**RESEARCH ARTICLE**

# OWBM: OSNR-Aware Wavelength Allocation and Branching Methods for Multicast Routing in Custom Topology-Based Optical Network-on-Chips

**YONG WOOK KIM[1], (Student Member, IEEE), AND TAE HEE HAN[ID][2], (Senior Member, IEEE)**
[1]Department of Electrical and Computer Engineering, Sungkyunkwan University, Suwon 16419, South Korea
[2]Department of Semiconductor Systems Engineering, Sungkyunkwan University, Suwon 16419, South Korea

Corresponding author: Tae Hee Han (than@skku.edu)

**ABSTRACT** In light of the rapid advancements in big data and artificial intelligence applications, heterogeneous computing (HGC) platforms that integrate diverse computing units have gained traction with the aim of achieving energy efficiency and high performance. A custom topology-based optical network-on-chip (ONoC) that provides unparalleled diversity between computing nodes is expected to be the next-generation communication infrastructure for meeting the bandwidth and energy efficiency requirements of HGC. One of the recent challenges in the field of ONoCs is to accelerate multicast routing via wavelength division multiplexing (WDM), which dispatches data parallelly across non-interfering wavelengths. The optimization of network throughput and laser power efficiency revolves around two factors: the number of wavelengths and optical signal-to-noise ratio (OSNR). Accordingly, we introduce OSNR-aware wavelength allocation and branching methods for multicast routing (OWBM) tailored to an HGC platform in a customized ONoC. OWBM increases the wavelength resource efficiency by establishing independent routing paths in the partitioned destination nodes such that each routing path is guaranteed to be prevented from overlapping among the partitions. Moreover, the adaptive branching mechanism of OWBM adaptively selects path-based routing and OSNR-aware routing on the fly according to the wavelength allocation cases, further augmenting the throughput and laser power efficiency. Consequently, OWBM outperformed conventional tree- and path-based multicast approaches by elevating the average throughput by 47.39% and curbing the laser power consumption by up to 35.92% in various convolutional neural network benchmarks. Compared with the existing ONoC wavelength allocation techniques, the OWBM demonstrated a maximum of 42.46% enhanced wavelength utilization on a 64-core HGC platform.

**INDEX TERMS** Deep learning kernel, heterogeneous computing platform, multicast routing, optical network-on-chip, wavelength allocation.

The associate editor coordinating the review of this manuscript and approving it for publication was Kathiravan Srinivasan[ID].

## I. INTRODUCTION

With the relentless growth of computer vision technology beyond extended realities (XRs) and the ubiquity of artificial

intelligence applications, the demand for energy-efficient computational systems is increasing. With the advancements in heterogeneous computing (HGC) platforms, general-purpose central processing unit (CPU) architectures have been developed to provide more versatile solutions tailored to address this multifaceted computational landscape [1], [2].

Optical network-on-chip (ONoC) has emerged as a promising innovation that provides a significant communication bandwidth of hundreds of gigabits per wavelength, which cannot be achieved using electrical network-on-chip (ENoC) [3], [4]. Furthermore, by harnessing the capability of wavelength-division multiplexing (WDM), ONoCs can provide scalable bandwidth with high energy efficiency for the upcoming application demands. The presence of diverse processors on an HGC platform invariably leads to asymmetric data access patterns. Custom topology-based ONoCs provide a robust solution to address these challenges. The custom topology envisages a structural modification approach by adapting an on-chip interconnect tailored to massive asymmetric traffic. In particular, as the importance of multicasting has increased in the latest deep neural network (DNN) applications, the capability of ONoCs to handle multicast requests has become essential for transmitting multicast packets efficiently. In other words, the importance of the ONoC multicast routing technique, which takes into account both throughput improvement through ONoC and multiple multicast response capabilities, has emerged.

Multicasting plays a pivotal role in accelerating parallel processing and ensuring synchronization in distributed computing systems. Several real-world applications, from programming to server workloads, such as TPC-H [5] and SPECweb99 [6], require concurrent management of multiple instances. Moreover, other applications, including distributed shared cache systems and various programming models inherently consist of multiple multicasts [7]. Hence, there is a growing demand for a network infrastructure that can effectively provide multicast capabilities across a wide range of applications, particularly in the context of the DNN, which involve aggressive multiple multicast transmissions.
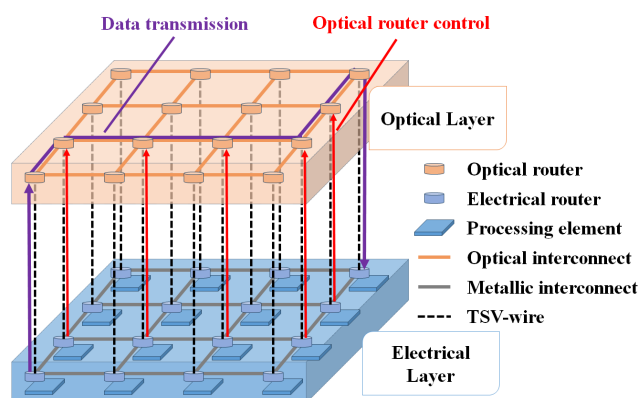


**FIGURE 1.** Example of data transmission in a hybrid ONoC.

Fig. 1 shows an example of data transmission in a hybrid ONoC [8]. The electrical layer includes routers that have a one-to-one correspondence with processing elements (PEs). Data transmission is conducted through an optical layer comprising optical routers controlled by electrical routers. In an electrical router, control packets are transmitted through metallic interconnects that physically connect the electrical routers to the path setup of the optical routers. Because the ONoC adheres to a circuit switching method that maintains the communication path until teardown once set up, it differs from packet switching-based ENoC routings [9], [10], [11]. In the electrical router of an ONoC, because a single packet contains only the address information of the source and destination necessary for controlling the optical elements, the packet length is relatively short compared to that of an ENoC. Accordingly, the deadlock problem caused by several packets occupying the buffers across multiple routers has not been addressed in ONoCs. The main challenges of ONoCs are minimizing contention in multiple multicasts and maximizing the worst-case optical signal-to-noise ratio (OSNR), which determines the laser source power.

- Reducing contention in multiple multicasts: When multiple source nodes request multicasts simultaneously, collisions in multicasting are attributed to the shared use of the same link at identical wavelengths. The conflict problems caused by multiple multicasts can be effectively resolved through dynamic wavelength allocation and adaptive branching [12]. The wavelength allocation technique, which assigns wavelengths for each optical path on the fly under limited wavelengths, efficiently mediates between overlapping paths. Moreover, adaptive routing technology that avoids overlapping routes as much as possible and suggests alternative routes is required.

- Enhancing worst-case (WC) OSNR: The attenuation experienced by the optical signal as it passes through the optical components and noise due to unintended light leakage are combined in the OSNR [13]. The laser source power is determined by the WC OSNR, photodetector sensitivity, and laser wall-plug efficiency ($L_e$), and accounts for a large portion of the ONoC power consumption [14]. Because $L_e$ and photodetector sensitivity are not controllable owing to hardware constraints, the WC OSNR is the main factor determining the laser source power [15]. Because the OSNR is affected by the optical components through which the optical signal passes, the WC OSNR is determined using the routing method under the given topology. Consequently, a routing strategy that includes branching schemes that consider WC OSNR is essential for the energy efficiency of ONoC multicasting.

Addressing the intertwined challenges of contention in multiple multicasts and WC OSNR is essential for achieving efficient and robust ONoC multicasting solutions. There is a pressing need for effective multicast routing solutions that minimize wavelength usage and maximize the OSNR. Accordingly, we introduce OSNR-aware wavelength

allocation and branching methods for multicast routing (OWBM), which is a novel approach tailored to custom topology-based ONoCs that serve various DNN applications causing aggressive multicasts. By separating the destination partitions for each source node, OWBM ensures that the routing paths for each partition are distinctly non-overlapping. The partitioning of destination nodes enhances the utilization of a single wavelength across multiple partitions.

On the other hand, the hybrid multicast routing mechanism of OWBM combines the strengths of both tree-based and path-based multicast routings. By concurrently deploying the Hamiltonian path-based routing table and the OSNR-aware adaptive branching table, OWBM enables the establishment of optical paths with an enhanced OSNR. Moreover, in the event of a blockage by another signal path at a branching point, the path-based routing table provides an alternative course via a labeled path. The hybrid routing method provides an excellent advantage in coping with frequent contentions in multiple multicasts by providing alternatives to blockages, while being oriented toward the communication of the shortest path. In summary, in OWBM, hybrid routing and on-the-fly wavelength allocation methods were proposed for the purpose of improving path diversity and WC OSNR, thereby enhancing throughput and reducing laser source power, while increasing wavelength utilization to augment energy efficiency. The main contributions of OWBM encompass three key areas.

1) For ONOC multicast, OWBM highlights the advantages of each routing method through path- and tree-based hybrid routing.
2) Control packets experience lower waiting time in multicasting processes by circumventing hot spots through on-the-fly wavelength allocation.
3) OWBM employs a two-stage OSNR-aware port filtering to select the port that maximizes the WC OSNR during the multicast path selection process.

The remainder of this paper is organized as follows: In Section II, related work and the background on NoC multicasting are described. A detailed description of the OWBM is presented in Section III. Section IV presents the evaluation of OWBM for various DNN benchmarks. Finally, conclusions are presented in Section V.

## II. RELATED WORK

Challenges with the WC OSNR and multiple multicast collisions are resolved by implementing routing strategies that include wavelength allocation and adaptive routing. The routing strategies that determine how paths are branched and routed are encapsulated within multicast routing methods. Historically, these methods have been broadly divided into two primary categories: tree- and path-based routing. In the following sections, we will provide a comprehensive review of notable studies that explore these two multicast routing schemes within the NoC framework.

### A. TREE-BASED ROUTING IN NoC

In tree-based routing, signals travel along a shared route to reach different destinations, and then they are duplicated and split by an optimal branching router. This type of routing has the advantage of reduced network latency owing to its shortest-path construction. However, a notable limitation is the potential increase in waiting time if other signals hinder replication at a specific branching point. Several techniques have been proposed under this category: Virtual circuit tree multicasting (VCTM) is proposed to form a multicast tree via a virtual circuit table [16]. Although VCTM ensures limited router latency during packet transmission, it poses power efficiency challenges because the virtual circuit table is pervasive across all routers.

An optimized tree (OPT) and a left-XY-right-optimized tree (LXYROPT) were introduced in [17], which ensure a deadlock-free environment by constraining the turn model in mesh-based topologies. Their strength lies in the low power consumption along with low latency achieved through a straightforward algorithm; however, they did not address multiple multicast environments. For example, because LXY-ROPT divides destinations into two groups and combines XY and YX routings, hold-and-wait packets frequently occur when two or more multicasts share the same axis. A switch tree-based algorithm (STBA) was proposed to configure a multicast tree in a reconfigurable NoC [18]. By leveraging both the Kruskal minimal spanning tree and the west-first routing algorithm, STBA enhances power efficiency and reduces latency.

### B. PATH-BASED ROUTING IN NoC

Path-based routing offers the benefits of a straightforward hardware implementation and prevents deadlocks by using Hamiltonian paths for multicasts. The Hamiltonian path is constructed via more routers rather than by establishing the shortest path, thereby increasing the path setup time and OSNR degradation. However, this drawback can be overcome by grouping destination nodes to form multiple paths. Several significant strategies under this mechanism include dual-path (DP), multi-path (MP), and column-path (CP) [19], [20]. DP is a fundamental method in path-based multicast routing, which divides destinations into two groups based on source labels. However, as the network size increases, DP exhibits limitations in terms of scalability. MP enhances the scalability of multicast routing by clustering the destination nodes into four groups. MP has an advantage of a shorter maximum path length than DP because of its subdivided destination groups. Moreover, CP further subdivides destinations using column addresses to achieve heightened parallelism and diminished packet latency.

Building on the foundations established by prior research on path-based multicast routing, a noteworthy contribution is the introduction of a deadlock-free adaptation of a dual-path multicast algorithm tailored for mesh-based ENoCs, as presented in [19]. However, DP has certain limitations in

terms of adaptability to multicast routing. A consequential downside is extended network latency stemming from the prolonged route creation process inherent in the fixed routing approach. To address these limitations, a routing protocol for multicast and unicast (HAMUM) was introduced in [21]. By harnessing the Hamiltonian path, HAMUM facilitates routing for both multicast and unicast traffic, thereby enhancing efficiency. A salient feature of the HAMUM is the network congestion awareness. HAMUM judiciously assesses the congestion state of the input port, ensuring that messages are routed through less crowded pathways to achieve a balanced load distribution.

Tiwari et al. proposed a dynamic partition merging (DPM)-based multicast routing algorithm. DPM was designed to address deadlock situations by dividing the physical network into high- and low-channel subnetworks [22]. The destination node partitioning method of DPM helps prevent mixed turns that could lead to deadlocks with improved throughput. Kang et al. developed a path-based multicast routing for neuromorphic processors (PMRNP), which segments NoC nodes into clusters to leverage Hamiltonian paths for efficient routing [23]. By managing multicast traffic, PMRNP results in a marked decrease in the maximum transmission latency while contributing to a substantial improvement in the average throughput.

Although the latency and energy limitations caused by heavy multicast in NoCs have been successfully addressed, previous studies on multicast routing methods have focused primarily on regular topologies, particularly around the mesh topology. Multicast routing for customized NoCs (MRCN) [24], lays out a multicast routing technique suitable for custom ENoC topologies. MRCN proposes a routing algorithm that adopts Hamiltonian path-based routing, labels the routers of a custom topology, and performs adaptive branching according to the buffer state. Although MRCN shows high performance in ENoCs with guaranteed deadlock prevention, it is difficult to apply it to ONoCs because MRCN does not address the wavelength allocation and OSNR challenges, which constitute a significant proportion of ONoC performance.

Research has been conducted on hybrid routing methods incorporating tree-based routing to mitigate the disadvantages associated with the longest paths in path-based routing. A tree- and path-based hybrid multicast routing algorithm for NoC was proposed to improve network connectivity [25]. This hybrid multicast routing algorithm primarily utilizes path-based routing by descending order of router labels. Tree-based branching is employed when the buffer of the downstream router has sufficient space. Similarly, Wu et al. introduced a hybrid multicasting approach with enhanced methods (HMEM) explicitly designed for mesh-based NoCs [26]. By employing a heuristic path balancing method (HPBM) and a node balancing (NB) strategy, HMEM improved throughput while avoiding deadlocks. However, hybrid routing approaches for ENoC still face the issue of long hop counts inherent to path-based routing, which is

their foundational routing strategy to circumvent deadlock problems.

In summary, while both tree- and path-based routing techniques have advanced the field, a glaring oversight has been the non-consideration of the OSNR, which considerably impacts ONoC power consumption and bandwidth. In ONoC environments with prevalent multicasts, overlapping communication paths, and substantial laser power consumption, a multicast routing approach that considers both wavelength allocation and OSNR is essential. Consequently, we proposed the OWBM that synergizes the strengths of both routing methods and emphasizes the OSNR.

## III. OWBM

In contrast to ENoCs, ONoCs exhibit significantly longer setup and teardown times, particularly in high-traffic scenarios including DNN applications. The primary reason for this is the circuit-switching characteristics of the ONoCs. Uneven traffic patterns on HGC platforms exacerbate the perceptible delays in ONoCs. Recognizing these inherent characteristics of the ONoC, we capitalized on the importance of increasing path diversity, thus increasing the throughput of the ONoC-based HGC platform while alleviating the challenges associated with circuit switching.

On the other hand, the laser source power level should be decreased as much as possible because it determines the overall power budget of the ONoC; however, there is a conflicting requirement to ensure a detectable OSNR at all receiving nodes. Consequently, OSNR is a crucial design consideration for ONoCs. The OSNR was determined by assessing two key factors: the insertion loss (IL), which quantifies the level of signal attenuation, and crosstalk noise, which arises from interfering optical signals as a result of the planar properties of the optical components. The IL and crosstalk noise are caused by the crossing of signals due to the planar structure of the optical elements, imperfections in the optical elements, and scattering of the light signals. Specifically, IL refers to the reduction in signal intensity that occurs as light travels through an optical component or medium, typically caused by absorption or reflection by the material. Conversely, crosstalk noise arises when optical signals overlap or interfere with one another, often because of their close proximity in dense optical networks. Interference due to crosstalk can significantly degrade the signal quality, affecting communication effectiveness and reliability.

Typically, IL is considered a more influential factor in the OSNR. Therefore, a routing strategy that prioritizes the shortest path through the fewest optical components is considered advantageous from the OSNR perspective. Nonetheless, a study published in [27] revealed that the shortest path is not always profitable because of the possibility that crosstalk from neighboring nodes dominates the OSNR. Hence, a routing strategy that considers both the IL and crosstalk can maximize the energy efficiency of ONoCs.
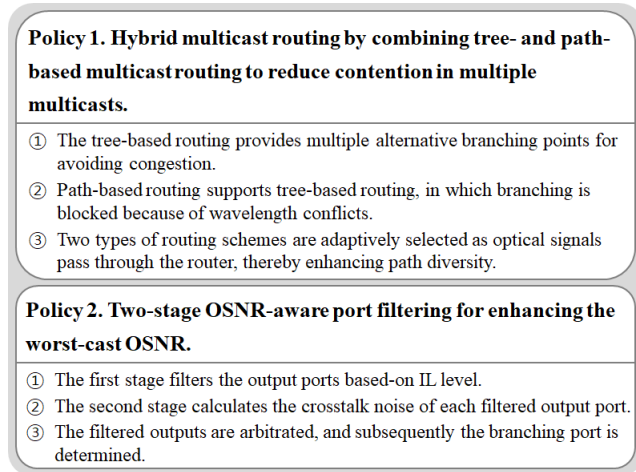
**Policy 1. Hybrid multicast routing by combining tree- and path-based multicast routing to reduce contention in multiple multicasts.**

① The tree-based routing provides multiple alternative branching points for avoiding congestion.
② Path-based routing supports tree-based routing, in which branching is blocked because of wavelength conflicts.
③ Two types of routing schemes are adaptively selected as optical signals pass through the router, thereby enhancing path diversity.

**Policy 2. Two-stage OSNR-aware port filtering for enhancing the worst-cast OSNR.**

① The first stage filters the output ports based-on IL level.
② The second stage calculates the crosstalk noise of each filtered output port.
③ The filtered outputs are arbitrated, and subsequently the branching port is determined.

**FIGURE 2.** OWBM policies to respond to the two major challenges.

The policies of OWBM designed to address these observations are illustrated in Fig. 2. Because the optical layer performs passive operations controlled by the electrical layer, OWBM focuses on designing electrical routers. In this context, OWBM deploys hybrid multicast routing-based adaptive branching as the first policy to resolve multiple multicast collisions. Tree-based routing transmits a signal by selecting one of the shortest paths from the source to the destination node. Path-based routing complements tree-based routing by providing branchable output ports when the branching of tree-based routing is disallowed owing to limited wavelengths. Consequently, path diversity can be enhanced through the adaptive selection of two routing schemes for each router, leading to a reduction in contention in multiple multicast.

In scenarios in which multiple sources attempt routing, waiting time can be significantly reduced if the sources are presented with numerous alternative paths [28], [29], [30]. This approach allows the effective bypassing of paths already set up, thereby minimizing the delays caused by waiting for path availability. In this context, OWBM, which adopts an adaptive selection of routing schemes, becomes particularly advantageous. The strategy of OWBM, which combines path- and tree-based routing schemes and leverages the concept of path diversity. By maximizing the number of potential detour paths and minimizing instances where occupied paths block other potential signal routes, OWBM effectively reduces the overall network latency. This is a critical aspect of OWBM, distinguishing it from traditional routing methods that lack such adaptive capabilities.

The second policy aims to search for a routing path with the maximum OSNR through two-stage OSNR-aware port filtering. The two filtering stages of OWBM pursue the routing path with the maximum OSNR by prioritizing the branching ports with the lowest noise, including the IL and crosstalk. Further detailed explanations regarding the operation of OWBM established in line with these policies are provided in the following subsections.

In summary, OWBM addresses the multicast routing method to improve the throughput and OSNR based on two policies. In addition, the network wavelength utilization increases through on-the-fly wavelength allocation, which determines and maps wavelengths that have not been assigned to adjacent routers. Because the static power consumption for laser and MR tuning accounts for most of the ONoC power consumption, increasing the wavelength utilization contributes to improving throughput and reducing packet delay with the same static power while enabling an energy-efficient multicast response in ONoCs.
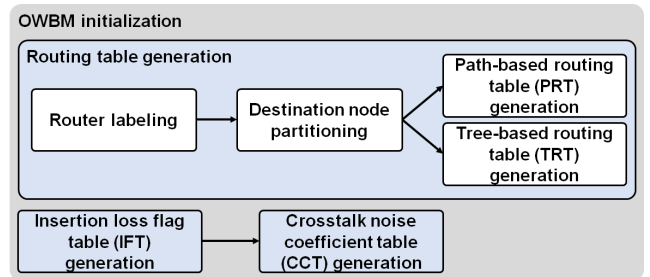


**FIGURE 3.** Procedure of OWBM initialization in design time.

### A. INITIALIZATION

Fig. 3 shows OWBM initialization procedure during design time. To facilitate a hybrid approach that combines path- and tree-based routing schemes, OWBM builds the required routing tables by labeling and partitioning the destination nodes. Next, the auxiliary tables, IL flag table (IFT) and crosstalk coefficient table (CCT), enumerating the IL and crosstalk are created. Routing and auxiliary tables are utilized to perform OSNR-aware routing by detouring congested paths.

#### 1) ROUTING TABLE GENERATION

First, router labeling was performed to facilitate node partitioning while determining the Hamiltonian path. The label of a router represents the order of the routing path from the source to the destination, and serves as an index for grouping the destination nodes.

Router labeling is performed to construct Hamiltonian paths that visit all destination nodes exactly once from the source node. As the number of links forming the network topology increases, path diversity is enriched, thereby increasing the probability of a Hamiltonian path being constructed. Thus, in router labeling, prioritizing the router with the lowest degree maximizes the number of links across the remaining unlabeled routers [24]. Accordingly, router labeling begin from the router with the lowest degree among routers adjacent to the source node.

Algorithm 1 presents the pseudocode for router labeling of OWBM. The router labeling algorithm receives a topology graph *TG* as an input consisting of vertices and links. Target source $i$ is designated as the first label of $lr_i$ denoting the

**Algorithm 1** Router Labeling of OWBM

---

**Input** Topology graph $TG$ {$V$(vertex set), $L$(link set)}
**Output** Labeled router set $LR$ {$lr_s | lr_s$ is labeled routers for source $s$}
1:  **for** $i = 1 to |V|$ **do** // all vertices selected as the target source $i$
2:      $lr_i(0) = i$
3:      $index(0) = 1$
4:      **while** ($lr_i < |V|$) **do**
5:          Insert vertices adjacent to the $v_{index-1}$ into the temporal
               vertex set $V_{tmp}$ ($v_i$ is $i^{th}$ node of NoC, $1 \leq I \leq |V|$)
6:          Sort $V_{tmp}$ in ascending order of router degree
7:          $check\_ham = 0$; // check if a Hamiltonian path formed
8:          **for** $j = 1 to |V_{tmp}|$ **do**
9:              **if** ($v_j \notin lr_i$)
10:                 $lr_i(index) = v_j$
11:                 $index + +$
12:                 $check\_ham = 1$
13:                 **break**
14:             **end if**
15:         **end for**
16:         **if** ($!check\_ham$) **break**
17:     **end while**
18: **end for**

---

labeled router set (Line 2). The router labeling for $i$ is repeated until a Hamiltonian path is formed (Line 4).

The routers adjacent to the most recently labeled router are sorted starting from the lowest degree (Lines 5-6). Because a Hamiltonian path is not guaranteed to exist in a custom topology, the $check\_ham$ variable is used to terminate router labeling when more than one router cannot be labeled (Line 7).

The unlabeled vertex, $v_i$, is selected for labeling according to the order of degrees using the pre-sorted vertex set $V_{tmp}$ (Lines 8-14). If the $check\_ham$ value is zero, this implies that all adjacent routers are already included in the labeled router set $lr_i$. Therefore, it is determined that forming a Hamiltonian path fails; thus, router labeling for $i$ is terminated.
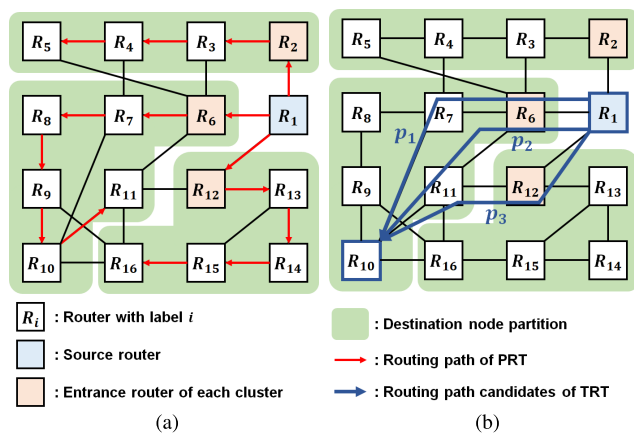


**FIGURE 4.** Example of (a) destination node partitioning and (b) routing path candidates of TRT.

Fig. 4(a) shows when destination node partitioning is performed with source router $R_1$. Given that the router adjacent to the source node serves as the entry router for the respective partition, the total number of partitions is determined by the degree of the source router. Because

the degree of the source router $R_1$ is three, the number of destination partitions for $R_1$ is also three. As the routing paths follow the order of the labels of the routers, the difference between the labels of consecutive entrance routers, depending on the order of the labels, is regarded as the associated partition size. For example, the partition size containing entrance router $R_6$ becomes six because the label difference between $R_{12}$ and $R_6$ is calculated to be six. Similarly, in the case of a partition that includes entrance router $R_2$, the partition size becomes four by considering that the label of the following entrance router is six. Since the labels of the routers already contain all the elements required for partitioning, destination node partitioning can be merged with router labeling, thereby shortening the initialization time.

After destination node partitioning, a path-based routing table (PRT) and tree-based routing table (TRT) were configured with partitioned routers. Such a dual-structured table is devised to augment autonomy in route selection so that OWBM offers alternative paths to mitigate congestion. Because destination partitions already entail routing paths, a PRT configuration is achieved without path search.

On the other hand, tree-based routing mandates data transmission between the source and destination nodes via the shortest path, which is distinct from path-based routing that necessitates a Hamiltonian path encompassing all routers in a sequential manner. In alignment with this criterion, a depth-first search is executed from each source node to construct the shortest path, thereby facilitating the generation of a TRT.

Inter-partition routing is not permitted to maximize wavelength utilization by using the same wavelength for each partition; thus, the routing path search of the TRT is confined to each individual partition. In addition, given that the packets are transmitted based on an adaptive selection between the TRT and PRT in each router, routing paths that violate the ascending order of labels are excluded from the path search. Hence, the routing path of the TRT complies with the ascending order of the labels within the partition, thereby ensuring that path-based routing remains executable after tree-based routing.

Fig. 4(b) shows an example of the routing path candidates of the TRT from $R_1$ to $R_{10}$. Routing paths $p_1$, $p_2$, and $p_3$ are the shortest paths passing through the fewest routers connecting $R_1$ and $R_{10}$. Only $p_1$, is accepted in the TRT because it belongs to the same partition while maintaining the ascending order of the labels.

By contrast, $p_2$ and $p_3$ violate the TRT rule because $p_2$ does not maintain the ascending order of the labels when moving from $R_{11}$ to $R_{10}$, and $p_3$ cannot be routed within the same partition. The route computation stage in the NoC pipeline can be simplified by selecting either the PRT or TRT.

### 2) IL FLAG TABLE GENERATION
When an optical signal propagates through a waveguide, it encounters optical components, such as crossings, bends, and MRs, resulting in an IL that attenuates the signal intensity. The corresponding attenuation ratios for a waveguide,

crossing, bending, OFF-state MR, and ON-state MR are 0.274 dB/cm, 0.005 dB/90°, 0.04, 0.005, and 0.5 dB, respectively [31], [32], [33], [34]. It has been observed that an optical signal passing through an ON-state MR causes the most significant attenuation, necessitating careful consideration of routing design. Therefore, to minimize the IL due to the ON-state MR, optical routers are generally designed to pass through at most one ON-state MR, depending on the connection between the input and output ports. For example, Crux [37], an optical router specialized for XY routing in mesh-based ONoCs, has 12 MRs corresponding to the 12 port connection cases, excluding YX turns.

Calculating the OSNR on the fly imposes a considerable computational burden owing to both the storage requirements for the floating-point IL coefficients and the computational load involved in the IL calculations. To resolve these challenges, OWBM employs an IL flag that serves as a simplified indicator, denoting whether the signal traverses an ON-state MR. If the optical signal passes through the ON-state MR, the corresponding IL flag is configured with a value of one. Otherwise, the IL flag is set to zero. The computational cost of the OSRN can be significantly reduced by using the composite IL flag values as a rough approximation for the IL. Moreover, the flag-based approach drastically reduces the storage space required to retain the IL coefficients as well as the combinational logic for IL calculations. For example, compared with the 16-bit floating-point IL coefficient, storing the IL flag can reduce the size of the IFT by 16 times.

### 3) CROSSTALK COEFFICIENT TABLE GENERATION

To calculate the crosstalk noise based on port mapping, OWBM utilizes the linear optical device model (LODM) and the computation method of optical signal propagation (CMOP) proposed by Kim [35]. OWBM employs LODM and CMOP to expedite the derivation of crosstalk noise coefficients, which is significantly faster than the traditional finite-difference time-domain (FDTD) methods used in optical simulators.

LODM streamlines the propagation of optical signals through components into a linear model, facilitating fast and precise signal strength calculations. CMOP applies an interpolation technique to reduce the extensive computations involved in MR resonance. The interpolation of the time difference between MR signal calculations per round trip adjusts for round-trip time discrepancies caused by computation rate synchronization.

To generate the CCT for OWBM, the $5 \times 5$ non-blocking optical Cygnus router [36] was modeled using the LODM approach. While LODM does not distinguish between noise and signals, signals exiting any port other than the intended output can be considered as noise. For the generation of CCT, we calculated the ratio of the output optical signal power to the input power in LODM using CMOP, and determined the corresponding ratio as the crosstalk noise coefficient.

When calculating crosstalk noise, retaining the coefficients for each potential port connection leads to significant area overhead and massive computing resources. Thus, CCT maintains the coefficients associated solely with the single port connection scenarios. The cumulative crosstalk was determined by referencing the active ports in the switch matrix and summing the associated interference values. Noise estimation involves the multiplication of the coefficients by the amplitude of the undesired input signals.

When the intensity of optical signal coming from the $i^{th}$ input port is $I_i$ ($1 \leq i \leq m$, $m$ is maximum port number), and the output port to which the $i^{th}$ input port is routed is $out(i)$, the sum of the crosstalk noise at the $j^{th}$ output port, $Xtalk_j$, for $CN_{i,out(i),j}$, which is the crosstalk noise coefficient from $i^{th}$ input port to $j^{th}$ output port, is calculated as follows:

$$Xtalk_j = \sum_{k=1}^{m} I_i \cdot CN_{i,out(i),j}, \ i,j \in \{0, 1, \cdots, m\}. \quad (1)$$

Given the inherent challenges in the optical domain, measuring the optical signal intensity directly at individual routers is impractical. A standardized intensity of 1 was adopted for all input signals to calculate the crosstalk. Consequently, a simplified $Xtalk_j$ is derived using (2).

$$Xtalk_j = \sum_{k=1}^{m} CN_{i,out(i),j}, \ i,j \in \{0, 1, \cdots, m\}. \quad (2)$$

Therefore, the runtime of adaptive branching (ADB) for accessing the crosstalk of each output port was significantly reduced because the multiplication of floating points was not required in the crosstalk calculation. As a result, fast derivation of crosstalk noise in an ADB is possible with only a small-sized CCT and logic. Distinct from OWBM evaluation, which requires precise crosstalk calculations, the simplification of input signal intensity in calculating crosstalk noise is used in the OSNR-aware port filtering phase of OWBM.
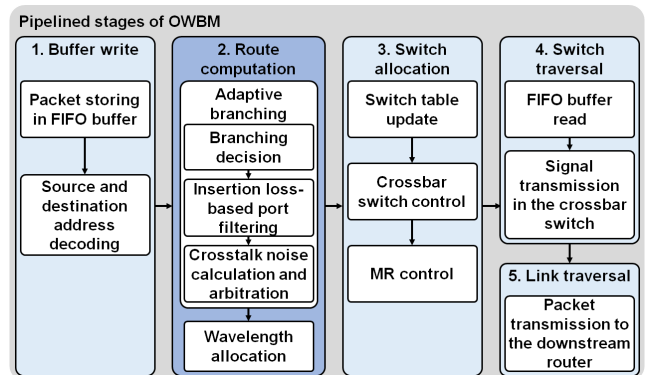


**FIGURE 5.** Overall procedure of five-stage pipelines of OWBM.

### B. ADB

The overall procedure for the pipeline operation of OWBM is illustrated in Fig. 5. Among the pipeline stages, OWBM

primarily focuses on the route computation stage, which is directly involved in multicast routing consisting of ADB and on-the-fly wavelength allocation (OWA). Based on the lookup tables created during the initialization stage, the OSNR and wavelength utilization were optimized by implementing the ADB and OWA procedures governed by the tables previously generated during the initialization phase. ADB and OWA are integral components of the route computation process within the pipeline stage of the router.

OWBM minimizes the path setup time and enhances the OSNR by prioritizing tree-based routing via TRT, where the optical signal passes through routers with minimum hop counts. For example, in a 4 × 4 mesh-based ONoC, the hop count of the longest path varies by up to nine depending on the routing algorithm. Thus, assuming a five-stage pipeline, OWBM can reduce the path setup time by up to 45 cycles compared to path-based routing. In the ADB process, the PRT provides alternative paths when the TRT is disallowed from branching owing to the wavelength constraints occupied by other multicast requests. Because the limitation in terms of packet delay of tree-based routing, which allows branching only at fixed nodes, is resolved by the adaptive selection of the PRT, OWBM can effectively enhance the throughput and energy efficiency of ONoC multicasts with on-the-fly wavelength allocation.

The ADB is articulated in three steps: 1) branching decision (BD), 2) IL-based port filtering (ILF), and 3) crosstalk noise calculation and arbitration (CCA). In the BD step, the port candidates in the TRT are searched to determine whether there are any branchable ports for each destination node of the multicast packet for the current source node. In the ILF step, the IL flags in the TRT are searched to filter the branchable output ports that undergo a relatively small IL. CCA is performed if multiple ports had the same IL flag. The crosstalk noise inserted from the undesired input ports into the filtered output port according to each port-switching state is calculated by referring to the CCT. Subsequently, CCA determines the branching port with the lowest crosstalk noise.

ADB is conducted using the four tables TRT, PRT, IFT, and CCT generated during the initialization phase of OWBM. Fig. 6 depicts the general ADB process by OWBM when the destination address is 7 among the multicast destination addresses. In the BD step, when the source and destination addresses and port number are 0, 7, and 0, respectively, the TRT value is determined as 11000, resulting in the initiation of the IFT step.

In the ILF stage, all the IL flags of output ports 11000 are identified as 1. Thus, the OSNR of the two output ports cannot be differentiated based on the IL level. The output ports filtered by the ILF step initiate the CCA phase. In the CCA stage, after calculating the crosstalk noise directed to the two filtered ports, the output port with the lowest crosstalk noise receives the highest priority during the CCA stage. Accordingly, OWBM ensures the branching of the signal along a path characterized by a superior OSNR.



**FIGURE 6.** Process of ADB when one of the multicast destinations is 7.

### 1) BD AND ILF

The first step of ADB determines whether a branchable port exists in the TRT. Because the bit position of the TRT value represents each branchable output port, if at least one bit in the TRT is 1, the BD step activates the ILF. Otherwise, a tree-based routing path does not exist in the current router. Thus, path-based routing is conducted by referring to the PRT. In conclusion, the DB step is compressed by searching the port candidates of the TRT and simply passing them through the OR gate, thereby enhancing the feasibility of the OWBM implementation.

**Insertion loss-based port filtering**



**FIGURE 7.** Flowchart of ILF.

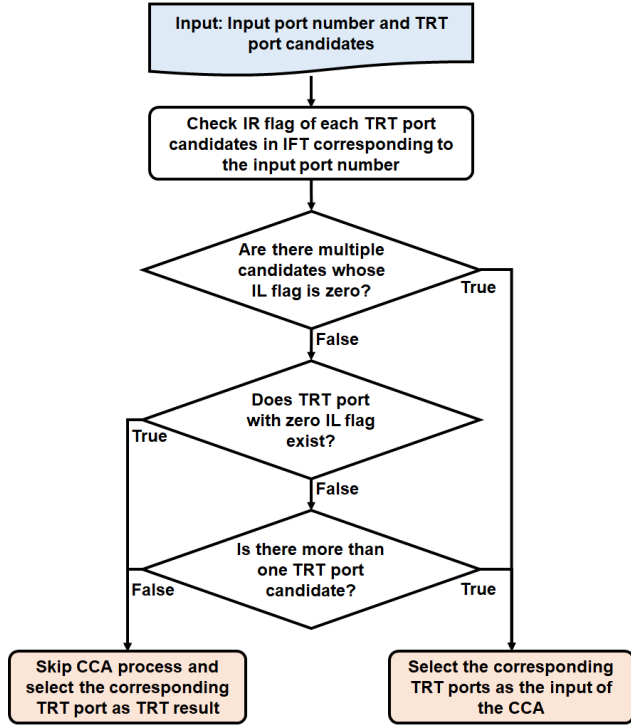The flowchart of the ILF is illustrated in Fig. 7. The objective of the ILF is to identify ports within the TRT that lead to minimal IL by utilizing IL flags, thereby contributing to the improvement of the OSNR. Beginning with a scan of IL flags associated with port candidates of the TRT for OSNR-sensitive filtering, multiple branchable ports with an IL flag value of zero receive prioritization as branching ports. A zero-valued IL flag indicates that an optical signal bypasses the ON-state MR, leading to reduced IL. If only one port remains after filtering, wavelength allocation is directly conducted, bypassing the CCA. Otherwise, if multiple filtered ports have the same IL flag, ranking the TRT port candidates using the IL flags becomes infeasible; thus, the filtered output ports proceed to the CCA stage.

### 2) CCA

In the CCA stage, OWBM calculates the crosstalk inserted from the undesired input ports to the output ports filtered by the ILF and preferentially selects one output port influenced by the lowest noise as the CCA result. The crosstalk noise coefficient in the CCT is determined using the filtered output port and switch table as reference indices. Given that the CCT contains only the crosstalk noise coefficient corresponding to each single port-mapping scenario, the CCT area can be significantly decreased. For example, storing only a single port-mapping scenario reduces the memory space in the CCT by approximately five times in a 5-port router compared to storing all the port-mapping scenarios. CCA prioritizes the

output port with the lowest crosstalk noise among the output ports filtered by the ILF. Finally, arbitration is performed based on the filtered output port and corresponding priorities. After the CCA phase, wavelength allocation is performed according to the branching port. If the wavelength allocation fails owing to other multicast requests, the CCA receives a feedback signal from the wavelength allocation logic to select the remaining output ports with subsequent priority.



**FIGURE 8.** Example of crosstalk noise calculation and arbitration when the filtered output port is 00011 and two output port 1 and 2 is occupied in 5-port router.

Fig. 8 depicts the CCA in a 5-port router when the filtered output port derived by the ILF is 00011 and two output ports, one and two, are occupied. As the two output ports are simultaneously occupied by other multicasting requests, the crosstalk noise coefficients corresponding to the output ports are derived from the CCT using the switch

table. Accordingly, $Xtalk_{00001}$, the crosstalk noise for the filtered output port 00001, is obtained by floating-point summation of $CN_{1,2,1}$ and $CN_{5,3,1}$. Likewise, $Xtalk_{00010}$ is obtained as the sum of $CN_{1,2,2}$ and $CN_{5,3,2}$. By comparing of $Xtalk_{00001}$ and $Xtalk_{00010}$, the priority of the filtered output port is determined and delivered to a 5-bit arbiter. The 5-bit arbiter then transmits the CCA results to the wavelength allocator using the priority of the filtered output.

If arbitration fails, tree-based routing through the TRT is not permitted in the given wavelength allocation states; thus, path-based routing is performed by referring to the PRT. If wavelength allocation is disallowed by the pre-assigned wavelengths for other multicast requests, the priority shift signal is fed back to the arbiter of the CCA to attempt to connect with another branching port.

## C. OWA

In the early stages of multicast packet transmission, the source node allocates wavelengths equal to the number of destinations for each partition. OWBM reduces signal contention by assigning different wavelengths using a round-robin arbiter, while also improving wavelength utilization by consistently allocating the same wavelength for each destination partition. By leveraging the capability of destination node partitioning, routing paths can be prevented from overlapping among destination partitions, which allows for identical wavelengths to be designated recurrently across multiple partitions. During signal transmission from the entrance router to the destination router, the OWA logic is utilized in the route computation stage to properly distribute the wavelength bundle assigned from the source node at each branching point.



**FIGURE 9.** OWA logic.

The OWA logic is depicted in Fig. 9 when the number of available wavelengths is $W$, and the number of router ports is $P$. The OWA logic is designed to map the input port wavelengths to the output ports based on the CCA result for adaptive branching. Therefore, the wavelength allocation states of the downstream router for each port are transmitted to the adjacent upstream router. OWA logic preserves only the wavelength allocation state corresponding to the CCA result using multiplexers that use the CCA results as control signals. Accordingly, the CCA results, which were composed of binary positional code, were encoded to binary digit code to control the multiplexers.

In the wavelength allocation state, the bit position corresponding to the allocated wavelength is denoted by one. Consequently, the wavelength allocation state is toggled using a $W$-bit inverter. Thus, the bit position represents the assignable wavelength. As the toggled wavelength allocation state and wavelengths of the input port pass through the AND gate, assignable wavelengths are selected from among the input port wavelengths. The selected wavelengths are passed through a round-robin arbiter. Then, an arbitrary wavelength is registered in the wavelength port mapping table. If the OWA fails because the wavelength arbitration result overlaps with the wavelengths pre-assigned by other multicasts, the wavelength port mapping table requests an alternative wavelength from the round-robin arbiter. If there are no more assignable wavelengths for the CCA results, a priority shift signal is transmitted to the arbiter in the CCA to request other branchable ports. These two types of feedback loops improve the capability of OWBM to cope with congestion by providing various wavelength alternatives along with detour routes.

For example, suppose that the $W$ and CCA results are four and 00001, respectively, in the 5-port router. In this case, the wavelength allocation state of the downstream router connected to the first port (00001) is loaded as the output of the multiplexer. If the wavelength allocation state of the corresponding downstream router is 1100 and the wavelength of the input port is 1111, the third and fourth wavelengths are blocked by passing the AND gate; thus, the arbiter receives 0011 as input. Consequently, the first and second wavelengths are determined as the available wavelengths to be allocated. If the priority of the first wavelength is higher than that of the second wavelength, OWBM attempts to register the first wavelength in the wavelength port mapping table. If the first wavelength is already assigned to other multicast requests, then the second wavelength with the following priority attempts to be recorded. If the second wavelength is available for allocation, the wavelength port mapping table is updated with the termination of the OWA.

Fig. 10 shows the entire route computation flow from ADB to OWA. The BD stage involves searching for the TRT and transmitting the resulting output ports to the ILF stage. If the TRT search cannot find any output port, tree-based routing is considered as infeasible, and path-based routing is performed instead. If only one TRT result exists, then the routing process
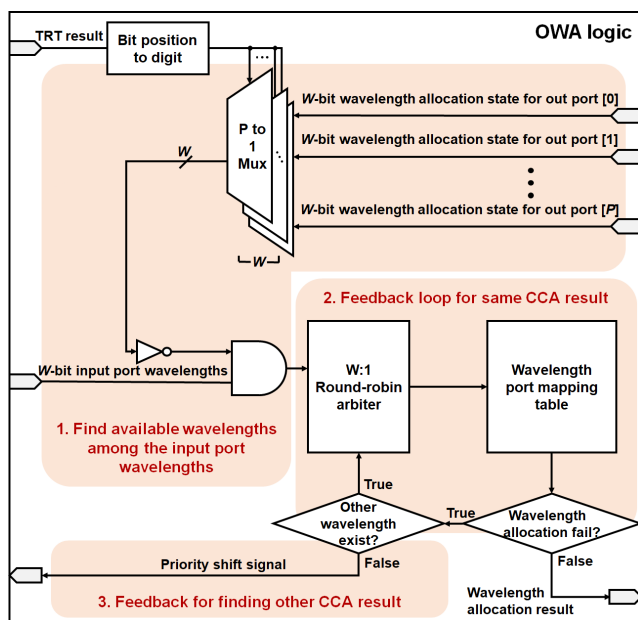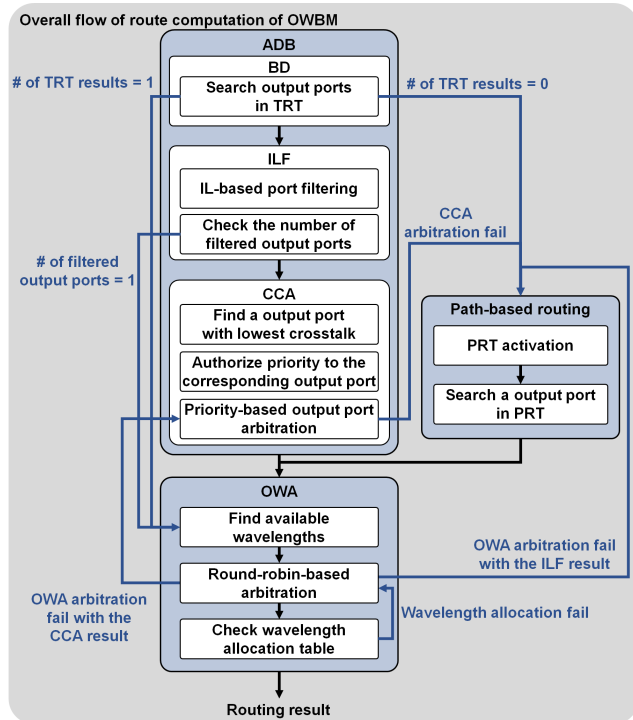
**FIGURE 10.** Overall flow of route computation of OWBM.

migrates directly to the OWA stage, bypassing the ILF and CCA stages. During the ILF stage, the TRT results are filtered based on the IL. If only one filtered result is obtained, the CCA stage is omitted. Otherwise, the process proceeds to the CCA stage. In the CCA stage, the crosstalk noise is calculated for each filtered output port, and the port with the lowest crosstalk is prioritized.

In the OWA stage, the process involves determining and arbitrating the available wavelengths. If a wavelength arbitration fails, there are two possible courses of action. If the failure is due to the CCA results, lower-priority ports from the CCA arbitration stage can be chosen as alternatives; thus, OWA require to repeat the CCA process. Alternatively, if wavelength arbitration fails based on the ILF results and there are no other TRT alternatives, path-based routing is executed immediately. The path-based routing stage activates the PRT and searches for the corresponding LUTs. Similar to the TRT results, the PRT results also underwent OWA, resulting in the final routing result.

## IV. EVALUATION

To assess OWBM, which proposes hybrid multicast routing with an on-the-fly wavelength allocation method, it was compared with conventional tree- and path-based routing schemes in terms of throughput, energy consumption, and wavelength utilization. ADB and OWA enhance throughput by addressing contention from multiple multicasts and optimizing energy efficiency. Therefore, throughput and energy consumption metrics were adopted as key indices to assess

the performance of the hybrid routing and OSNR-aware two-stage port filtering in OWBM. LXYROPT [17], routing and wavelength assignment for distribution-based multiple multicasts (RWADMM) [38], routing and wavelength assignment algorithms using multiple multicast relationship graphs (RWMMRG) [39], MP [19], and MRCN [24] were adopted as comparison groups. The comparison groups encompass diverse tree-based and path-based routing techniques, facilitating a comprehensive assessment of OWBM performance. LXYROPT, MP, and MRCN, which focus on multicasting for ENoCs, addressed routing algorithms irrelevant to the wavelength allocation. Therefore, the comparison groups LXYROPT, MP, and MRCN were applied to the ONoC by assuming that the optimal wavelengths were allocated.

CNN models, with their frequent matrix operations, high data reusability, and repetitiveness, exhibit the highest intensity of multicast among AI applications, ensuring an appropriate evaluation of the multicast algorithms of the OWBM and comparison groups. The visual geometry group (VGG) at Oxford University developed two types of CNNs: VGG-16 with 16 layers and VGG-19 with 19 layers [40]. VGG-16, VGG-19, and ResNet [41], which are widely used in computer vision for tasks such as image classification, object detection, video recognition, and medical image analysis, have more convolution layers than other CNNs, resulting in a higher multicast frequency. Hence, VGG networks and ResNet were adopted as application benchmarks for evaluating the multicast routing algorithms. ResNet uses skip connections in the residual blocks to effectively counteract vanishing gradients.

In CNNs, the selection of data to be stored in the PE buffer for data reusability leads to different data flow strategies: weight stationary (WS), output stationary (OS), and row stationary (RS) [42]. In the RS strategy, every PE is equipped with three types of buffer that can store input feature maps (ifmaps), filters, and partial sums. The stored ifmaps and filters in PEs are sequentially used in convolution operations, which is a characteristic feature of RS that results in high data reusability. In particular, RS exhibited a tendency to increase the multicast intensity as the network size increased, with a proportional increase in PE buffers storing identical multiple ifmaps. Thus, the RS data flow was adopted to evaluate the multicast routing schemes, including the OWBM.

The gem5-gpu [43], which integrates gem5 with GPGPU-Sim, was used to analyze the HGC platform, which combines CPUs and GPUs in terms of throughput, energy consumption, and wavelength utilization. The gem5 [44] is a configurable C++-based simulator that simulates the operation of a full system consisting of a CPU, memory, and an I/O device. GPGPU-Sim [45] was specifically tailored to model the parallel execution behavior and performance of general-purpose computations running on GPUs. The extended WC OSNR searching algorithm (EWOSA) [15], a C++-based OSNR estimator, was combined with gem5-gpu to evaluate the WC OSNR owing to massive multicast. The adaptive scan range of the EWOSA is a key feature that efficiently

filters out signals that significantly influence the IL and crosstalk in the current optical signal path. The EWOSA streamlines identify the worst-case OSNR and reduce the required computational time. Accordingly, the flexibility and scalability of EWOSA are particularly suitable for multicast-supporting ONoCs using WDM, thus assisting in evaluating the WC OSNR performance of multicast routing techniques, including OWBM. Consequently, the EWOSA provides a reliable foundation for the performance analysis of the OWBM in the ONoC. This integration was facilitated by the common use of C++ in both gem5 and EWOSA, which simplified the addition of the WC OSNR analysis code of EWOSA to the C++ headers of the RUBY network.

However, integrating Verilog code for area and power analysis poses a challenge because of its hardware design language nature, which differs significantly from the programming languages in our simulation environment. Consequently, we implemented the electrical router separately in the register-transfer level (RTL) code, underscoring the hardware implementation complexity involved in the logic area and power consumption.

**TABLE 1.** Specifications of optical devices.

| Parameter | Value | Reference |
|---|---|---|
| MR heating power | $5\,\mu\text{W}$ | [14] |
| Photodetector sensitivity ($P_s$) | $-20\,\text{dBm}$ | [46] |
| Power margin of laser | $13\,\text{dB}$ | [47] |
| Laser wall plug efficiency ($L_e$) | 0.25 | [47] |
| Transceiver bandwidth per wavelength | $10\,\text{Gb/s}$ | [4] |
| Maximum number of wavelengths | 4 | |

The specifications of the optical devices for WDM are listed in Table 1. Considering the wavelength utilization and frequency of multiple multicast requests, the maximum number of wavelengths was determined to be four. The lowest wavelength is assumed to be 1550 nm, the free spectral range (FSR) is 30 nm, and the Q factor is 9000 in a WDM-based ONoC using four wavelengths. Thus, the interference intensity between neighboring MRs can be obtained as $5.30 \times 10^{-4}$. Owing to the absence of a library of optical devices using consistent technology, the optical layer was simulated by incorporating studies on each optical device. The hardware constraints affecting the laser source power decisions are the photodetector sensitivity ($P_s$), power margin, and $L_e$. To ensure reliability, the optical signal intensity must be higher than the photodetector sensitivity. If we consider the total worst-case insertion loss $IL_{total}$ and crosstalk noise $CN_{total}$, the signal intensity $S_i$ emitted from the laser must exceed the sensitivity of the photodetector after accounting for $IL_{total}$ and $CN_{total}$ as follows:

$$S_i \cdot (1 - IL_{total}) - CN_{total} \geq P_s. \tag{3}$$

$CN_{total}$ indicates the combined crosstalk noise induced throughout the entire optical signal propagation path, including the IL of crosstalk. In this scenario, $S_i$ can be derived from the product of laser power consumption and $L_e$. Because IL and crosstalk are the only variables that can be optimized through the routing algorithm, reducing the $IL_{total}$ and $CN_{total}$ is the primary focus of our proposed two-stage OSNR-aware port filtering technique in OWBM.

In addition, the MR heater undergoes heat-based tuning such that it can resonate at an appropriate wavelength. Therefore, the MR heating power is calculated as a static value proportional to the total number of MRs in the network.

For the MR modulator, attenuation $\alpha$ is 0.972, coupling ratio $r$ is 0.954, and the length is 10.5 $\mu$m in the structure used in [48], whose insertion loss is $-2$ dB in the ON-state and $-12$ dB in the OFF-state. For the MR switch, the attenuation of each phase shifter is 0.997, the coupling ratio is 0.984, and the length is 10.0 $\mu$m in the structure used in [34] and [49], whose insertion losses were $-0.04$ dB at the through port and $-0.5$ dB at the drop port.

An optical transceiver with a bandwidth of 10 Gb/s per wavelength is used as the optical network interface [4]. Except for OWBM, the remaining comparison groups used an electrical router with 5-stage pipelines consisting of buffer write, route computation, switch allocation, switch traversal, and link traversal, with an operating clock frequency of 1 GHz. Considering the computation time of the ADB and OWA with clock frequency, route computation was separated into two stages so that the six-pipelined electrical router was assumed for OWBM.

**TABLE 2.** Loss and noise coefficients of optical components.

| Loss Parameter | Value | Reference | Noise Parameter | Value | Reference |
|---|---|---|---|---|---|
| Waveguide | $-0.274$ dB/cm | [31] | Optical Terminator | $-50$ dB | [50] |
| Waveguide Bending | $-0.005$ dB/90° | [32] | Waveguide Crossing Reflection | $\approx 0$ | [33] |
| Waveguide Crossing | $-0.04$ dB | [33] | Waveguide Crossing | $-40$ dB | [33] |
| Crossing MR Through | $-0.04$ dB | [34] | | | |
| Crossing MR Drop | $-0.5$ dB | [34] | | | |
| Parallel MR Through | $-0.005$ dB | [34] | Parallel MR Through | $-20$ dB | [34] |
| Parallel MR Drop | $-0.5$ dB | [34] | Parallel MR Drop | $-25$ dB | [34] |

The IL and crosstalk noise coefficients of the optical components are listed in Table 2. To generate the CCT,

we established an optical signal transmission environment among the optical components in a single router using the Verilog-AMS. The crosstalk noise coefficients associated with the single port mapping cases were configured in the lookup table for the CCA. Because the IFT was created by considering the ON-state MR without referring to the IL coefficients, the IL coefficients were used only for the estimation of the WC OSNR.

**TABLE 3.** Processing unit and memory specifications.

| Parameter | Configuration |
|---|---|
| GPU / Shader clock frequency | 0.7 GHz / 1.4 GHz |
| SIMT width | 8 |
| GPU Private L1 I / D cache size | 64 / 64 KB |
| CPU clock frequency | 2.5 GHz |
| CPU Private L1 I / D cache size | 64 / 64 KB |
| Shared L2 cache size | 1 MB per MC |
| DRAM | 2 GB (DDR3-1600 MHz) |

**TABLE 4.** HGC platform architecture configurations.

| Parameter | HGC platform | | |
|---|---|---|---|
| | 16 cores | 32 cores | 64 cores |
| Number of CPU tile | 1 | 2 | 4 |
| Number of GPU tile | 14 | 28 | 56 |
| Number of MC tile | 1 | 2 | 4 |

Table 3 lists the processing unit and memory specifications of multicore CPUs, GPUs, and memory controllers (MCs). Each GPU tile comprised four unified shaders, and it was assumed that the last-level cache of the MC functioned as an L2 cache shared by all CPU and GPU tiles. The configuration of the HGC platform is listed in Table 4. To evaluate the OWBM in a high multicast intensity environment, the ratio of GPU tiles performing CNN operations was assumed to be relatively higher than that of other tiles. Because all multicast routing techniques, except OWBM and MRCN, are specialized for the mesh topology, an HGC platform based on the mesh topology was assumed for a fair comparison. Accordingly, Cygnus [36], a non-blocking $5 \times 5$ switch for the mesh topology, was used as the optical router in the HGC platform simulations.

Additional experiments were conducted on the torus topology to demonstrate the topological versatility of the OWBM beyond the mesh structure. Among the comparison groups of OWBM, the routing algorithms extendable to the torus include RWADMM and MP. RWADMM defines four theorems based on the x- and y-axis relationships between the source and destination using routing paths such as XY, YX, XYX, and YXY. The theorems of the RWADMM remain applicable to torus expansion. Similarly, the destination node

partitioning strategy of the MP based on the source router coordinates exhibits strengths in the torus topology, where the center and corner nodes are indistinguishable. The MP in the torus topology maintained cluster sizes at a quarter of the network size, offering shorter and longer paths than the mesh.

### A. INITIALIZATION TIME ANALYSIS

The initialization of the OWBM refers specifically to the pre-configuration phase of the design stage. Thus, the initialization stage is distinct from the operational initialization of the routers within the network. This process involves creating routing tables beforehand using a C++-based dedicated program. The initialization simulation time was evaluated on a robust server setup, which featured a 22-core Intel CPU with a 3.6 GHz clock frequency and 512 GB of DRAM.

**TABLE 5.** Simulation time for generating required tables in initialization.

| LUT | Simulation time (sec) | | |
|---|---|---|---|
| | 16 cores | 32 cores | 64 cores |
| PRT | 280.31 | 1029.74 | 3235.01 |
| TRT | 137.69 | 396.97 | 1131.54 |
| IFT | 10.94 | | |
| CCT | 35.74 | | |

Table 5 lists the time required to generate the four tables during the initialization phase. The path labeling and destination node partitioning necessary for PRT and TRT generation are encompassed within the PRT generation time, thus leading to a longer simulation time for the PRT than for the TRT. TRT generation, based on a depth-first search with O(n) complexity in a mesh topology where each node has a degree of five, demonstrated a tendency of tripling the simulation time with each doubling of the core count. This significant increase in the O(n) complexity is attributed to the more complex and varied task scenarios encountered when searching for the shortest path as the core count increases.

The generation times for the IFT and CCT are consistent across all network sizes and are solely dependent on the type of optical router used. The IFT configured based on the MR state mapping without separate coefficient calculations required the least time for generation. In addition, the generation of the CCT, facilitated by a linear optical device model enabling rapid crosstalk coefficient calculations for a $5 \times 5$ optical router, requires less than a minute.

In summary, the OWBM facilitated the generation of all the necessary tables for the design of electrical routers in a 64-core HGC platform within just a few dozen minutes.

### B. ONoC PERFORMANCE ESTIMATIONS REGARDING THROUGHPUT AND ENERGY EFFICIENCY

The multicast routing schemes LXYROPT, RWADMM, RWMMRG, MP, MRCN, and OWBM were evaluated in

terms of throughput, WC OSNR, and wavelength utilization. Wavelength utilization refers to the proportion of total bandwidth occupied by the wavelengths employed across all links, normalized over the simulation time. This metric quantifies the efficiency of the available wavelengths for transmitting data within the network. The WC OSNR was calculated as a factor that determines the laser source power. The laser source power was integrated with the MR heating power and throughput to derive the energy per bit of the NoC. All multicast routing methods were simulated using gem5-gpu by recording data movement as a task graph while repeatedly running the given CNN benchmarks on a mesh-based HGC platform.

As shown in Table 6, the throughput, WC OSNR, and wavelength utilization in the 16-, 32-, and 64-core HGC platforms with the VGG-16, VGG-19, and ResNet benchmarks were determined using the gem5-gpu simulator. As the number and size of convolutional layers increased from VGG-16 to ResNet, the multicast intensity increased. In the mesh-based VGG-16 benchmark of the 16-core HGC platform, the OWBM improved the throughput by 43.36, 16.05, 11.17, 33.01, and 5.88% over LXYROPT, RWADMM, RWMMRG, MP, and MRCN, respectively. As the multicast intensity increased, the throughput improvement of the OWBM increased. In the ResNet benchmark of the mesh-based 16-core HGC platform, the OWBM throughput increased by 49.13, 14.83, 13.16, 19.44, and 10.98% compared to those of LXYROPT, RWADMM, RWMMRG, MP, and MRCN, respectively. In the process of OWBM checking the wavelength allocation state of adjacent routers and assigning connectable wavelengths on the fly, the throughput was improved by increasing the wavelength utilization as the multicast intensity deepened. Meanwhile, coupled with the balance of computing, memory, and communication factors, ONoC throughput tendencies vary nonlinearly with network size, thus illustrating some atypical observations. For instance, in the VGG-16 benchmark with LXYROPT, the 32-core HGC platform surpassed the 64-core HGC platform in terms of throughput.

An increase in the number of cores in the mesh-based HGC platform caused an increase in the number of CPU and GPU tiles, thereby expanding the range and strength of multiple multicasts. Accordingly, the throughput enhancement of OWBM with high contention capability tended to intensify as the network size increased. In the mesh-based 16-core HGC platform, the throughput of the OWBM improved by an average of 46.63, 16.45, 14.86, 22.39, and 5.81% over those of LXYROPT, RWADMM, RWMMRG, MP, and MRCN, respectively. In the mesh-based 64-core HGC platform with a larger network size, the throughput of OWBM improved by an average of 55.00, 28.21, 11.99, 20.26, and 2.75% over those of LXYROPT, RWADMM, RWMMRG, MP, and MRCN, respectively. In summary, OWBM showed the highest capability in response to multiple multicasts, with the highest scalability of the network size compared to the other multicasting techniques. Moreover, owing to the circuit

switching of the ONoC, where the bandwidth of the optical transceiver has a dominant effect on the throughput rather than the path setup time, the overhead due to the increase in the pipeline stage of OWBM was negligible. The throughput analysis underscores the efficacy of OWBM, balancing the intricate trade-offs between the OSNR, delay, and loss ratio to optimize network performance.

For the torus-based HGC platform with the VGG-16 benchmark, the throughput of the OWBM improved by an average of 27.55, 38.97, and 7.45% compared to those of RWADMM, MP, and MRCN, respectively. The throughput enhancement was more pronounced in high multicast intensity ResNet applications.

LXYROPT, RWADMM, and RWMMRG, which use tree-based routing via the shortest paths, have a lower WC OSNR than MP and MRCN because of their relatively small IL. In the VGG-16 benchmark of the mesh-based 16-core HGC platform, the WC OSNR of LXYROPT, RWADMM, RWMMRG was 4.04, 1.51, and 0.93 dB higher than that of MP, respectively. The OWBM, which integrates the advantages of tree- and path-based routing, generally showed a higher WC OSNR than the tree-based routings RWADMM and RWMMRG. In the mesh-based 16-, 32-, and 64-core HGC platforms, OWBM had WC OSNRs that were 1.07, 0.89, and 0.28 dB higher on average than RWMMRG, respectively. Overall, in the mesh-based 32-core HGC platform, OWBM showed an average of 0.13, 0.89, 1.89, and 1.33 dB higher WC OSNR than RWADMM, RWMMRG, MP, and MRCN, respectively. The OSNR-aware port filtering of the OWBM resulted in a higher WC OSNR compared to the tree-based routings RWADMM and RWMMRG, which transmit packets only over the shortest paths. In contrast, LXYROPT had the lowest crosstalk noise because it had the fewest routing paths that could communicate simultaneously, owing to low path diversity. Additionally, shorter signal paths in LXYROPT imply that fewer optical components are involved, thereby reducing the insertion loss and further enhancing the WC OSNR. As a result, LXYROPT, with fewer simultaneous signal connections and emphasis on shortest-path communications, achieved a higher WC OSNR than the OWBM.

In contrast to the mesh topology, the WC OSNR gap between the multicast routing techniques was less pronounced for the torus topology. Because of the interconnected edge routers in the torus topology, which significantly shorten the longest path compared to the mesh topology, the shortened signal path results in diminished IL and crosstalk noise. In conclusion, if the cost associated with increased waveguide lengths in the torus topology is manageable, a torus topology with higher WC OSNR can offer better energy efficiency than a mesh topology, particularly for reducing the laser source power.

By contrast, on the torus-based 64-core HGC platform, the wavelength utilization of OWBM was improved by 30.74, 36.85, and 16.51% compared with those of RWADMM, MP, and MRCN, respectively. The higher path diversity

**TABLE 6.** Network performance comparison in 16-, 32-, 64-core HGC platform with VGG-16, VGG-19, ResNet applications.

| Mesh | | 16 cores | | | | 32 cores | | | | 64 cores | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Throughput | WC OSNR | Wavelength utilization | MI | Throughput | WC OSNR | Wavelength utilization | MI | Throughput | WC OSNR | Wavelength utilization | MI |
| VGG-16 | LXYROPT | 31.997584 | 11.495325 | 0.075895 | | 39.772365 | 8.223973 | 0.048599 | | 39.234723 | 4.328437 | 0.030277 | |
| | RWADMM | 39.526162 | 8.962813 | 0.105477 | | 38.522124 | 7.722755 | 0.071129 | | 49.123156 | 3.954330 | 0.041944 | |
| | RWMMRG | 41.264247 | 8.385725 | 0.106753 | 11.50% | 40.235871 | 6.966152 | 0.072547 | 16.23% | 47.920296 | 3.220632 | 0.047726 | 18.88% |
| | MP | 34.486531 | 7.452176 | 0.078221 | | 39.071647 | 5.948417 | 0.064177 | | 45.219896 | 2.832430 | 0.035684 | |
| | MRCN | 43.325904 | 7.234514 | 0.106342 | | 47.974599 | 6.304363 | 0.068337 | | 57.964821 | 3.481304 | 0.042616 | |
| | OWBM | 45.871789 | 8.999794 | 0.117935 | | 51.438835 | 8.134535 | 0.078596 | | 56.990181 | 4.251215 | 0.052616 | |
| VGG-19 | LXYROPT | 34.805362 | 11.354296 | 0.080379 | | 39.267857 | 7.352800 | 0.057660 | | 39.168833 | 3.880786 | 0.035791 | |
| | RWADMM | 43.299719 | 9.034345 | 0.097037 | | 47.296750 | 7.632420 | 0.068589 | | 44.436895 | 3.488287 | 0.037818 | |
| | RWMMRG | 42.656708 | 7.972122 | 0.120705 | 15.46% | 46.898449 | 6.610865 | 0.074286 | 20.03% | 56.065349 | 3.666674 | 0.044412 | 22.57% |
| | MP | 44.724558 | 6.780817 | 0.100496 | | 43.270755 | 5.620776 | 0.059932 | | 51.261507 | 2.916719 | 0.039689 | |
| | MRCN | 51.004935 | 7.009608 | 0.123972 | | 54.417878 | 6.191480 | 0.067973 | | 55.938293 | 3.017114 | 0.044271 | |
| | OWBM | 51.299570 | 9.692778 | 0.134088 | | 59.422994 | 6.650312 | 0.085549 | | 56.832271 | 3.482726 | 0.050883 | |
| ResNet | LXYROPT | 38.279711 | 11.586939 | 0.086274 | | 39.261421 | 8.636722 | 0.057414 | | 39.395385 | 3.640078 | 0.042268 | |
| | RWADMM | 49.713818 | 9.338762 | 0.123935 | | 51.465350 | 6.821518 | 0.088345 | | 48.895067 | 3.533584 | 0.046974 | |
| | RWMMRG | 50.446479 | 8.053731 | 0.112222 | 22.48% | 60.103132 | 6.327028 | 0.071742 | 25.53% | 59.476344 | 3.482429 | 0.049975 | 27.13% |
| | MP | 47.794784 | 8.604838 | 0.102812 | | 49.971459 | 5.332798 | 0.068571 | | 55.531262 | 3.031449 | 0.047290 | |
| | MRCN | 51.438179 | 7.136585 | 0.134592 | | 63.156179 | 6.076041 | 0.090058 | | 63.501863 | 3.069556 | 0.050782 | |
| | OWBM | 57.086795 | 8.918681 | 0.145081 | | 65.831671 | 7.791535 | 0.093159 | | 68.799136 | 3.481912 | 0.052841 | |

| Torus | | 16 cores | | | | 32 cores | | | | 64 cores | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Throughput | WC OSNR | Wavelength utilization | MI | Throughput | WC OSNR | Wavelength utilization | MI | Throughput | WC OSNR | Wavelength utilization | MI |
| VGG-16 | RWADMM | 48.875568 | 12.052450 | 0.095760 | | 48.571349 | 10.537979 | 0.055718 | | 62.369112 | 8.032330 | 0.032317 | |
| | MP | 41.609904 | 11.333833 | 0.067900 | 11.50% | 47.072131 | 9.856612 | 0.049210 | 16.23% | 57.999547 | 7.661310 | 0.025778 | 18.88% |
| | MRCN | 56.045607 | 11.551188 | 0.095485 | | 60.630654 | 10.069413 | 0.053697 | | 73.034979 | 7.761964 | 0.032921 | |
| | OWBM | 60.344470 | 12.255166 | 0.108171 | | 67.295035 | 10.696417 | 0.063419 | | 76.205018 | 8.103732 | 0.042768 | |
| VGG-19 | RWADMM | 54.227633 | 12.149715 | 0.087839 | | 60.928560 | 10.458248 | 0.053540 | | 56.822532 | 7.531715 | 0.027731 | |
| | MP | 53.862010 | 11.313440 | 0.091464 | 15.46% | 53.902120 | 9.777968 | 0.044871 | 20.03% | 63.663179 | 7.216833 | 0.029341 | 22.57% |
| | MRCN | 65.494184 | 11.609784 | 0.113934 | | 71.319631 | 9.987092 | 0.052995 | | 71.762854 | 7.328521 | 0.034520 | |
| | OWBM | 66.938715 | 12.392328 | 0.123307 | | 77.077483 | 10.572319 | 0.071381 | | 77.446059 | 7.596541 | 0.040549 | |
| ResNet | RWADMM | 61.091356 | 12.514705 | 0.114428 | | 64.008744 | 9.628092 | 0.072483 | | 64.544028 | 7.603442 | 0.036664 | |
| | MP | 55.722804 | 11.640801 | 0.093816 | 22.48% | 60.990937 | 9.025428 | 0.053529 | 25.53% | 68.792628 | 7.274560 | 0.037279 | 27.13% |
| | MRCN | 61.567596 | 11.931417 | 0.125069 | | 78.660198 | 9.234015 | 0.075030 | | 79.013718 | 7.390872 | 0.041086 | |
| | OWBM | 75.539324 | 12.703224 | 0.135055 | | 87.317017 | 9.793361 | 0.078892 | | 89.908571 | 7.670697 | 0.043126 | |

MI:Multicast intensity ratio for the total communication.

of OWBM allows for a varied selection of connectable paths along with an on-the-fly wavelength allocation method that efficiently identifies available wavelength candidates, leading to the highest wavelength utilization.

In summary, the enhanced connectivity of the torus topology, in contrast to the mesh topology, improves the throughput and wavelength utilization capabilities of OWBM. However, as the longest path became shorter in the torus than in the mesh, the OSNR improvement of OWBM decreased.

The normalized wavelength utilizations of OWBM with the comparison group on the mesh-based 16-, 32-, and 64-core HGC platforms are shown in Fig. 11. OWBM, which allocates wavelengths on the fly to each branch point in the route computation stage, exhibits the highest wavelength utilization among multicast routing schemes. In the VGG-16 benchmark of the 16-core HGC platform, the wavelength utilization of OWBM was 35.65, 10.56, 9.48, 33.67, and 9.83% higher than those of LXYROPT, RWADMM, RWMMRG, MP, and MRCN, respectively.
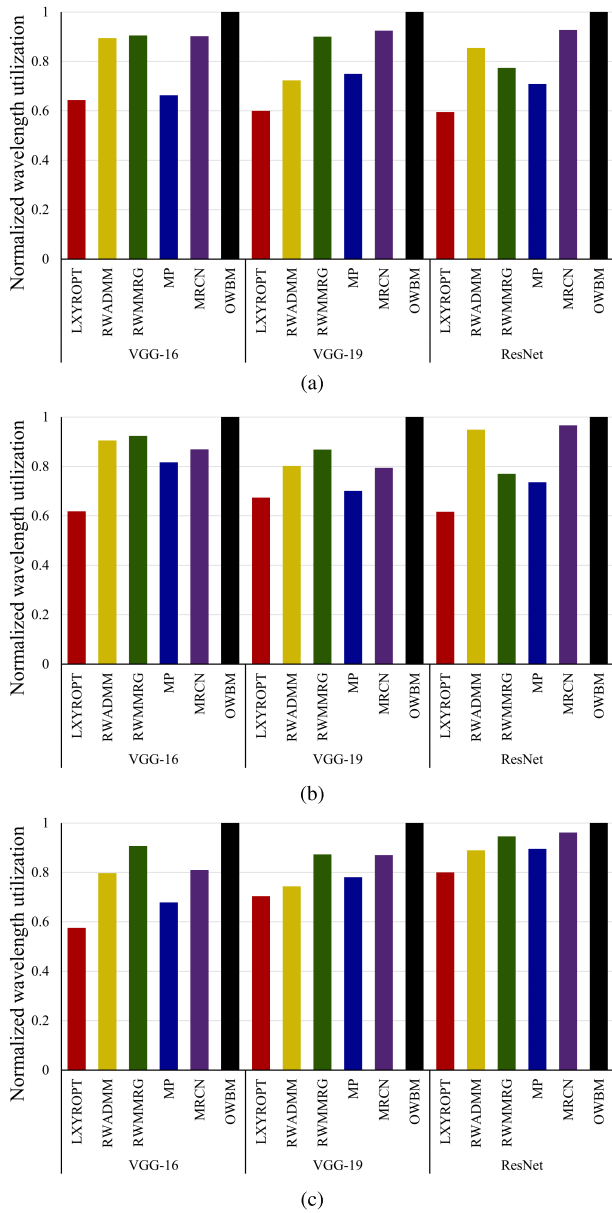
**FIGURE 11.** Normalized wavelength utilization comparison in (a) 16-core (b) 32-core (c) 64-core HGC platform with the mesh topology.

LXYROPT and MP, which communicate through a fixed routing path according to the destination group, have low path diversity and severe overlap of multiple paths; thus, in most cases, a single wavelength cannot be used for multiple paths. As a result, the wavelength utilization of the LXYROPT and MP tended to be significantly lower than that of the OWBM. RWADMM and RWMMRG proposed four routing paths, XY, YX, XYX, and YXY, in addition to pre-allocating all wavelengths at the source before attempting routing. Thus, they showed higher wavelength utilization than LXYROPT. However, RWADMM and RWMMRG cannot surpass the wavelength utilization of OWBM, which searches all available paths and assigns wavelengths on the fly.

In particular, the improvement in the wavelength utilization of OWBM was intensified on the HGC platform with a large network size, where the multicast intensity and network complexity increased. In the VGG-16 benchmark of the 64-core HGC platform, OWBM had 42.46, 20.28, 9.29, 32.18, and 19.01% higher wavelength utilization than LXYROPT, RWADMM, RWMMRG, MP, and MRCN, respectively. A hybrid routing technique combining tree- and path-based routing contributes to OWBM scalability as the network size increases.

Meanwhile, as the multicast intensity increases from the VGG networks to the ResNet benchmark, the gap in wavelength utilization between the multicast routing schemes decreases. In the VGG-16 benchmark, OWBM had averages 38.76, 13.45, 8.82, 28.07, and 13.96% higher than those of LXYROPT, RWADMM, RWMMRG, MP, and MRCN, respectively. In the ResNet application, the wavelength utilization of OWBM increased by an average of 32.97, 10.28, 17.02, 22.01, and 4.82% compared with those of LXYROPT, RWADMM, RWMMRG, MP, and MRCN, respectively. This decrease in wavelength utilization differences occurred because as the number of links using the four maximum wavelengths increased, the overall network reached saturation, and the gap between each multicast routing technique decreased. Because the MRCN assumed optimal wavelength allocation, the MRCN characterized by high path diversity showed a slight difference from OWBM. If a practical wavelength allocation logic such as OWBM is inserted into the MRCN, the difference in wavelength utilization is expected to widen further.

Fig. 12 shows the normalized latency per bit for the mesh-based 16-, 32-, and 64-core HGC platforms. The MRCN and OWBM, which propose many alternative routes through adaptive branching along with destination partitioning, generally exhibit lower latencies than tree-based routing. This tendency intensified in ResNet, which had the highest multicast intensity. For example, in the 32-core HGC platform within the VGG-16 application, the latencies of MRCN and OWBM were 18.40 and 23.90% lower on average, respectively, than those of the tree-based routing schemes LXYROPT and RWADMM. In the ResNet benchmark, MRCN and OWBM showed average latencies that were 28.17 and 31.09% lower than those of LXYROPT and RWADMM, respectively.

Meanwhile, the latency of OWBM was lower than that of the MRCN by an average of 5.34, 6.41, and 2.52% on the 16-, 32-, and 64-core HGC platforms, respectively. OWBM exhibited a relatively lower latency than the MRCN because it maximized the path diversity through tree-based routing compared with MRCN, which only performs adaptive branching based on path-based routing. Overall, OWBM showed an average of 32.32, 22.45, 17.19, 25.11, and 6.41% lower latency than LXYROPT, RWADMM, RWMMRG, MP, and MRCN, respectively, on the 32-core HGC platform. The latency decrease in OWBM is considered to be an effect of ADB which deploys hybrid multicast routing.
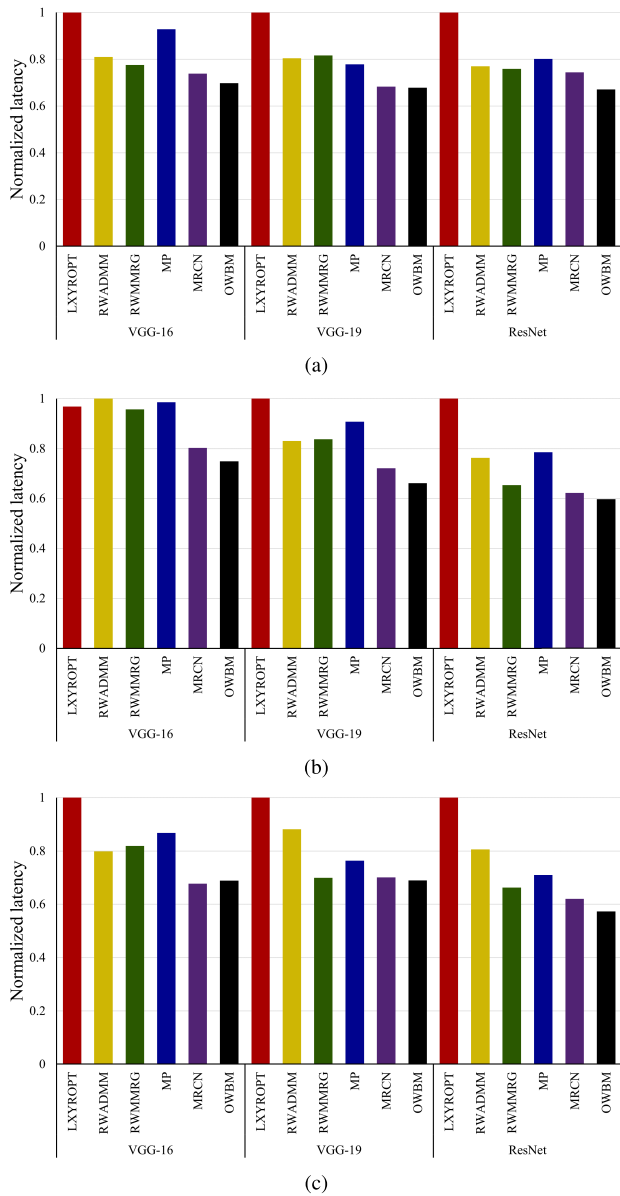
**FIGURE 12.** Normalized latency per bit comparison in (a) 16-core (b) 32-core (c) 64-core HGC platform with the mesh topology.



**FIGURE 13.** Energy per bit comparison in (a) 16-core (b) 32-core (c) 64-core HGC platform with the mesh topology.

An energy per bit comparison for the 16-, 32-, and 64-core HGC platforms with mesh topology is described in Fig. 13. Tree-based multicast routing schemes, which aim for communication along the shortest path, have a higher OSNR than the MP and MRCN schemes, contributing to a reduction in the laser source power. Accordingly, on the 16-core HGC platform, the energies per bit of LXYROPT, RWADMM, and RWMMRG were lower than those of MP and MRCN. However, as the number of cores increased, the packet latency of MRCN decreased dramatically, as more alternative routes were presented. Hence, the MRCN, which presents high congestion response ability, showed better energy efficiency than tree-based multicast routing schemes. For example, on the 64-core HGC platform, the energy per
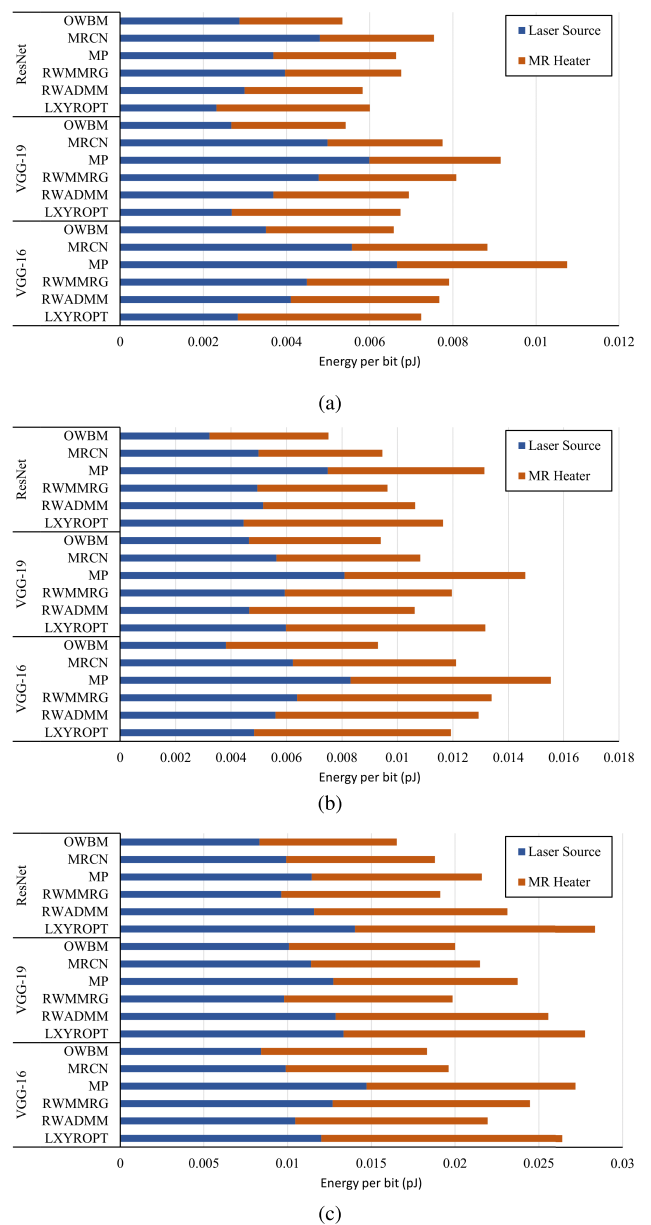
bit of the MRCN in ResNet was 25.67 and 10.61% lower than that of LXYROPT and RWADMM, respectively.

In the 16-core HGC platform, the OWBM showed an average of 28.23, 33.00, 23.56, 11.64, and 13.21% lower energies per bit than LXYROPT, RWADMM, RWMMRG, MP, and MRCN, respectively. When the network size is small, the congestion owing to multiple multicasts is relatively low. Hence, the advantages of OWBM and MRCN, which avoid congestion through adaptive branching on the fly, are not extensively highlighted in small-sized networks. However, on the 64-core HGC platform, OWBM showed an average of 8.53, 23.94, 12.64, 22.26, and 33.39% lower energy per bit than LXYROPT, RWADMM, RWMMRG,

MP, and MRCN, respectively. Consequently, the decrease in the average packet latency due to the adaptive branching of OWBM contributed to the overall energy consumption reduction.

Meanwhile, as the size of the convolutional layer increased from VGG-16 to ResNet and the multicast intensity increased, the reduction in packet latency due to contention avoidance through ADB was maximized. On the 32-core HGC platform, the energy consumption of OWBM was 21.43 and 22.11% lower than that of RWMMRG for VGG-16 and VGG-19, respectively. In the ResNet benchmark of the 32-core HGC platform, the OWBM had 30.57% less energy per bit than the RWMMRG. In conclusion, OWBM showed the highest energy efficiency among the multicast routing methods.

### C. AREA AND POWER ESTIMATIONS OF ELECTRICAL ROUTERS

The Synopsys Design Compiler® tool with the TSMC 40 nm process library was used to analyze the electrical routers of OWBM and baseline with and without multicast routing support. Dynamic power analysis was performed on the components of an electrical router for the VGG-16 application using the switching activity interchange format (SAIF). We measured the area and power consumption of 4- and 5-port electrical routers, which are widely used in both regular and customized topologies. The electrical router architecture of OWBM was estimated by comparing it with the wormhole routing-based electrical router (WER) proposed by Mo et al. as the baseline [8]. WER involves the transmission of control packets in the electrical router prior to the path setup in the optical router. The WER then sends MR control signals to map the input and output ports of the optical router according to the corresponding ports. For a fair comparison, the WER was expanded to resemble a typical electrical router by adding the logic required for switch allocation and route computation, thereby maintaining a 5-stage pipeline.

Table 7 lists the area measurements by component for the 4- and 5-port routers. The electrical router of the ONoC for MR control does not contain data in the packet. Therefore, the first-in-first-out (FIFO) buffer depth of the ONoC router is relatively lower than that of the router in ENoC. Because PRT, TRT, IFT, and CCT were inserted into the routing logic of OWBM, the area of OWBM increased by 2.26 and 2.24 times at 4- and 5-port routers, respectively, compared to the baseline. The increase in the number of ports led to an increase in the area of LFT and CCT due to the expansion of the port mapping scenarios. However, because the routing table does not depend on the number of ports, the area overhead of the routing logic is alleviated.

Because the wavelength port mapping table is involved in the switch allocator by storing the wavelength allocation states, OWBM occupied areas 13.60 and 7.90% larger than the switch allocator of the baseline at ports 4 and 5, respectively. In the crossbar switch of OWBM, the branching

**TABLE 7.** Area analysis of the electrical router of baseline and OWBM.

| Port | Component | Baseline Area ($\mu m^2$) | Baseline Proportion | OWBM Area ($\mu m^2$) | OWBM Proportion |
|---|---|---|---|---|---|
| 4 | FIFO buffers | 15766.49 | 46.55% | 16569.47 | 29.45 % |
| | Routing logic | 1707.34 | 5.04 % | 3861.60 | 6.86 % |
| | Wavelength allocator | 0 | 0 % | 961.65 | 1.71 % |
| | Switch allocator | 6594.61 | 19.47 % | 7491.45 | 13.31% |
| | Crossbar | 3347.69 | 9.88% | 4531.58 | 8.05 % |
| | Total | 27416.13 | 100% | 33415.74 | 100 % |
| 5 | FIFO buffers | 20124.07 | 38.62% | 20769.07 | 25.19 % |
| | Routing logic | 2036.13 | 3.91 % | 4571.03 | 5.54 % |
| | Wavelength allocator | 0 | 0 % | 1148.91 | 1.39 % |
| | Switch allocator | 8183.47 | 15.71 % | 8829.83 | 10.71% |
| | Crossbar | 3844.35 | 7.38% | 5248.87 | 6.37 % |
| | Total | 34188.02 | 100% | 40567.71 | 100 % |

unit is integrated for packet duplication, so that identical packets can be transmitted to multiple ports simultaneously. Hence, the branching unit increased the area of OWBM by approximately 35.36% compared with that of the baseline. Overall, the total area of the electrical router in OWBM increased by 21.88 and 18.66% at ports 4 and 5, respectively, compared with that of the baseline. Because OWBM includes routing tables and OWA logics that are not proportional to the number of ports, the overhead of multicast routing tends to be alleviated as the number of ports increases.

**TABLE 8.** Power comparison of baseline and OWBM.

| Port | Baseline Static power ($\mu W$) | Baseline Dynamic power ($\mu W$) | OWBN Static power ($\mu W$) | OWBN Dynamic power ($\mu W$) |
|---|---|---|---|---|
| 4 | 136.32 | 6900.10 | 168.11 | 7877.77 |
| 5 | 195.25 | 8536.99 | 203.46 | 9511.13 |

The static and dynamic power consumptions of the routers at ports 4 and 5 are presented in Table 8. In the area estimation of the electrical router, OWBM uses router components with a larger area than the baseline. According to this trend, the static and dynamic power consumptions are also higher than that of the baseline. In the 4-port router, the static and dynamic power consumptions of the OWBM electrical router were 11.84 and 46.48% higher than those of the baseline, respectively. Because the LFT, CCT, and branching units are proportional to the number of ports, the power consumption of OWBM increases as the number of ports increases. In the 5-port router, the static and dynamic power consumptions of OWBM increased by 4.20 and 11.41%, respectively, compared with those of the baseline. The reduction in the power consumption overhead in 5-port router is interpreted

as an effect of the routing table, which is independent of the number of ports.

## V. CONCLUSION

For energy-efficient multicasting of the HGC platform based on a custom topology, we proposed OWBM, which performs ADB by integrating tree- and path-based multicast routing and assigns wavelengths on the fly by referring to the state of adjacent routers. OWBM achieved up to 67.68% throughput improvement and 42.79% energy consumption reduction on a 32-core HGC platform compared to existing multicast routing techniques. The area and power overheads of the electrical router owing to the added LUTs were 18.66 and 11.41%, respectively, compared with the 5-port baseline router. OWBM improves the throughput and energy efficiency by allocating wavelengths and performing branching on the fly in the route computation stage based on OSNR-aware hybrid routing. OSNR-aware 2-stage port filtering technique contributes to the optimization of the laser power by improving the WC OSNR. In essence, OWBM represents a significant advancement in the HGC platform, striking a balance between enhanced performance and minimal resource overhead, thus providing energy-efficient communication in multicast-intensive benchmarks. In particular, the OWBM methodology holds significant potential for extension and application in the realm of generative AI models using transformers, which cause aggressive multicast intensity for input data parallelism.

## REFERENCES

[1] S. Mittal and J. S. Vetter, "A survey of CPU-GPU heterogeneous computing technique," *ACM Comput. Surv.*, vol. 47, no. 4, p. 69, Jul. 2015.

[2] M. J. Schulte, M. Ignatowski, G. H. Loh, B. M. Beckmann, W. C. Brantley, S. Gurumurthi, N. Jayasena, I. Paul, S. K. Reinhardt, and G. Rodgers, "Achieving exascale capabilities through heterogeneous computing," *IEEE Micro*, vol. 35, no. 4, pp. 26–36, Jul. 2015.

[3] P. Dong, Y.-K. Chen, T. Gu, L. L. Buhl, D. T. Neilson, and J. H. Sinsky, "Reconfigurable 100 Gb/s silicon photonic network-on-chip [invited]," *J. Opt. Commun. Netw.*, vol. 7, no. 1, pp. A37–A43, Jan. 2015.

[4] K.-H. Lee, D. J. Shin, H.-C. Ji, K. W. Na, S. G. Kim, J. K. Bok, Y. S. You, S. S. Kim, I. S. Joe, S. D. Suh, J. H. Pyo, Y. H. Shin, K. H. Ha, Y. D. Park, and C. H. Chung, "10Gb/s silicon modulator based on bulk-silicon platform for DRAM optical interface," in *Proc. Opt. Fiber Commun. Conf. Expo. Nat. Fiber Optic Eng. Conf.*, Mar. 2011, pp. 1–3.

[5] P. Boncz, T. Neumann, and O. Erling, "TPC-H analyzed: Hidden messages and lessons learned from an influential benchmark," in *Proc. Technol. Conf. Perform. Eval. Benchmarking.* Cham, Switzerland: Springer, Aug. 2013, pp. 61–76.

[6] The Standard Performance Evaluation Corporation. (1999). *SPECWeb99.* [Online]. Available: http://www.spec.org/osg/web99

[7] M. P. Malumbres, J. Duato, and J. Torrellas, "An efficient implementation of tree-based multicast routing for distributed shared-memory multiprocessors," in *Proc. 8th IEEE Symp. Parallel Distrib. Process. (SPDP)*, Oct. 1996, pp. 186–189.

[8] K. H. Mo, Y. Ye, X. Wu, W. Zhang, W. Liu, and J. Xu, "A hierarchical hybrid optical-electronic network-on-chip," in *Proc. IEEE Comput. Soc. Annu. Symp. VLSI*, Jul. 2010, pp. 327–332.

[9] L. Liu and Y. Yang, "Energy-aware routing in hybrid optical network-on-chip for future multi-processor system-on-chip," in *Proc. ACM/IEEE Symp. Architectures Netw. Commun. Syst. (ANCS)*, Oct. 2010, pp. 1–9.

[10] H. Li, H. Gu, Y. Yang, and X. Yu, "A hybrid packet-circuit switched router for optical network on chip," *Comput. Elect. Eng.*, vol. 39, no. 7, pp. 2197–2206, 2013.

[11] Y. Ye, X. Wu, J. Xu, M. Nikdast, Z. Wang, X. Wang, and Z. Wang, "System-level analysis of mesh-based hybrid optical-electronic network-on-chip," in *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS)*, May 2013, pp. 321–324.

[12] A. Ding and G. S. Poo, "A survey of optical multicast over WDM networks," *Comput. Commun.*, vol. 26, no. 2, pp. 193–200, 2003.

[13] Y. Xie, J. Xu, J. Zhang, Z. Wu, and G. Xia, "Crosstalk noise analysis and optimization in 5×5 hitless silicon-based optical router for optical networks-on-chip (ONoC)," *J. Lightw. Technol.*, vol. 30, no. 1, pp. 198–203, Jan. 2012.

[14] S. Werner, J. Navaridas, and M. Luján, "A survey on optical network-on-chip architectures," *ACM Comput. Surv.*, vol. 50, no. 6, pp. 1–37, Dec. 2017.

[15] Y. W. Kim, J. H. Lee, and T. H. Han, "Extended worst-case OSNR searching algorithm for optical network-on-chip using a semi-greedy heuristic with adaptive scan range," *IEEE Access*, vol. 8, pp. 125863–125873, 2020.

[16] N. E. Jerger, L. S. Peh, and M. Lipasti, "Virtual circuit tree multicasting: A case for on-chip hardware multicast support," in *Proc. 35th IEEE Int. Symp. Comput. Archit. (ISCA)*, Jun. 2008, pp. 229–240.

[17] W. Hu, Z. Lu, A. Jantsch, and H. Liu, "Power-efficient tree-based multicast support for networks-on-chip," in *Proc. 16th Asia South Pacific Design Autom. Conf. (ASP-DAC)*, Jan. 2011, pp. 363–368.

[18] F. Nasiri, H. Sarbazi-Azad, and A. Khademzadeh, "Reconfigurable multicast routing for networks on chip," *Microprocessors Microsyst.*, vol. 42, pp. 180–189, May 2016.

[19] E. A. Carara and F. G. Moraes, "Deadlock-free multicast routing algorithm for wormhole-switched mesh networks-on-chip," in *Proc. IEEE Comput. Soc. Annu. Symp. VLSI*, Apr. 2008, pp. 341–346.

[20] M. Daneshtalab, M. Ebrahimi, T. C. Xu, P. Liljeberg, and H. Tenhunen, "A generic adaptive path-based routing method for MPSoCs," *J. Syst. Archit.*, vol. 57, no. 1, pp. 109–120, Jan. 2011.

[21] M. Ebrahimi, M. Daneshtalab, P. Liljeberg, and H. Tenhunen, "HAMUM—A novel routing protocol for unicast and multicast traffic in MPSoCs," in *Proc. 18th Euromicro Conf. Parallel, Distrib. Netw.-Based Process.*, Feb. 2010, pp. 525–532.

[22] B. Tiwari, M. Yang, Y. Jiang, and X. Wang, "Efficient on-chip multicast routing based on dynamic partition merging," in *Proc. 28th Euromicro Int. Conf. Parallel, Distrib. Netw.-Based Process. (PDP)*, Mar. 2020, pp. 274–281.

[23] Z. Y. Kang, S. M. Li, S. Y. Wang, L. H. Qu, R. Gong, W. Shi, W. X. Xu, and L. Wang, "Path-based multicast routing for network-on-chip of the neuromorphic processor," *J. Comput. Sci. Technol.*, vol. 38, no. 5, Sep. 2023, pp. 1098–1112.

[24] Y. S. Lee, Y. W. Kim, and T. H. Han, "MRCN: Throughput-oriented multicast routing for customized network-on-chips," *IEEE Trans. Parallel Distrib. Syst.*, vol. 34, no. 1, pp. 163–179, Jan. 2023.

[25] B. Thenral, S. Kumari, P. Iswarya, and K. Arthi, "A hybrid multicast routing algorithm for networks-on-chip for improving connectivity," *Int. J. Innov. Trends Eng.*, vol. 14, no. 2, pp. 61–65, Feb. 2016.

[26] C.-W. Wu, K.-J. Lee, and A. P. Su, "A hybrid multicast routing approach with enhanced methods for mesh-based networks-on-chip," *IEEE Trans. Comput.*, vol. 67, no. 9, pp. 1231–1245, Sep. 2018.

[27] Y. Xie, T. Song, Z. Zhang, C. He, J. Li, and C. Xu, "Formal analysis of crosstalk noise in mesh-based optical networks-on-chip with WDM," *J. Lightw. Technol.*, vol. 34, no. 15, pp. 3550–3562, Aug. 2016.

[28] R. Tian, Q. Zhang, Z. Xiang, Y. Xiong, X. Li, and W. Zhu, "Robust and efficient path diversity in application-layer multicast for video streaming," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 15, no. 8, pp. 961–972, Aug. 2005.

[29] A. Ben Hassouna, H. Koubaa, and L. A. Saidane, "Multi-user diversity wireless multicast: A survey," *Comput. Netw.*, vol. 175, Jul. 2020, Art. no. 107282.

[30] S. Jain and S. R. Das, "Exploiting path diversity in the link layer in wireless ad hoc networks," *Ad Hoc Netw.*, vol. 6, no. 5, pp. 805–825, Jul. 2008.

[31] P. Dong, W. Qian, S. Liao, H. Liang, C.-C. Kung, N.-N. Feng, R. Shafiiha, J. Fong, D. Feng, A. V. Krishnamoorthy, and M. Asghari, "Low loss silicon waveguides for application of optical interconnects," in *Proc. IEEE Photon. Soc. Summer Topicals*, Jul. 2010, pp. 191–192.

[32] F. Xia, L. Sekaric, and Y. Vlasov, "Ultracompact optical buffers on a silicon chip," *Nature Photon.*, vol. 1, no. 1, pp. 65–71, Dec. 2007.

[33] W. Ding, D. Tang, Y. Liu, L. Chen, and X. Sun, "Compact and low crosstalk waveguide crossing using impedance matched metamaterial," *Appl. Phys. Lett.*, vol. 96, no. 11, pp. 111114–111116, 2010.

[34] J. Chan, G. Hendry, K. Bergman, and L. P. Carloni, "Physical-layer modeling and system-level design of chip-scale photonic interconnection networks," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 30, no. 10, pp. 1507–1520, Oct. 2011.

[35] M. S. Kim, Y. W. Kim, and T. H. Han, "System-level signal analysis methodology for optical network-on-chip using linear model-based characterization," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 39, no. 10, pp. 2761–2771, Oct. 2020.

[36] H. Gu, K. H. Mo, J. Xu, and W. Zhang, "A low-power low-cost optical router for optical networks-on-chip in multiprocessor systems-on-chip," in *Proc. IEEE Comput. Soc. Annu. Symp. VLSI*, May 2009, pp. 19–24.

[37] Y. Xie, M. Nikdast, J. Xu, W. Zhang, Q. Li, X. Wu, and W. Liu, "Crosstalk noise and bit error rate analysis for optical network-on-chip," in *Proc. 47th Design Automat. Conf.*, Jun. 2010, pp. 657–660.

[38] W. Yang, Y. Chen, Z. Huang, and H. Zhang, "RWADMM: Routing and wavelength assignment for distribution-based multiple multicasts in ONoC," in *Proc. IEEE Int. Symp. Parallel Distrib. Process. with Appl. IEEE Int. Conf. Ubiquitous Comput. Commun. (ISPA/IUCC)*, Dec. 2017, pp. 550–557.

[39] F. Gao, C. Yu, Y. Chen, and B. Gao, "Routing and wavelength assignment algorithm for mesh-based multiple multicasts in optical network-on-chip," *Theory Comput. Syst.*, vol. 1, no. 1, pp. 1–18, Oct. 2022.

[40] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014, *arXiv:1409.1556*.

[41] S. Targ, D. Almeida, and K. Lyman, "ResNet in ResNet: Generalizing residual architectures," 2016, *arXiv:1603.08029*.

[42] Y.-H. Chen, T. Krishna, J. S. Emer, and V. Sze, "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks," *IEEE J. Solid-State Circuits*, vol. 52, no. 1, pp. 127–138, Jan. 2017.

[43] J. Power, J. Hestness, M. S. Orr, M. D. Hill, and D. A. Wood, "gem5-GPU: A heterogeneous CPU-GPU simulator," *IEEE Comput. Archit. Lett.*, vol. 14, no. 1, pp. 34–36, Jan./Jun. 2015.

[44] N. Binkert, B. Beckmann, G. Black, S. K. Reinhardt, A. Saidi, A. Basu, J. Hestness, D. R. Hower, T. Krishna, S. Sardashti, R. Sen, K. Sewell, M. Shoaib, N. Vaish, M. D. Hill, and D. A. Wood, "The gem5 simulator," *ACM SIGARCH Comput. Archit. News*, vol. 29, no. 2, pp. 1–7, 2011.

[45] A. Bakhoda, G. Yuan, W. Fung, H. Wong, and T. M. Aamodt, "Analyzing CUDA workloads using a detailed GPU simulator," in *Proc. IEEE Int. Symp. Perform. Anal. Syst. Softw.*, Apr. 2009, pp. 163–174.

[46] F. E. Doany, B. G. Lee, S. Assefa, W. M. J. Green, M. Yang, C. L. Schow, C. V. Jahnes, S. Zhang, J. Singer, V. I. Kopp, J. A. Kash, and Y. A. Vlasov, "Multichannel high-bandwidth coupling of ultradense silicon photonic waveguide array to standard-pitch fiber array," *J. Lightw. Technol.*, vol. 29, no. 4, pp. 475–482, Feb. 2011.

[47] I. O'Connor and G. Nicolescu, *Integrated Optical Interconnect Architectures for Embedded Systems*. Berlin, Germany: Springer, 2012.

[48] P. Dong, S. Liao, H. Liang, W. Qian, X. Wang, R. Shafiiha, D. Feng, G. Li, X. Zheng, A. V. Krishnamoorthy, and M. Asghari, "High-speed and compact silicon modulator based on a racetrack resonator with a 1 V drive voltage," *Opt. Lett.*, vol. 35, no. 19, pp. 3246–3248, Oct. 2010.

[49] S. Xiao, M. H. Khan, H. Shen, and M. Qi, "Multiple-channel silicon micro-resonator based filters for WDM applications," *Opt. Exp.*, vol. 15, no. 12, pp. 7489–7498, Jun. 2007.

[50] G. Zhou, X. Li, and N. Feng, "Design of deeply etched antireflective waveguide terminators," *IEEE J. Quantum Electron.*, vol. 39, no. 2, pp. 384–391, Feb. 2003.

**YONG WOOK KIM** (Student Member, IEEE) received the B.S. degree in electronic and electrical engineering from Sungkyunkwan University, Suwon, South Korea, in 2016, where he is currently pursuing the M.S. and Ph.D. degrees in electrical and computer engineering. His research interests include NoC, machine learning, and computer architecture.

**TAE HEE HAN** (Senior Member, IEEE) received the B.S., M.S., and Ph.D. degrees in electrical engineering from Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 1992, 1994, and 1999, respectively. From 1999 to 2006, he was with the Telecom Research and Development Center, Samsung Electronics, where he developed 3G wireless, mobile TV, and mobile WiMax handset chipsets. Since March 2008, he has been with Sungkyunkwan University, Suwon, South Korea, as a Professor. From 2011 to 2013, he was a full-time Advisor on system ICs for Korean Government. His current research interests include SoC/chiplet architectures for AI, advanced memory architecture, network-on-chip, and system-level design methodologies.

• • •