

# Multicast-Aware High-Performance with Encrypt Wireless Network-On-Chip Architectures

A.Padma Priya<sup>1</sup>, M.Ashok Kumar<sup>2</sup>,

M.E. VLSI Design, Adhiyamaan College of Engineering, Hosur, TN, India<sup>1</sup>

Asst. Professor, Adhiyamaan College of Engineering, Hosur, TN, India<sup>2</sup>

**ABSTRACT:** Conventional NoCs are designed predominantly for unicast data exchanges. In such NoCs, the multicast traffic is generally handled by converting each multicast message to multiple unicast transmissions. Hence, applications dominated by multicast traffic experience high queuing latencies and significant performance penalties when running on systems designed with unicast-based NoC architectures. Various multicast mechanisms such as XY-tree multicast and path multicast have already been proposed to enhance the performance of the traditional wireline mesh NoC incorporating multicast traffic. However, even with such added features, the multihop nature of the wireline mesh NoC leads to high network latencies and thus limits the achievable system performance. In this paper, to sustain the high-bandwidth and high-throughput requirements of emerging applications, we propose the design of a wireless NoC (WiNoC) architecture incorporating necessary multicast support. By integrating congestion-aware multicast routing with network coding, the WiNoC is able to efficiently handle heavy multicast injections. For applications running with a broadcast-heavy Hammer cache coherence protocol, the proposed multicast-aware WiNoC achieves an average of 47% reduction in message latency compared with the XY-tree-based multicast-aware mesh NoC. This network level improvement translates into a 26% saving in full-system energy delay product.

**KEYWORDS:** Multicast communication, network coding (NC), networks on chip (NoCs), system-on-chip, wireless NoCs (WiNoC), XILINX ISE, MODELSIM

## I.INTRODUCTION

UNICAST leads to following problems;

- 1) local congestion on the source node;
- 2) poorly guaranteed QoS;
- 3) global congestion due to competition for the same network resource among repeated unicast packets;
- 4) power consumption overhead due to repeated injections

Therefore, without properly supported by the underlying NoC architectures and routing protocols, the overall network performance deteriorates sharply even as the percentage of collective communication increases by a small amount. This prevents NoCs to exploit massive-scale parallelism in many critical large-scale applications, which are inherently embedded with considerable collective traffic requirements. Neural-network-based computing real-time object recognition processing, genetic-algorithm-based protein folding analysis, and neuromorphic computing are a few of the SoC applications that exhibit significant multicast traffic patterns. Moreover, many SoC control functions such as passing global states, managing, and configuring the on-chip networks and implementing: 1) operand exchange networks; 2) nonuniform cache architectures; and 3) cache

coherency protocols also require efficient multicast data exchange. To elaborate on this, as examples, we consider the traffic patterns associated with two different cache coherence protocols, directory and Hammer protocols. Both directory and Hammer protocols can give rise to traffic hotspots and generate significant number of long-range traffic patterns warranting low-latency NoC links, even among the physically far apart processing cores. However, unlike the directory protocol that involves only sparse multicasts, the Hammer protocol involves dense broadcast injections and hence stresses the NoC affecting the overall performance when implemented on systems predominantly designed for unicast-based transmissions. Moreover, the multicast traffic arising due to the cache coherence protocols is inherent to the many-core platform and can be observed irrespective of the executed applications. On the other hand, the on-chip area overhead associated with implementing the Hammer protocol is very low, and hence, it is highly preferable for large many-core chips than the traditional directory protocol. Hence, designing a high-performance multicast-aware NoC that enables an efficient implementation of Hammer cache coherence protocol is a highly relevant research topic. As we show later, cache coherence protocols also exhibit multicast bursts, leading to high multicast traffic volumes at a particular instance of time and thus necessitating a multicast and congestion-aware.

Wireless NoC (WiNoC) is an emerging paradigm to design high-bandwidth and energy-efficient communication backbone for many-core chips. The previous works have already investigated the feasibility of the on-chip wireless communication. The viability of on-chip wireless communication has been demonstrated through prototype developments. In addition, a recent study on emerging on-chip interconnects concluded that the radio-frequency interconnects such as the millimeter-wave (mm-wave) wireless links and surface-wave interconnects are also more power and cost efficient than the on-chip optical links. Moreover, compared with surface-wave interconnects, the on-chip wireless technology is more mature and is completely CMOS compatible. The wireless channels of the WiNoC have inherent broadcast capability. Hence, WiNoC provides a natural solution for designing efficient multicast-aware NoCs.

In this paper, we present the design of a multicast-aware wireless NoC incorporating congestion-aware routing and network coding (NC) technique to handle high volumes of multicast injections. The resultant WiNoC improves the overall execution time of the many-core chip and lowers the total energy dissipation. We demonstrate the efficacy of the proposed WiNoC through case studies using the above-discussed Hammer and directory cache coherence protocols.

## II. RELATED WORK

Conventional NoCs convert multicast messages into multiple unicast messages. This causes sharp increases in traffic injection rates, which ultimately lead to high queuing delays and poor system performance. To eliminate these queuing delays and improve the latency in delivering multicast messages in mesh NoCs, path multicast methodologies have been proposed. In a path multicast methodology, for each multicast message, the overall mesh network is divided into multiple destination regions. The number of such regions and the size of the each destination region is determined by the source node of the message and the variation of the adopted path multicast algorithm. For each destination region, one copy of the original multicast message is created and this copy is routed across all the destinations of the region in ascending/descending order of the destination node addresses. As we show later, for cache-coherence-induced traffic patterns, the path multicast methodology can lead to large destination regions, which require a high number of hops to distribute the multicast messages. This ultimately leads to high network latency and poor system performance. A dynamic path multicast mechanism where the network is recursively divided into multiple small destination regions is proposed in. However, for all path multicast mechanisms, at each destination along a path, the header flit is needed to be repackaged with updated list of destinations. Hence, path multicast mechanism requires complex router architectures.

The XY-tree multicast mechanism for mesh NoCs has been proposed in, and. Under this methodology, the multicast message is first forwarded from the source node to all the intermediate nodes lying in the same row. The message is then replicated at these intermediate nodes and a copy of the original message is forwarded to

all the destinations lying in the same column. Unlike the path-based multicast mechanism, XY-tree multicast does not require a complex router and involves simple message forwarding

Recently, NoC architectures incorporating wires with clock less repeaters have been proposed (referred to as SMART NoCs). A multicast-aware SMART NoC is also proposed. The performance advantages of the SMART NoCs mainly come from the SMART control mechanism and this principle can be integrated with any NoC architecture. However, the SMART control mechanism restricts the operating frequency of interconnect (1–2 GHz) and involves high control overheads and requires more complex router design than the traditional NoC routers.

### CRYPTOGRAPHY

**Cryptography** or **cryptology** (from Greek κρυπτός *kryptós*, "hidden, secret"; and γράφειν *graphein*, "to write", or *-λογία -logia*, "study", respectively) is the practice and study of techniques for secure communication in the presence of third parties called adversaries. More generally, cryptography is about constructing and analyzing protocols that prevent third parties or the public from reading private messages; various aspects in information security such as data confidentiality, data integrity, authentication, and non-repudiation are central to modern cryptography. Modern cryptography exists at the intersection of the disciplines of mathematics, computer science, electrical engineering, communication science, and physics. Applications of cryptography include electronic commerce, chip-based payment cards, digital currencies, computer passwords, and military communications.

Cryptography prior to the modern age was effectively synonymous with *encryption*, the conversion of information from a readable state to apparent nonsense. The originator of an encrypted message shared the decoding technique needed to recover the original information only with intended recipients, thereby precluding unwanted persons from doing the same. The cryptography literature often uses the name Alice ("A") for the sender, Bob ("B") for the intended recipient, and Eve ("eavesdropper") for the adversary. Since the development of rotor cipher machines in World War I and the advent of computers in World War II, the methods used to carry out cryptology have become increasingly complex and its application more widespread.

### CHIP IMPLEMENTATION

The chip occupies  $2100 \times 2100 \mu\text{m}^2$  with  $1210 \times 1210 \mu\text{m}^2$  core area, including four  $512 \times 13$  single-port SRAMs, four  $8 \times 13$  caches, and four  $128 \times 13$  ROMs. The function of the chip is verified, and its performance is measured using a digital test station. This chip can operate at maximum 51 MHz with 33.3 mW. The original message is forwarded to all the destinations lying in the same column. Unlike the path-based multicast mechanism, XY-tree multicast does not require a complex router and involves simple message forwarding.

For cache-coherence-induced traffic patterns, each multicast message is associated with a set of acknowledgement (ACK) messages that are transmitted from each multicast destination back to the source node. To meet the requirements of the cache coherence communication, an XY-tree multicast NoC incorporating an ACK aggregation network is proposed in. A mesh NoC performing broadcasts over uncongested trees and incorporated with ACK aggregation and single-cycle multiport replication of broadcast flits is presented in. These works demonstrate that with an efficient replication mechanism and with the use of ACK gathering networks, the performance of the many-core systems incorporating Hammer cache coherence protocol can be greatly enhanced. However, in these systems, the underlying NoC is still a mesh network. In a mesh NoC, the network latency is usually high due to the inherent multihop nature of the system. High network latencies cause undesired delays in forwarding the multicast messages as well as in collecting the ACKs, leading to stalled processor cycles. In addition, high multicast injection rates can lead to heavy congestion in a mesh network where the XY-tree multicast follows a single fixed path to transmit data.

The use of the NC scheme to eliminate the queuing latencies caused by message arbitration in a router is proposed in. This router architecture called as the NoX router is shown to significantly outperform the traditional sequential arbitration routers. The application of the NC for the multicast transmissions in NoCs has been proposed in and these NoC designs use the NC scheme to efficiently distribute multiple multicast packets that are injected within a very short time window.

Recently, NoC architectures incorporating wires with clock less repeaters have been proposed (referred to as SMART NoCs). A multicast-aware SMART NoC is also proposed. The performance advantages of the SMART NoCs mainly come from the SMART control mechanism and this principle can be integrated with any NoC architecture. However, the SMART control mechanism restricts the operating frequency of interconnect (1–2 GHz) and involves high control overheads and requires more complex router design than the traditional NoC routers.

### CACHE COHERENCE TRAFFIC PATTERNS

In a chip multiprocessor (CMP) platform, cache coherence protocols are used to ensure consistency among the multiple cached copies of shared data. Cache coherence protocols play

a vital role in determining the nature of the on-chip network traffic. In general, selecting an appropriate cache coherence protocol for a CMP involves analyzing the area and traffic tradeoffs associated with the protocol .

In this paper, as case studies, we consider two cache coherence protocols: the two-level MESI directory protocol and the AMD's Hammer-based HyperTransport (HT) protocol with probe filters. Between these two protocols, Hammer requires the least area overhead while generating high broadcast traffic. Directory exhibits lower traffic intensity while occupying high area overheads.

#### A. Directory Protocol

In a directory-protocol-based CMP, each translation looks-side buffer entry is assigned to one of the cache controllers, called as the home node. The home node maintains the list of processing cores that have the most recent copy of the data block. Whenever a processing core  $S$  needs to update a data block or encounters a cache miss, it first communicates with the home node. For a read miss, if the requested data block is present on the chip, the home node forwards this request to the core  $D$  that has the most recent copy of the requested data. Core  $D$  responds by forwarding the data to the requested node ( $S$ ).

On receiving a write miss from node  $S$ , the home node sends an invalidation signal to all the nodes that are having a copy of the data block. Thus, this invalidation signal generates a multicast traffic. Upon receiving the invalidation signal, the cores respond with an ACK signal to the home directory.

The necessity to maintain a list of sharing cores for each data block causes high memory overheads in systems using the directory protocol. Furthermore, the fraction of this overhead (in total silicon area) linearly increases as the number of cores increases and thus limits the scalability of the system.

#### B. Hammer Protocol

In the Hammer protocol, similar to the directory protocol, each data block is assigned a home node. However, unlike the directory protocol, the home nodes in hammer protocols do not maintain the list of cores that are sharing the most recent copy of the data block. Instead, the home node broadcasts the read/write requests to all the cores in the system. This enables the Hammer protocol to have a low logic and memory overhead. Hence, the Hammer protocol enhances the scalability of the many-core system as it does not require large memory structures to keep the list of shared cores typically needed in directory-based protocols. However, the systems operating.

The average number of destinations associated with a multicast message in both directory and Hammer protocols. The destination count for the Hammer protocol broadcasts is 63. For the directory protocol multicasts, the average destination count varies from 12 to 22 among the considered applications. Fig. 5 presents the distributed nature of multicast destinations for the Hammer and directory protocols. Since Hammer multicasts

are mainly broadcasts, almost all the multicast messages need to be transferred to destinations in four different quadrants [quadrants are shown later in Fig. 9(a)]. Even with the directory protocol, about 65% of the total multicast messages need to be communicated to destinations in four different quadrants. Thus, the cache-coherence-induced multicasts create heavy long-range communication requirements. Fig. 6 shows the percentage of total traffic exchanged between NoC routers that are separated by  $h$ -mesh NoC hops under the directory protocol. All intertile communications such as the traffic flow caused by last-level caches, directories, and memory controllers are taken into account. From Fig. 6, it can be observed that even in the unicast heavy directory protocol, the traffic patterns are “long-range heavy” and about 70% of the total injected messages are transferred among routers that are greater than three hops apart ( $>7.5$  mm in a  $20$  mm  $\times$   $20$  mm die). These results (Figs. 4–6) indicate the necessity for employing a multicast-aware NoC that is efficient for long-range communication. In this scenario, wireless channels provide an efficient solution. A wireless channel is a broadcast communication medium capable of establishing one-hop links even between physically distant nodes. Hence, it can simultaneously transfer the data to destinations lying in all four quadrants with low latencies.

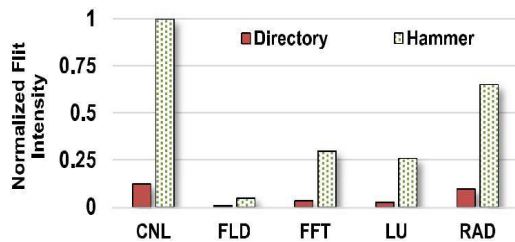


Fig. 2. Percentage of traffic arising due to multicast injections.

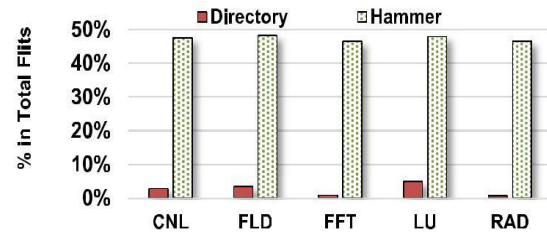


Fig. 1. Average flit intensity with the Hammer and directory protocols

## MULTICAST-AWARE WiNoC NETWORK DESIGN

As discussed earlier, in this paper, we strive to design congestion-aware low-hop-count multicast-enabled NoC. In this context, we present the design of an SW-network-based multicast-enabled WiNoC architecture.

### A. WiNoC Design Constraints

In our WiNoC, each processing core is connected to a router and the routers are interconnected following a power-law-based SW connectivity. Essentially, SW networks incorporate multiple short-range links and a few long-range shortcuts so as the overall inter-router hop count ( $H_{avg}$ ) is minimized. However, while adding links in SW NoCs, we need to follow certain restrictions. First, we should restrict the average number of ports ( $K_{avg}$ ) per router to four so that the WiNoC does not introduce any additional router port overhead compared with a conventional mesh. Next, we should restrict the maximum number of ports in a router ( $K_{max}$ ) so that no particular router becomes unnecessarily large. For a system size of 64 cores, the SW-based networks achieve highest throughput with lowest energy dissipation when the maximum port count in a router is restricted to seven. Finally, as we explain later in Section VI-B, for link lengths exceeding 6.25 mm in 28 nm CMOS technology node, wire-less links are more efficient than the traditional metal wires in terms of energy delay product (EDP). Hence, WiNoCs should prefer using wireless links instead of wireline links for data exchanges exceeding a communication distance of 6.25 mm. However, depending on the available wireless resources and allowed area overhead, we can make only a limited number of the longest links wireless, while the others need to remain wireline considering these facts, for our WiNoC, we follow a region-based wireless node placement strategy. In this strategy, the whole system is divided into multiple equal-sized nonoverlapping regions and each region is provided with a set of WIs. Principally, in this design, the short-range communication (within a

region) is handled through the wireline links, while the long-range communication (inter-region data exchange) is handled using wireless links. In each considered region, we place the WIs such that the communication distance from all the nonwireless nodes to their nearest WIs is minimized. The number of WIs allowed in a region is equal to the number of available nonoverlapping wireless channels. So far, it has been demonstrated that using on-chip mm-wave wireless links, it is possible to create three nonoverlapping channels. For a system of 64 nodes, it is already shown that placing the WIs at the center of each region minimizes the communication distance from any node to its nearest WI, and hence it ensures that the utilization of wireless medium is maximized. Fig. 9(a) shows an example where a system of 64 cores is divided into four regions (quadrants), with each quadrant having three WIs in its center.

#### *E. Distributed Routing and Deadlock Freedom*

Though MALASH multicast routing includes multiple features such as forwarding toward destination region WIs, congestion-aware wireline distribution and NC, it is distributed in nature and simple to implement. We use two additional fields  $C_m$  and  $R_m$  in the multicast headers (a total of four additional bits per message) to accomplish the distributed routing.  $C_m$  is a 2-bit field that indicates the channel ID of the source region WI to which the message is first forwarded from the source node.  $R_m$  is a 2-bit field that guides the wireline transmissions associated with the message. An  $R_m$  value of 0 would indicate that the message is being forwarded from the original source node to the source region WI. All  $R_m$  values that greater than 0 indicate that the message is undergoing regional wireline distribution. Tables I and II list the routing rules followed in MALASH to enable a distributed routing. As it can be observed from Table II, the MALASH routing can be simplified to four parallel conditional cases.

Aside from enabling the distributed routing, the  $C_m$  field can also be helpful in identifying the necessity for NC.

### III. EXPERIMENTAL RESULTS

#### *A. Experimental Setup*

In this paper, we evaluate the performance of the proposed WiNoC architecture against three other multicast-aware NoC architectures, XY-tree mesh NoC, path multicast mesh NoC, and the wireless mesh (WiMesh). These three architectures are shown in Fig. 14. We have already explained the salient features of the XY-tree mesh NoC and the path multicast mesh in Section II.

### IV. RESULT ANALYSIS

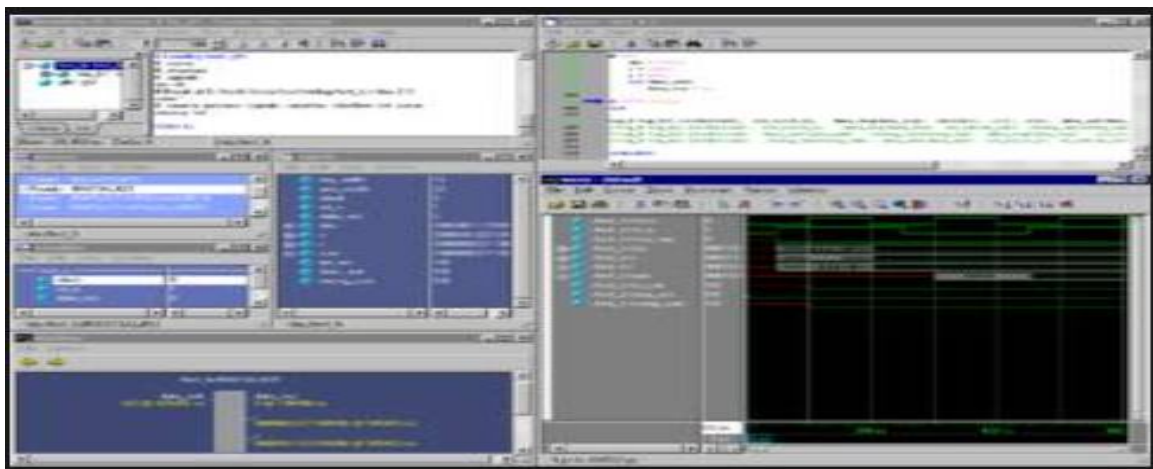
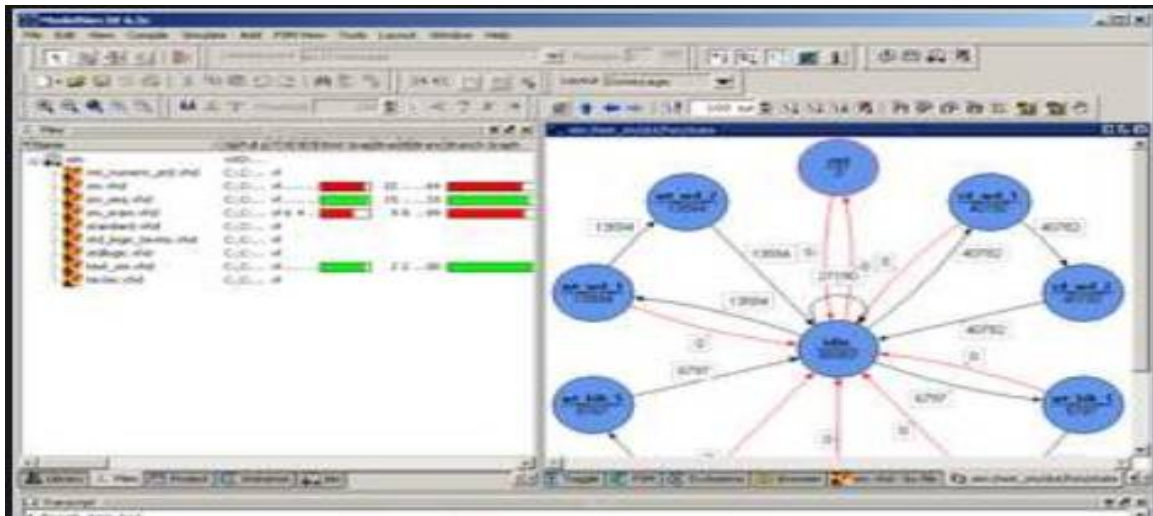
#### RESULT ANALYSIS

##### A. SYNTHESIS IS DONE IN XILINX ISE 14.4

###### TOOL

Family	Spartan3e
Part	xc3s500e
Package	fg320
Temp Grade	Commercial
Process	Typical
Speed Grade	-4
Characterization	PRODUCTION,v1.2,06-23-09

MODEL SIM



XILINX SIMULATION

V. CONCLUSION

With an ever-expanding application pool, today's many-core architectures are expected to handle a high diversity of on-chip traffic patterns. The Hammer cache coherence protocol is one such system-on-chip application that stresses on chip networks with heavy multicast injections. In this paper, we presented a multicast-aware WiNoC architecture that can efficiently handle multicast-heavy cache coherence communications. Incorporated with a congestion-aware mul-ticast routing and NC, WiNoC eliminates the initial and intermediate queuing latencies seen in conventional wireline mesh NoCs. Moreover, using wireless shortcuts, the WiNoC achieves significant reductions in network latencies leading to improved system performances.

**International Journal of Innovative Research in Science, Engineering and Technology**

*An ISO 3297: 2007 Certified Organization*

*Volume 7, Special Issue 1, March 2018*

**6<sup>th</sup> National Conference on Frontiers in Communication and Signal Processing Systems (NCFCSPS '18)**

**13<sup>th</sup>-14<sup>th</sup> March 2018**

**Organized by**

**Department of ECE, Adhiyamaan College of Engineering, Hosur, Tamilnadu, India**

**REFERENCES**

- [1] A. Karkar, N. Dahir, R. Al-Dujaily, K. Tong, T. Mak, and A. Yakovlev, "Hybrid wire-surface wave architecture for one-to-many communication in networks-on-chip," in Proc. IEEE Design Autom. Test Eur. Conf. Exhibit. (DATE), Mar. 2014, pp. 1–4.
- [2] D. Vainbrand and R. Ginosar, "Network-on-chip architectures for neural networks," in Proc. 4th ACM/IEEE Int. Symp. Netw.-Chip (NOCS), May 2010, pp. 135–144.
- [3] J.-Y. Kim, J. Park, S. Lee, M. Kim, J. Oh, and H.-J. Yoo, "A 118.4 GB/s multi-casting network-on-chip with hierarchical star-ring combined topology for real-time object recognition," IEEE J. Solid-State Circuits, vol. 45, no. 7, pp. 1399–1409, Jul. 2010.
- [4] Y. Xue and P. Bogdan, "User cooperation network coding approach for NoC performance improvement," in Proc. 9th ACM Int. Symp. Netw.-Chips, 2015, Art. No. 17.
- [5] Neuromorphic Computing: From Materials to Systems Architecture, Report of a Roundtable Convened to Consider Neuromorphic Computing Basic Research Needs, U.S. Dept. Energy, Gaithersburg, MD, USA, and Oct. 2015.
- [6] N. E. Jerger, L.-S. Peh, and M. Lipasti, "Virtual circuit tree multicasting: A case for on-chip hardware multicast support," in Proc. 35th Int. Symp. Comput. Archit. (ISCA), Jun. 2008, pp. 229–240.
- [7] A. Ros, M. E. Acacio, and J. M. García, "Dealing with traffic-area tradeoff in direct coherence protocols for many-core CMPs," in Proc. Int. Workshop Adv. Parallel Process. Technol., Berlin, Germany: Springer, Aug. 2009, pp. 11–27.