# A Flexible Hybrid Interconnection Design for High-Performance and Energy-Efficient Chiplet-Based Systems

Md Tareq Mahmud<sup>®</sup>, Graduate Student Member, IEEE, and Ke Wang<sup>®</sup>, Member, IEEE

Abstract—Chiplet-based multi-die integration has prevailed in modern computing system designs as it provides an agile solution for improving processing power with reduced manufacturing costs. In chiplet-based implementations, complete electronic systems are created by integrating individual hardware components through interconnection networks that consist of intra-chiplet network-onchips (NoCs) and an inter-chiplet silicon interposer. Unfortunately, current interconnection designs have become the limiting factor in further scaling performance and energy efficiency. Specifically, inter-chiplet communication through silicon interposers is expensive due to the limited throughput. The existing wired Networkon-Chip (NoC) design is not good for multicast and broadcast communication because of limited bandwidth, high hop count and limited hardware resources leading to high overhead, latency and power consumption. On the other hand, wireless components might be helpful for multicast/broadcast communications, but they require high setup latency which cannot be used for one-to-one communication. In this paper, we propose a hybrid interconnection design for high-performance and low-power communications in chiplet-based systems. The proposed design consists of both wired and wireless interconnects that can adapt to diverse communication patterns and requirements. A dynamic control policy is proposed to maximize the performance and minimize power consumption by allocating all traffic to wireless or wired hardware components based on the communication patterns. Evaluation results show that the proposed hybrid design achieves 8% to 46% lower average end-to-end delay and 0.93 to 2.7× energy saving over the existing designs with minimized overhead.

*Index Terms*—Computer architecture, chiplets, network-onchip (NoC), hybrid interconnection, wireless.

#### I. INTRODUCTION

HIPLET-BASED systems have become a prevalent technique in modern System-on-Chip (SoC) design to mitigate the limitations of fabrication, and enhance the design flexibility to achieve better performance and lower energy consumption. In a chiplet-based design [1], a monolithic chip is split into isolated chips connected through the interposer for inter-chip data sharing. Each chip is integrated with the conventional NoC for efficient on-chip communication. As the technology scales with more shared components integrated into the chiplet-based circuits, inter-chiplet and intra-chiplet communications become critical factors for further scaling the performance and energy efficiency of the computing systems.

Received 13 April 2024; revised 28 July 2024; accepted 3 October 2024. Date of publication 9 October 2024; date of current version 19 November 2024. This work was supported in part by the National Science Foundation under Grant CCF-2245950 and Grant CNS-2321225. (Corresponding author: Md Tareq Mahmud.)

The authors are with the Department of Electrical and Computer Engineering, University of North Carolina at Charlotte, Charlotte, NC 28223 USA (e-mail: mmahmud1@charlotte.edu; ke.wang@charlotte.edu).

Digital Object Identifier 10.1109/LCA.2024.3477253

Existing interconnection designs including the silicon interposer design and the NoC design are not sufficient. Interposers have huge bandwidth limitations, and the throughput can be very low because of the nature of the 2.5D packaging for chipletbased systems. The overall runtime frequency for the interposer can not be very high due to the power wall issues, and hotspots can be created. Besides, conventional NoC design might not be suitable for all application requirements and communication patterns due to its lower bandwidth, higher hop count and limited hardware resources (e.g., buffer). For example, in Artificial Intelligence (AI) oriented applications that involve lots of multicasting and broadcasting traffic, the traditional NoC will induce excessive timing overhead for multi-hop routing computation, resource allocation, and suffers from congestion [2]. Alternatively, wireless components can supplement the convention NoC for faster broadcasting and can create a second one-hop communication channel for inter-chiplet communication. However, simply integrating the wireless components into chiplet-based systems is challenging because they have high setup latency and cannot be used for one-to-one communication, even though they may be useful for multicast/broadcast communications [3]. It will have some area consumption also the routing mechanism should be re-designed.

We propose a flexible hybrid interconnection design that integrates electrical interconnections and wireless interconnections (WIs) that can adapt to diverse packet patterns and application requirements for efficient inter-chiplet and intra-chiplet communication. Specifically, we design a low-cost architecture that integrates both wired interconnects and wireless components, and can automatically select transmission channels according to communication patterns and available resources to achieve improved performance and power. In each chiplet, we integrate a novel hybrid router designed with some lightweight wireless components. Besides, we design a novel flow control policy and routing mechanism to adaptively allocate the hardware resources for different traffic according to their communication patterns (e.g., broadcasting, multicasting, and unicasting). Simulation results show that the proposed flexible hybrid design achieves an average of 22% lower average end-to-end delay and average 1.8 × energy savings when compared to conventional interconnects.

# II. PROPOSED INTERCONNECTION DESIGN

# A. Chiplet-Based Hybrid Interconnection

The proposed hybrid interconnection design integrates both wired and wireless networks. The wired network handles both intra-chiplet and inter-chiplet communication, while the wireless network is limited to inter-chiplet communication. The design, shown in Fig. 1, uses a  $4\times4$  2D mesh topology on an active interposer, connecting 16 chiplets, each with 16 processing

1556-6056 © 2024 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.

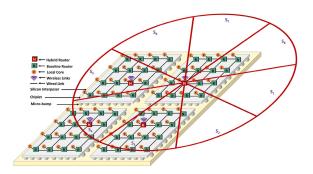


Fig. 1. Proposed flexible hybrid interconnection design with reconfigurable four-element planar array antenna.

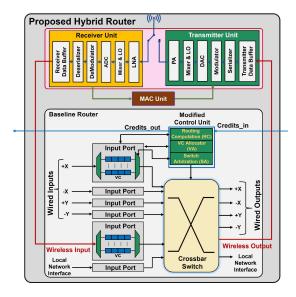


Fig. 2. Proposed hybrid router architecture.

elements (PEs). Among these, 15 PEs connect to baseline routers (BRs), and one PE connects to a hybrid router (HR). The HR operates like a BR for intra-chiplet communication and manages wireless inter-chiplet communication to reduce congestion and enable parallel transmission among other HRs. The flexible design allows adjustments in chiplet size, number, and HR placement to meet application-specific performance demands.

# B. Hybrid Router Architecture

We design a novel HR that consists of a baseline wormhole router, wireless interconnect (WI) and a MAC table. It has a modified control unit to adapt to the proposed novel flow control and routing operation. One wireless input port and one output port are added to the crossbar switch to enable wireless communication for this hybrid design. The transceiver module of the WI has a transmitter unit and a receiver unit (as shown in Fig. 2). The MAC table stores the wireless channel states, and is utilized for a fair channel access mechanism.

# C. Proposed Antenna Model

Each HR uses a single radio transceiver, operating in half-duplex mode. A reconfigurable four-element planar array antenna provides 360° main beam scanning, with beam adjustments in 45° increments for communication across eight directions (east, west, north, south, and the four diagonals) between chiplets (as shown in Fig. 1) [4]. The antenna can



Fig. 3. Flit format for wireless MAC operation.

switch between unique frequency channels but access only one at a time. It supports both omnidirectional and directional modes, with directional transmission enabling high-bandwidth, interference-free single-hop communication [5].

### D. Medium Access Control (MAC) Policy

Each wireless interface (WI) in the hybrid design is assigned a unique address for managing flow, routing, and access to shared wireless channels. The multi-channel approach allows nodes to communicate simultaneously on different frequencies, boosting capacity and reducing latency [6]. Dual-channel communication supports both broadcasting and unicast via omnidirectional and directional links. In this design, two non-overlapping channels (CH-1 and CH-2) are used: CH-1 handles omnidirectional communication, broadcasting data or control packets, with control packets managing collision-free access to CH-2 for directional communication. Each WI maintains a MAC table to store channel access information for collision-free communication. Transceivers are usually set to CH-1; WIs use CSMA/CA to access it. Before broadcasting, a WI senses CH-1. If it's free, the WI checks its MAC table to confirm all receivers are tuned to CH-1. If so, the data is broadcast; if not, the sender waits for a small time ( $\Delta$ ) and retries the CSMA/CA process until successful.

If a wireless interface (WI) wants to communicate directly with another, both WIs must switch their antennas to CH-2. Before switching, they ensure no other WIs communicate in the same direction on CH-2. To coordinate, they first negotiate on CH-1 by broadcasting Request/Reply control packets (as shown in Fig. 3). Other WIs, upon overhearing these packets, wait for the specified channel access time. Once the directional communication on CH-2 ends, both WIs switch back to CH-1. Notably, data packet broadcasts have priority over control packet broadcasts on CH-1.

# III. FLOW CONTROL AND ROUTING STRATEGY

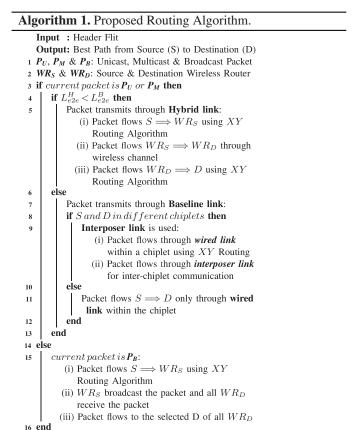
The proposed flow control and routing strategy for intrachiplet and inter-chiplet data communication of the proposed design depends on the position of the destination of a packet and the end-to-end (e2e) latency for the transmission of that packet from its source to the destination. A packet can be transmitted either through the baseline wired link or through the hybrid link which is composed of both wired and wireless links. This decision is made by the source by calculating and comparing the e2e latency for both transmission paths.

We assume that  $\xi$  is the current packet injection rate.  $\mho_W$  and  $\mho_{RF}$  are the current utilization of the wired and wireless link, respectively. Hence, the queuing delay for the communication using the baseline links and hybrid links is as follows:

$$L_Q^B = f(\xi, \mho_W) \tag{1}$$

$$L_Q^H = f(\xi, \mho_W, \mho_{RF}) \tag{2}$$

When an unicast or multicast packet is injected by a router into the network, the routing computation unit estimates the e2e latency for transmission that packet. The source router calculates both the e2e latency for baseline transmission  $L_{e2e}^B$  and the transmission through the hybrid link  $L_{e2e}^H$  using 1 and 2,



respectively. The baseline link is composed of wired links and interposer links, whereas the hybrid link is composed both of wired and wireless links.

If the e2e latency for hybrid link transmission is less than baseline transmission (i.e.,  $L_{e2e}^H < L_{e2e}^B$ ), then data is transmitted through the hybrid link. Packet flows from the source baseline router (S) to the source wireless router  $(WR_S)$  of the source chiplet using the XY Routing. Then the  $WR_S$  transmits the packet to the destination  $(WR_D)$  of the destination chiplet through the wireless channel. After receiving the data packet by the  $WR_D$ , the data packet flows to the final destination baseline router (D) using the XY Routing.

Otherwise, data packet transmission occurs through wired baseline links (when  $L_{e2e}^{H} > L_{e2e}^{B}$ ). While using the baseline link, if the baseline source (S) and destination (D) routers are in the same chiplet then the packet flows only through a wired link within the chiplet from S to D. Again if, S and D are in different chiplets then packet flows through the wired link within the chiplet using XY Routing, and through interposer link for inter-chiplet communication. Fig. 4 shows the flit format for flow control. If the packet is a broadcast packet, then S directly transmits using the hybrid links. Algorithm 1 contains the pseudo-code of the proposed routing algorithm.

#### IV. EVALUATION AND ANALYSIS

To evaluate the performance of our proposed flexible hybrid interconnection design, we use a customized version of the cycle-accurate simulator *BookSim2.0* [7] which integrates the proposed hybrid router, MAC policy, flow control, and routing strategy. Benchmark traces from Netrace [8] capture workloads using PARSEC [9] applications. For energy analysis, ORION 3.0 [10] is employed. Simulations consider a 4 × 4 chiplet

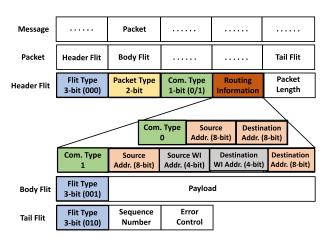


Fig. 4. Flit format for flow control.

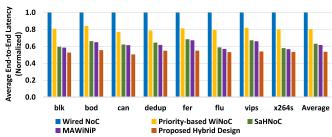


Fig. 5. Comparison of average end-to-end latency (normalized).

configuration connected via a  $16 \times 16$  2D mesh network with 16 hybrid routers, where each hybrid router has 4 virtual channels with buffer size 4 flits. The Tx/Rx buffer size of the wireless transceiver is 16 flits. We also considered 2 wireless channels with 10 Gbps maximum data rate per channel achieved utilizing 256-QAM modulation. We compare our design with traditional wired interconnect design and hybrid designs like Priority-based WNoC [11], SaHNoC [12] and MAWiNiP [3] where we developed the baselines of the compared designs as chiplet-based baselines and these chiplets are also attached to a conventional interposer design. Moreover, we justify the energy efficiency of the proposed hybrid design based on chiplet sizes  $(n \times n)$ .

The comparison of normalized average end-to-end (e2e)packet latency for different benchmark applications execution is shown in Fig. 5, and the normalization is done to the baseline wired NoC. Compared to the baseline wired NoC, proposed hybrid design and MAWiNiP obtain a 46% and 38% lower average e2e latency, respectively. Here, the proposed hybrid design achieves 8% reduced e2e latency than MAWiNiP design because even if for the communication between two edge PEs of two different chips MAWiNiP uses wireless communication and costs additional e2e latency for channel negotiation/setup and data transmission. The priority-based WNoC achieves a 19% lower average e2e latency than wired NoC due to employing a priority-based dynamic MAC mechanism. However, priority-based WNoC suffers from a 27% increased average e2e latency than the proposed design for its static architecture without considering the optimal placement of a WR. Besides, the proposed hybrid design obtains a 9% reduced average e2elatency over the SaHNoC; because ShHNoC lacks the adaptive flow control and routing mechanisms.

strategy. Benchmark traces from Netrace [8] capture workloads using PARSEC [9] applications. For energy analysis, ORION 3.0 [10] is employed. Simulations consider a  $4 \times 4$  chiplet Authorized licensed use limited to: University of North Carolina at Charlotte. Downloaded on August 29,2025 at 14:18:32 UTC from EEE Xplore. Restrictions apply.

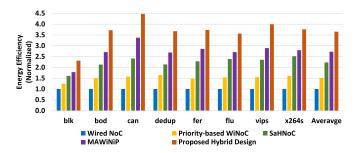


Fig. 6. Comparison of energy efficiency (normalized).

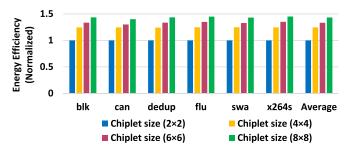


Fig. 7. Comparison of energy efficiency based on different chiplet size of proposed hybrid NoC.

is comprised of the overall power consumption. A benchmark application's complete execution time  $(T_{Ex})$  is multiplied by the overall power consumption for that benchmark application's execution resulting in the total energy consumption. The energy efficiency  $(E_{Ef})$  estimated as:

$$E_{Ef} = \frac{1}{T_{Ex} \times (P_S + P_D)} \tag{3}$$

When compared to other hybrid designs, proposed flexible hybrid design's unique architecture with compact wireless components, reliable flow management and routing mechanism ensures transmission reliability while using fewer hardware resources for long-distance wireless communication along with one-time energy consumption for multicast and broadcast. As depicted in Fig. 6), the proposed flexible hybrid design provides an average of 2.7 times, 2.17 times, 1.43 times and 0.93 times improved energy efficiency than the baseline wired NoC, Priority Based WNoC, SaHNoC and MAWiNiP, respectively.

Fig. 7 shows the impact on the energy efficiency of the proposed flexible hybrid design with different chiplet sizes. When we use  $8 \times 8$  size chiplets proposed design gains 10%, 18% and 43% energy efficiency than using the  $6 \times 6$ ,  $4 \times 4$ and  $2 \times 2$  size chiplets, respectively. Because smaller chiplets fit numerous components into a limited area, they could have greater power density and result in excessive heat generation. To mitigate this overheating issue, the total energy efficiency may be impacted for integrating complex cooling systems. Besides, smaller chiplet sizes require longer communication links among chiplets resulting in increased capacitance and resistance, which raises energy consumption due to attenuation. Hence, it needs more power for signal amplification or regeneration to overcome the attenuation. By reducing the energy needed for data transmission between chiplets, larger chiplets can result in higher power efficiency.

The hardware cost of the proposed hybrid interconnection design is evaluated by synthesizing both the baseline and hybrid routers using 45 nm technology with the Synopsis Design

Compiler. The baseline router, with five ports and four virtual channels per port, has a buffer size of 389488  $\mu m^2$ . The key components include a switch allocator (6589  $\mu m^2$ ), virtual channel allocator (9062  $\mu m^2$ ), and crossbar (29806  $\mu m^2$ ). The area required for channels is 95675  $\mu m^2$ . The total area for 256 baseline routers is  $135.88 \times 10^6 \mu m^2$ , plus an interposer area of  $20.98 \times 10^6 \mu m^2$ , making the total area required for baseline NoC is  $156.8 \times 10^6 \mu m^2$ . The main overhead for the hybrid design is due to wireless transceiver components and MAC unit. Each hybrid router includes an antenna (1.125  $\times 10^6 \mu m^2$ ), transceiver components (0.46  $\times 10^6 \mu m^2$ ) and a MAC unit (193.25  $\mu m^2$ ). With 16 hybrid routers, the additional area is  $25.36 \times 10^6 \mu m^2$ . Thus, the total area for the hybrid design becomes  $182.16 \times 10^6 \mu m^2$ , with wireless components adding a 13.9% area overhead.

#### V. CONCLUSION

This paper proposes a flexible hybrid interconnection design for high-performance and low-power communications in chiplet-based systems. The proposed design consists of both wired and wireless interconnects that can adapt to diverse communication patterns and requirements. A dynamic control policy is proposed to maximize performance and minimize power consumption by allocating all traffic to wireless or wired hardware components based on communication patterns. Finally, the simulation results show that the proposed hybrid design achieves 22% average end-to-end delay reduction and 1.8  $\times$  energy saving over the existing designs.

#### REFERENCES

- H. Zheng, K. Wang, and A. Louri, "A versatile and flexible chiplet-based system design for heterogeneous manycore architectures," in *Proc. 57th* ACM/IEEE Des. Automat. Conf., 2020, pp. 1–6.
- [2] S. S. Rout et al., "2DMAC: A sustainable and efficient medium access control mechanism for future wireless nocs," *J. Emerg. Technol. Comput.* Syst., vol. 19, no. 3, pp. 1–25, jun 2023.
- [3] M. M. Ahmed, N. Mansoor, and A. Ganguly, "An asymmetric, one-to-many traffic-aware mm-wave wireless interconnection architecture for multichip systems," *IEEE Trans. Emerg. Topics Comput.*, vol. 10, no. 1, pp. 324–338, First Quarter, 2022.
- [4] P. Baniya, A. Bisognin, K. L. Melde, and C. Luxey, "Chip-to-chip switched beam 60 GHz circular patch planar antenna array and pattern considerations," *IEEE Trans. Antennas Propag.*, vol. 66, no. 4, pp. 1776–1787, Apr. 2018.
- [5] S. H. Gade et al., "Energy efficient chip-to-chip wireless interconnection for heterogeneous architectures," ACM Trans. Des. Autom. Electron. Syst., vol. 24, no. 5, Jul. 2019.
- [6] D. Zhao et al., "cm3WiNoCs: Congestion-aware millimeter-wave multichannel wireless noc," *IEEE Access*, vol. 8, pp. 24098–24107, 2020.
- [7] N. Jiang et al., "A detailed and flexible cycle-accurate NoC simulator," in Proc. IEEE Int. Symp. Perf. Analys. Syst. Soft., 2013, pp. 86–96.
- [8] J. Hestness et al., "Netrace: Dependency-driven trace-based network-onchip simulation," in *Proc. 3rd Int. Workshop Netw. Chip Architectures*, NY, USA: ACM, 2010, pp. 31–36.
- [9] C. Bienia et al., "The PARSEC benchmark suite: Characterization and architectural implications," in *Proc. ACM 17th Int. Parallel Architectures Arch. Compilation Techn.*, NY, USA, 2008, pp. 72–81.
- [10] A. B. Kahng, B. Lin, and S. Nath, "ORION3.0: A comprehensive NoC router estimation tool," *IEEE Embed Syst. Lett.*, vol. 7, no. 2, pp. 41–45, Jun. 2015.
- [11] Y. Ouyang et al., "Architecting a priority-based dynamic media access control mechanism in wireless network-on-chip," *Microelectronics J.*, vol. 116, pp. 1–9, 2021.
- [12] A. Alagarsamy et al., "Sahnoc: An optimal energy efficient hybrid networks-on-chip architecture," *J. Supercomput.*, vol. 79, no. 6, pp. 6538–6559, Nov. 2022.