

Hierarchical Network-on-Chip Design for Interposer-Based Systems and DNN Accelerators

by

Hesam Eddin Shabani

A Dissertation

Presented to the Graduate and Research Committee

of Lehigh University

in Candidacy for the Degree of

Doctor of Philosophy

in

Electrical Engineering

Lehigh University

(May 2023)

PREVIEW

© 2023 Copyright
Hesam Eddin Shabani

Approved and recommended for acceptance as a dissertation in partial fulfillment
of the requirements for the degree of Doctor of Philosophy

Hesam Eddin Shabani
Hierarchical Network-on-Chip Design for Interposer-Based Systems and DNN Accelerators

Defense Date

Dissertation Director

Approved Date

Committee Members:

Prof. Wujie Wen, Committee Chair

Prof. Xiaochen Guo

Prof. Zhiyuan Yan

Prof. Mahdi Nikdast

Acknowledgements

I would like to express my sincere gratitude to Prof. Guo, who allowed me to join the ECE-CompArch Lab research team and conduct my research under her supervision. She has kindly supported me during this project, and her insightful comments helped me to proceed with my research in the right direction leading to publications in prestigious conferences. I would also like to thank my committee members Prof. Wujie Wen, Prof. Zhiyuan Yan, and Prof. Mahdi Nikdast, for their suggestions and feedbacks.

”Dedicated to my beloved parents and my lovely brother” for their love, endless support, encouragement, and sacrifices.

Contents

Acknowledgements	iv
List of Tables	viii
List of Figures	ix
Abstract	1
1 Introduction	3
2 ClusCross	6
2.1 Introduction	7
2.2 Background and Related Work	9
2.2.1 Interposer-Based Interconnection Networks	9
2.2.2 Conventional NoC Topologies	11
2.2.3 Silicon Interposer-Based Topologies	12
2.3 ClusCross Topologies	14
2.3.1 ClusCross Topologies Designed for Interposer-Based Systems .	15
2.3.2 General-Purpose ClusCross Topology	17
2.4 Evaluation Results	18
2.4.1 Experimental Setup	19
2.4.2 System Performance Evaluation Using Injected Traffic Patterns	22
2.4.3 System Performance Evaluation Using PARSEC	24

2.4.4	Power and Area Evaluation	25
2.5	Conclusion	27
3	HIRAC: A <u>H</u>ierarchical <u>A</u>ccelerator with Sorting-based Packing for SpGEMMs in DNN Applications	29
3.1	Introduction	30
3.2	Background and Related Work	33
3.2.1	Preprocessing Methods	33
3.2.2	Accelerator Designs for DNNs and SpGEMMs	35
3.3	System Design Overview	38
3.4	The SorPack Algorithm	39
3.4.1	Impact of the SorPack Steps on the Hardware Design	41
3.5	The HIRAC design	45
3.5.1	PE Array	47
3.5.2	Interconnection Network	49
3.5.3	On-Chip SRAM	50
3.6	Experimental setup	50
3.7	Evaluation Results	52
3.7.1	The SorPack Algorithm Evaluation	52
3.7.2	The HIRAC Architecture	53
3.7.3	Sensitivity Study	58
3.7.4	End-to-End Evaluations of a DNN Workload	60
3.7.5	The GPUSorPack Version	62
3.8	Conclusion	64
4	Heterogeneous Accelerator for Different Sparsity Ranges	66
4.1	Introduction and Background	66
4.2	Heterogeneous HIRAC	70
4.3	Analytical Results	72

4.4 Future Work	74
Bibliography	76
Biographical Sketch	87

PREVIEW

List of Tables

2.1	A comparison of interposer-based topologies.	16
2.2	A comparison of general-purpose topologies.	18
2.3	Network Parameters.	20
2.4	Architecture Parameters.	22
2.5	PARSEC V2.1 applications.	22
2.6	Average packet latency of memory traffic in reply-and-request batch mode.	24
3.1	The same-cycle/col merging percentage comparison of SorPack without sorting vs. with sorting.	45
3.2	Architectural parameters of HIRAC.	52
3.3	Common DNN workloads dimensions.	52
3.4	Area, Power, and Cycle runtime for different PE subarray sizes. . . .	60
4.1	Area, Power, and Cycle runtime estimation of the heterogeneous design (Sparse PE array %= 83, Dense PE array %= 17) over the HIRAC for the dynamic sparse attention matrix.	74

List of Figures

2.1	An illustration of a 64-core system composed of four 16-core processor chips and four HBM DRAMs.	10
2.2	Illustrations of existing misaligned interposer-based topologies. (a) FoldedTorus x and (b) ButterDonut x.	14
2.3	Illustrations of two versions of ClusCross topology. (a) ClusCross x-v1 and (b) ClusCross x-v2.	15
2.4	An illustration of a 64-node general-purpose ClusCross topology. . . .	18
2.5	Saturation throughput of topologies for different numbers of VCs with the shortest-path routing algorithm under coherence traffic.	21
2.6	Average packet latency and saturation throughput of different network topologies for coherence traffic.	23
2.7	Average packet latency and saturation throughput of different network topologies for memory traffic.	24
2.8	Total simulation runtime normalized to CMesh x.	25
2.9	Average packet latency normalized to CMesh x.	26
2.10	Power consumption breakdown of topologies.	26
2.11	Area breakdown of different topologies.	27
3.1	An overview of proposed HW/SW co-design architecture composed of preprocessing and accelerator parts.	37
3.2	An example of applying the SorPack in the streaming and stationary matrices.	43

3.3	Condensing factor and the percentage of partial sums with the same-cycle/col merging for different matrix partitioning sizes. The result is from the 100×100 matrix size, and the sparsity of the stationary and streaming matrices are 50% and 70%, respectively.	44
3.4	Examples of partial sums to be merged produced under different situations. A -1 col or row id means a zero element that needs to be skipped.	46
3.5	An illustration of the PE subarray of the HIRAC.	48
3.6	An illustration of the optimization to compute the tiled output matrix. The red block represents submatrices in level 1, and the blue block represents submatrices in level 2. Shaded submatrices in (a) and (b) are the ones to load and stream to compute tiles 1 and 2.	51
3.7	The runtime comparison of the SorPack and the collision-aware algorithm for different matrix sizes (a-c) and sparsities.	54
3.8	Effects of the Sorting step on system performance. The matrix size is 100×100	55
3.9	Speedup comparison of the HIRAC and SIGMA over the Google TPU for representative matrices in DNN workloads (a-f). Sta:80% Str:90% means the stationary matrix has 80% zeros, and the streaming matrix has 90% zeros.	56
3.10	Area breakdown of the HIRAC.	57
3.11	Speedup over the SIGMA for different matrix Partitioning size P of the SorPack in two different matrix sizes of (a) 128×128 and (b) the Set 3 workload. The sparsity of the stationary and streaming matrices are 50% and 70%, respectively.	59
3.12	Results of cycle runtime, SRAM bank conflicts and area for the different number of the banks.	60
3.13	An end-to-end runtime evaluation using GNMT v2 [26].	61

3.14	The runtime comparison of running (a) SorPack on the CPU and (b) GPUSorPack on the GPU for packing the dynamic sparse attention matrix.	64
3.15	The runtime percentage breakdown of running (a) SorPack on the CPU and (b) GPUSorPack on the GPU for packing the dynamic sparse attention matrix.	65
4.1	Tensor sparsity ranges of different workload domains [21].	67
4.2	An illustration of the Heterogeneous HIRAC design.	70
4.3	The Heterogeneous PE subarray.	71
4.4	The runtime estimation is when only the sparse PE arrays exist compared to the only dense PE arrays.	73
4.5	The runtime estimation of the heterogeneous design over the HIRAC from the analytical simulator. The portion of Dense/Sparse PE array is matched with the portion of Dense/Sparse partitions.	74

Abstract

Network-on-Chip (NoC) is a crucial chip multiprocessor component to communicate between many nodes. Continued increases toward multicore and manycore scalability have led to performance challenges of NoCs because of the increasing network diameter. Also, up to 30% of the chip's overall power budget is contributed by NoCs in modern chips [25], and on-chip power consumption exceeds the total power budget by increasing cores in the general-purpose chip multiprocessors. The hierarchical design approach is a promising solution to offer straightforward paths to improve performance and minimize power consumption. Hierarchical on-chip interconnection design is suitable for large systems by providing routes with shorter hop counts in the network. The hierarchical design approaches generally require inter-chip communication; however, as the number of small chips increases, the chip-to-chip communication becomes a performance bottleneck. Therefore, the interconnection network should be carefully designed to provide the shortest paths for as many source-destination pairs and avoid network congestion and minimum area and power consumption overhead. Furthermore, many widespread applications like modern AI systems require a large amount of data to support the computation, creating considerable data movement for on-chip and off-chip communications. Therefore, general-purpose on-chip network designs could not be appropriate for providing power efficiency in large-scale AI systems. Application-specific on-chip networks are proposed to leverage in the embedded systems to address the mentioned challenges in the general purpose. Therefore, hierarchical interconnect approaches such as tile-based architectures are well

studied and applied frequently in deep-learning accelerator designs. In this dissertation, three projects have been proposed. The first work proposes a new hierarchical topology design, ClusCross, to improve multicore interconnection networks on silicon interposer-based systems. The key idea is to treat each small chip as a cluster and use cross-cluster long links to increase bisection width and decrease average hop count without increasing the number of ports in the routers. The second work proposes a HW/SW co-design architecture to compute SpGEMM efficiently without requiring complex interconnection networks. A novel fast-packing algorithm, SorPack, is proposed to convert a sparse matrix into a dense matrix that increases PE utilization. Additionally, The HIRAC, a novel hierarchical accelerator, is proposed for executing Sparse GEMM and provides a scalable system that maximizes the parallelism of the PEs. The last chapter presents the heterogeneous design approach that can be used in applications requiring both sampled SpGEMM and Highly SpGEMMs and efficiently covering the higher sparsity ranges. The proposed heterogeneous design achieves 24% faster runtime estimation over HIRAC for a dynamic sparse attention matrix extracted from a state-of-the-art sparse attention model layer.

Chapter 1

Introduction

The increasing number of cores challenges the scalability of chip multiprocessors due to the requirements of high compute throughput systems. The hierarchical design approach is a promising solution to offer straightforward paths to improve scalability.

For example, multi-chip-modules packaging approaches such as silicon interposer-based systems [48] apply the idea of disintegration by partitioning a large chip into multiple smaller chips and using silicon interposer-based integration (2.5D) or organic substrates to connect these smaller chips. Since a small chip's verification, logic, and physical design are more convenient than a large chip, design costs are reduced compared to a large monolithic die. Also, this approach can improve overall yield because a smaller chip has fewer components, resulting in less likelihood of catching defects. In addition, modularity provides multiple smaller chips instead of a big chip; hence a tiny defective chip can be replaced at a lower cost when re-integrated through interposers. Moreover, as the number of cores grows, Network-on-chip (NoC)'s performance becomes limited because of the increasing network diameter. Hence, hierarchical design for on-chip interconnection provides routes with shorter hop counts from source to destination, which is more proper for large systems.

The hierarchical design approaches generally require inter-chip communication; however, as the number of small chips increases, the chip-to-chip communication

becomes a performance bottleneck. Therefore, the interconnection network should be carefully designed to provide the shortest paths for as many source-destination pairs, avoid network congestion, and minimize area and power consumption overhead.

Meanwhile, on-chip power consumption exceeds the total power budget by increasing cores in the general-purpose chip multiprocessors. The reason is the limitation of the power delivery network and thermal dissipation capability. Moreover, many widespread applications like modern AI systems require a large amount of data to support the computation, creating considerable data movement for on-chip and off-chip communications. As a result, it causes more energy consumption because data movement can consume more energy than computation. Therefore, general-purpose on-chip network designs could not be appropriate for providing power efficiency in large-scale AI systems.

Furthermore, application-specific on-chip networks are proposed to leverage in the embedded domain, which generates on-chip interconnection corresponding to the application’s communication graph to address the challenges in general-purpose one. For example, the interconnection network in the specialized accelerator’s design, tailored to the dataflow in Deep Neural Networks (DNNs) applications, is essential to meet the system requirements. The goal is to enable a large-scale system and extract the maximum possible parallelism from the available processing elements (PEs). Therefore, hierarchical interconnect approaches such as tile-based architectures are well studied and applied frequently in deep-learning accelerator designs. [7], [8], [45], [48].

For instance, Simba [48] uses a hierarchical interconnection design consisting of a Mesh NoC topology and a network-on-package (NoP). The Mesh NoC connects multiple processing elements (PEs) efficiently on the same chiplet. The NoP connects chiplets on the same package to provide the design for a large-scale system.

Additionally, SIGMA [45] leverages hierarchical design to interconnect different PEs to make the maximum possible parallelism of the in-use multipliers. SIGMA design includes NoC design and a combination of Flex-DPE units to construct a

Flex-DPU, which each Flex-DPU is for running one general matrix-matrix multiplication (GEMM). Hence, Multiple Flex-DPUs can consider in parallel to run multiple GEMMs. Also, the NoC is responsible for providing interconnection among the Flex-DPEs like the idea of the other tile-based architectures [17], [48]. However, the SIGMA has a 37.7% area overhead and 82% more power consumption as compared to the TPU because of the high flexibility and complexity of interconnection networks in the non-blocking distribution and reduction networks. Additionally, the utilization of SIGMA is determined by the sparsity of the streaming matrix. Thereby, the SIGMA design is not performance efficient for a sparser streaming matrix.

Chapter 2

ClusCross

As the number of small chips increases in the hierarchal design of the interposer-based system, chip-to-chip communication becomes a performance bottleneck. Hence, the interconnection network design should be a target to improve system performance carefully.

This work proposes a new network topology, ClusCross, to improve multicore interconnection networks on silicon interposer-based systems. The key idea is to treat each small chip as a cluster and use cross-cluster long links to increase bisection width and decrease average hop count without increasing the number of ports in the routers. Synthetic traffic patterns and real applications are simulated on a cycle-accurate simulator. Network latency reduction and saturation throughput improvement report compared to previously proposed topologies. Two versions of the ClusCross topology are presented. One version of ClusCross has a 10% average latency reduction for coherence traffic compared to the state-of-the-art network-on-interposer topology, the misaligned ButterDonut. The other version of ClusCross has a 7% and a 10% reduction in power consumption as compared to the FoldedTorus and the ButterDonut topologies, respectively.

2.1 Introduction

As the number of transistors increases, more processor cores can be integrated into a Chip Multi-Processor (CMP) to boost the computation throughput. With the invention of High Bandwidth Memories (HBMs) [42], memory bandwidth can be significantly improved by connecting multiple 3D-stacked DRAMs to processor chips through silicon interposers to satisfy the overall demands from the processor’s cores. Each processor core, however, might need to access multiple memory locations and the increased number of cores also escalate coherence traffic among the cores. The on-chip networks are facing fundamental challenges to enable the scalability of the CMPs and to satisfy both the coherence and memory traffic demands. Meanwhile, with the increasing number of cores, on-chip power consumption is about to exceed the total power budget due to the limitation of the power delivery network and thermal dissipation capability. On-chip network designs have to be power efficient to meet the system power constraint.

Inspired by silicon interposer-based memory integration (*e.g.*, HBM), which is also referred to as 2.5D integration, recent studies [32] proposed the idea of disintegration by taking apart a large system into smaller parts by using the interposer-based integration to improve overall yield. This is because a smaller chip has fewer components and is less likely to catch defects. Having multiple smaller chips instead of a big chip also provides modularity, and a small defective chip can be replaced at a lower cost when re-integrated through interposers. Nevertheless, as multiple smaller chips are integrated through the interposers, the amount of chip-to-chip communications increases. Heavy traffic through the interposers can become a performance bottleneck [32]. Moreover, any processor core can access different parts of the on-chip memories. Hence, the memory traffic also needs to pass across different chips through the interposers. Even though disintegration can improve the yield and reduce the fabrication cost, the interconnection network can become a performance bottleneck if it is not carefully designed to overcome the challenges posed by the interposer-based

multi-chip systems.

Topology is one of the most important elements in interconnection network design, which directly influences network performance. In interposer-based systems, memory traffic can compete with coherence traffic for bandwidth [32]. The network topology should be designed to reduce such contention by increasing the number of links and bandwidth on segments critical to both memory and coherence traffic.

This work proposes a new interconnection network topology, ClusCross, for silicon interposer-based multi-chip systems. This topology is based on the idea of clustering. In order to decrease the network diameter and increase the cross-chip bandwidth, ClusCross maps a cluster of routers onto each small chip and increases the number of cross-cluster long links. As a result, the proposed topology can increase path diversity and bisection bandwidth, which can help to reduce contentions between memory and coherence traffic. In addition, cross-cluster long links can effectively reduce the hop counts for long-distance communication in both memory and coherence traffic.

The main contributions of this work include the following:

- Two versions of ClusCross on-chip network topology are proposed to improve network performance in NoC-on-interposer systems through decreasing average hop count and increasing cross-chip bandwidth by leveraging long links.
- Performance and cost of the proposed ClusCross topologies are evaluated and compared against other existing topologies using synthetic memory and coherence traffic.
- System performance of ClusCross topologies is evaluated using the PARSEC suite traces appropriate for CMP assessment.

The rest of the chapter is organized as follows: In Section 2.2, a brief overview of the interconnection networks based on silicon interposers is provided, and related work for both conventional NoC topologies and topologies for silicon interposer systems is

discussed. Section 2.3 presents the structure of the ClusCross and two versions of this topology. In Section 2.4, evaluation results are shown using both synthetic traffic patterns and real applications. The proposed topologies are compared against other topologies designed for interposer-based systems. Section 2.5 concludes the work.

2.2 Background and Related Work

2.2.1 Interposer-Based Interconnection Networks

Technology scaling does not benefit wires as much as it does transistors [11]. On-chip communication becomes a bottleneck for both power consumption and performance. Three-dimensional (3D) integration promises to bring processing elements and memory components physically close to each other to reduce wire distance and overcome the communication bottleneck. True 3D integration, however, requires through-silicon vias (TSVs), which are complicated to implement on processor dies and might introduce severe thermal issues and die yield reduction [32], [14]. As an alternative, individual chips can be connected to the silicon interposer layer through micro-bumps. Hence, memory and processor chips can be connected through a layer of silicon interposers on a substrate die to increase memory bandwidth.

Since interposer integration does not need TSVs in the silicon interposer layer, higher die yield and additional routing capabilities are provided for the system [44]. In addition, interposer-based systems have lower manufacturing and R&D costs as compared to the true 3D integration [44]. Although the physical design of the interposer integration also has technology-related challenges, such as thermal management and pin assignment [59], these challenges are solvable in the near term [44]. Consequently, interposer-based systems are the most promising near-term solution for die-stacking integration. Several commercial products of interposer-based ICs are already on the market [38], [39]. For example, the HBM uses TSVs to integrate stacks of DRAM dies and connects the DRAM stacks to the processor die using silicon interposers. Multi-