

Guest Editorial: IEEE TC Special Issue On Communications for Many-core Processors and Accelerators

Zhonghai Lu^{ID}, *Senior Member, IEEE*

COMMUNICATIONS in various forms have become increasingly important with the advent of diverse computing platforms, such as many-core CPUs, GPUs, FPGAs, machine learning accelerators, and other domain-specific processors. Meanwhile, data-intensive workloads pose greater challenges to both on-chip and off-chip communications, whereas emerging technologies offer new opportunities for interconnection networks. These trends on architecture, application and technology require innovative communication designs for the next generation of computing systems.

This special issue of IEEE TRANSACTIONS ON COMPUTERS (TC) explores academic and industrial research on all topics related to the communication issues in general-purpose and domain-specific processors. Following an open call for papers, we received 32 submissions from authors in over 10 countries on a broad range of topics related to communications in many-core processors and accelerators. A review committee was formed by inviting top experts worldwide in the field to conduct rigorous professional reviews. Nearly 100 reviews were sought in the first round based on the specific topics of the submissions and reviewer expertise. Another 58 reviews were performed in the second round to evaluate the revised submissions. Eventually, 10 high quality manuscripts have been selected for this special issue. To provide a more informative overview of the special issue to readers, we provide below a summary of the contribution of each of the papers.

The article “S-SMART++: A Low-Latency NoC Leveraging Speculative Bypass Requests” by Pérez *et al.* proposes speculative SMART++, leveraging previous SMART (Single-Cycle Multi-hop Asynchronous Repeated Traversal) techniques. It is a speculative mechanism that can acquire the multi-hop bypass in advance to decrease the base latency and the latency sensitivity to the maximum number of by-pass hops per cycle. The evaluations show its gains in reducing logic resources and dynamic power.

The article “HAM: Hotspot-Aware Manager for Improving Communications With 3D-Stacked Memory” by Wang *et al.* proposes a Hotspot-Aware Manager (HAM) design for 3D-stacked memory devices. It optimizes memory access streams via request aggregation, hotspot detection, and in-memory prefetching. As a result of optimized communication between host processor and 3D stacked memory and within the 3D-stacked memory itself, HAM improves the overall memory system performance by 21.81% and reduces the power consumption by 35.07% on a set of representative data-intensive benchmarks.

The article “OPTWEB: A Lightweight Fully Connected Inter-FPGA Network for Efficient Collectives” by Mizutani *et al.* proposes an optical all-to-all network for coordinating collective operations in multiple FPGA systems, resulting in greatly simplifying the network stack. Based on real hardware implementation, it demonstrates very small start-up latency, high bandwidth, and lower cost compared to a conventional packet-switched network.

The article “A New Optoelectronic Hybrid Network Based on Scheduling Optimization of Optical Links” by Shao *et al.* proposes a software-defined network accelerator (sDNA) that leverages optical links to offload traffic and dynamically balance resources. Optical link candidates are dynamically chosen to maximize traffic offloading revenue and then selected for adaptive routing based on flexible objectives. Evaluation shows that sDNA maintains a throughput of more than 80% bandwidth while reducing communication delay.

The article “PIT: Processing-In-Transmission With Fine-Grained Data Manipulation Networks” by Xia *et al.* proposes a novel computation paradigm, Processing-In-Transmission, in which data can be manipulated during transmission. Based on this concept, the authors design a processing-in-transmission architecture (PITA) that performs data sorting, reordering, and multicast. Evaluation shows that PITA significantly improves performance in matrix inversion tasks while achieving nearly 100% PE efficiency with sparse CNN computations.

The article “Opportunistic Caching in NoC: Exploring Ways to Reduce Miss Penalty” by Das *et al.* proposes a method to opportunistically cache data in Network-on-Chip buffers. These buffers consist of router input buffers, which are underutilized except during high network congestion, and trace buffers, which are traditionally unused following post-silicon validation. Evaluations show that the proposed architecture can significantly reduce cache miss penalties and provide an average of 14% improvement in system performance.

The article “Computing En-Route for Near-Data Processing” by Huang *et al.* proposes an in-network computing architecture for near-data processing. Computation in the memory network is supported by a novel mechanism to dynamically construct topology-oblivious routing trees, thereby allowing reduction operations to exploit massive memory-level parallelism. Evaluations show that the proposed architecture achieves an average of 60% performance improvement while reducing the energy-delay product by 80% compared to the prior state-of-the-art.

The article “Efficient Pipelined Execution of CNNs Based on In-Memory Computing and Graph Homomorphism Verification” by Dazzi *et al.* presents an efficient communication fabric for in-memory computing cores based on a graph topology well-suited for convolutional neural networks (CNNs). It facilitates the pipelined execution of CNNs while showing the existence of a homomorphism between the graph representations of CNNs and the communication fabric. Extensive evaluations including a mapping case study expose the advantages and potential of this novel approach.

The article “Optimizing Vertex Pressure Dynamic Graph Partitioning in Many-Core Systems” by McCrabb and Bertacco proposes a novel heuristic to repartition dynamic graphs in processing-near-memory architectures. Repartitioning decisions made by the heuristic are based on past vertex-update messages, rather than costly queries, so introduce no performance overhead while significantly reducing costly inter-vault messages. Consequently, evaluations show that the heuristic provides a 1.9x speedup, on average, over several graph algorithms and datasets when compared against a baseline without graph repartitioning.

The article “Plasticity-on-Chip Design: Exploiting Self-Similarity for Data Communications” by Xiao *et al.* presents a novel Plasticity-on-Chip (PoC) design flow to mine and exploit self-similarity of parallel programs. The methodology includes communication modeling, parallelization discovery, graph feature (fractal dimension and degree distribution) extraction, and cluster mapping & irregular NoC synthesis. In particular, the PoC framework is general enough to be applied to both heterogeneous platforms and reconfigurable systems.

We would like to thank all the authors for their submissions. We also would like to express our sincere gratitude to all reviewers in the expert review committee for their dedication and hard work in providing high quality reviews in time. Special thanks to the co-Guest Editor on the TC side, Associate Editor Prof. Lizhong Chen, for his significant efforts in various stages of this special issue, and the Corresponding Topical Editor, Associate Editor Prof. Cristina Silvano, for handling submissions that we have conflict of interests with. Finally, we would like to thank the Editor-in-Chief Prof. Ahmed Louri, the Associate Editor Prof Avinash Karanth, and the staff of the IEEE Transactions on Computers for their professional support and guidance throughout the entire special issue process.



Zhonghai Lu (Senior Member, IEEE) received the BS degree in radio and electronics from Beijing Normal University, Beijing, China, in 1989, and the MS degree in system-on-chip design and the PhD degree in electronic and computer system design from the KTH Royal Institute of Technology, Stockholm, Sweden, in 2002 and 2007, respectively. He was an engineer in the area of electronic and embedded systems from 1989 to 2000. He is currently a professor with the School of Electrical Engineering and Computer Science, KTH Royal Institute of Technology. He has authored about 200 peer-reviewed articles. His current research interests include interconnection network, computer architecture, design automation, and real-time systems.