacob, Philip et al., "Predicting the Performance of a 3D Processor-Memory Stack." *IEEE Design and Test of Computers*, Nov. 2005, pp. 540-547.

[15]     C. Addo-Quaye, "Thermal-aware Mapping and Placement for 3D NoC Designs", *Proceedings of the IEEE International SoC Conference*, 2005, pp. 25-28.

[16]     Li Shang, L. Peh, A. Kumar, and N.K. Jha, "Temperature-Aware On-Chip Networks", *IEEE Micro*, vol. 26, no. 1, pp. 130-139.

[17]     H. Yu, Y. Shi, L. He, T. Karnick, "Thermal Via Allocation for 3D ICs Considering Temporally and Spatially Variant Thermal Power", *Proceedings of the 2006 ISLPED*, pp. 156-161.

[18]     Vangal et al., "An 80-Tile 1.28TFLOPS Network-on-Chip in 65nm CMOS", *Proceedings of IEEE International Solid-State Circuits Conference (ISSCC)*, 2007, pp. 98-99.

[19]     W.J. Dally and C.L. Seitz, "The Torus Routing Chip," Technical Report 5208:TR: 86, Computer Science Dept., California Inst. of Technology, pp. 1-19, 1986.

[20]     R. I. Greenberg and L. Guan, "An Improved Analytical Model for Wormhole Routed Networks with Application to Butterfly Fat Trees", *Proceedings of the International Conf. on Parallel Processing (ICPP 1997)*, pp. 44-48.

[21]     C. Grecu et al. "A Scalable Communication-Centric SoC Interconnect Architecture." *Proceedings of the 5th International Symposium on Quality Electronic Design*, 2004, pp. 343-348.

[22]     P. Guerrier and A. Greiner, "A Generic Architecture for On-Chip Packet-Switched Interconnections," *Proc. of Design and Test in Europe (DATE)*, pp. 250-256, Mar. 2000.

[23]     C. Grecu, P. P. Pande, A. Ivanov, R. Saleh, "Timing Analysis of Network on Chip Architectures for MP-SoC Platforms", *Microelectronics Journal*, Elsevier, Vol. 36, issue 9, pp. 833-845.

[24]   P. P. Pande, C. Grecu, M. Jones, A. Ivanov, R. Saleh, "Performance Evaluation and Design Trade-offs for Network on Chip Interconnect Architectures", *IEEE Transactions on Computers*, Vol. 54, no. 8, August 2005, pp. 1025-1040.

[25]   J. Duato, S. Yalamanchili, L. Ni, *Interconnection Networks – An Engineering Approach*, Morgan Kaufmann, 2002.

[26]   K. Park and W. Willinger, *Self-Similar Network Traffic and Performance Evaluation*. John Wiley & Sons, 2000.

[27]   D. R. Avresky, V. Shubranov, R. Horst, P. Mehra, "Performance Evaluation of the ServerNetR SAN under Self-Similar Traffic" *Proceedings of 13th International and 10th Symposium on Parallel and Distributed Processing*, April 12-16th, 1999, pp. 143-147.

[28]   Varatkar, G.V.; Marculescu, R., "On-chip traffic modeling and synthesis for MPEG-2 video applications", *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, Volume 8, Issue 3, June 2000 Page(s):335 - 339

[29]   Circuits Multi-Projects. http://cmp.imag.fr/

[30]   K.C. Saraswat et al., "Technology and Reliability Constrained Future Copper Interconnects—Part II: Performance Implications," *IEEE Trans. Electron Devices*, vol. 49, no. 4, pp. 598-604, Apr. 2002.

[31]   Pande, P.P.; Grecu, C.; Jones, M.; Ivanov, A.; Saleh, R.; "Effect of traffic localization on energy dissipation in NoC-based interconnect" *Proceedings of IEEE International Symposium on Circuits and Systems*, 23rd-26th May 2005, pp. 1774-1777.

[32]   J. A. Davis, R. Venkatesan, A. Kaloyeros, M. Beylansky, S. J. Souri, K. Banerjee, K. C. Saraswat, A. Rahman, R. Reif, and J. D. Meindl, "Interconnect Limits on Gigascale Integration (GSI) in the 21st Century" *Proceedings of the IEEE*, vol. 89, no 3, Mar. 2001, pp. 305-324.

[33]   D. Meeks., "Fundamentals of Heat Transfer in a Multilayer System," *Microwave Journal*, vol. 1, no. 1, Jan. 1992, pp. 165-172.

[34]   W. Huang, K. Sankaranarayanan, R. J. Ribando, M. R. Stan and K. Skadron, "An Improved Block-Based Thermal Model in HotSpot 4.0 with Granularity Considerations." *Proceedings of the Workshop on Duplicating, Deconstructing, and Debunking*, in conjunction with the *34th International Symposium on Computer Architecture (ISCA)*, June 2007.

[35]    M. F. Chang et al., "RF/Wireless Interconnect for Inter- and Intra-Chip Communications", *Proceedings of the IEEE*, vol. 89, no. 4, Apr. 2001, pp. 456-466

[36]    P. P. Pande, A. Ganguly, B. Belzer, A. Nojeh, A. Ivanov, "Novel Interconnect Infrastructures for Massive Multicore Chips – An Overview" *Proceedings of the 2008 ISCAS*, 2008.

[37]    R. G. Gallager. "Low-Density Parity-Check Codes", *IEEE Transactions on Information Theory*, Jan. 1962, pp. 21-28.

[38]    D. Mackay R. Neal. "Near Shannon limit performance of low density parity check codes", *IEEE Electronics Letters*, Vol.33, no. 6, March 1997, pp 457-458.

[39]    T. Theocharides, G. Link, N Vijaykrishnan, and M. J. Irwin, "Implementing LDPC Decoding on Network-on-Chip" *Proceedings of 18$^{th}$ International Conference on VLSI Design (VLSID '05)*, 2005.