



HAL
open science

Optics for Disaggregating Data Centers and Disintegrating Computing

Nikos Terzenidis, Miltiadis Moralis-Pegios, Stelios Pitris, Charoula Mitsolidou,
George Mourgias-Alexandris, Apostolis Tsakyridis, Christos Vagionas,
Konstantinos Vyrsokinos, Theoni Alexoudi, Nikos Pleros

► **To cite this version:**

Nikos Terzenidis, Miltiadis Moralis-Pegios, Stelios Pitris, Charoula Mitsolidou, George Mourgias-Alexandris, et al.. Optics for Disaggregating Data Centers and Disintegrating Computing. 23th International IFIP Conference on Optical Network Design and Modeling (ONDM), May 2019, Athens, Greece. pp.274-285, 10.1007/978-3-030-38085-4_24 . hal-03200687

HAL Id: hal-03200687

<https://inria.hal.science/hal-03200687v1>

Submitted on 16 Apr 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Optics for Disaggregating Data Centers and Disintegrating Computing

Nikos Terzenidis^{1,2}, Miltiadis Moralis-Pegios^{1,2}, Stelios Pitris^{1,2}, Charoula Mitsolidou^{1,2}, George Mourgias-Alexandris^{1,2}, Apostolis Tsakyridis^{1,2}, Christos Vagionas^{1,2}, Konstantinos Vyrsokinos³, Theoni Alexoudi^{1,2} and Nikos Pleros^{1,2}

¹ Dept. of Informatics, Aristotle University of Thessaloniki, Thessaloniki, Greece

² Center for Interdisciplinary Research and Innovation, Balkan Center, Thessaloniki, Greece

³ Department of Physics, Aristotle University of Thessaloniki, Thessaloniki, Greece

npleros@csd.auth.gr

Abstract. We present a review of photonic Network-on-Chip (pNoC) architectures and experimental demonstrations, concluding to the main obstacles that still impede the materialization of these concepts. We also propose the employment of optics in chip-to-chip (C2C) computing architectures rather than on-chip layouts towards reaping their benefits while avoiding technology limitations on the way to many-core set-ups. We identify multsocket boards as the most prominent application area and present recent advances in optically enabled multsocket boards, revealing successful 40Gb/s transceiver and routing capabilities via integrated photonics. These results indicate the potential to bring energy consumption down by more than 60% compared to current Quick-Path Interconnect (QPI) protocol, while turning multsocket architectures into a single-hop low-latency setup for even more than 4 interconnected sockets, which form currently the electronic baseline.

Keywords: computing architectures, disintegrated computing, Network-on-Chip, silicon photonics.

1 Introduction

Workload parallelism and inter-core cooperation are forcing computing to rely at a constantly growing degree on data movement. That led to an upgraded role for the on-chip and off-chip communication infrastructures that support low-power and high-bandwidth interconnect technologies. This came almost simultaneously with the revolutionary advances triggered in the field of optical interconnects [1] and silicon photonics [2]. The last 20 years, optical interconnects were transformed to a mature technology for rack-to-rack [3] and board-to-board communications [4], supporting also the emerging concepts of disaggregated computing [5] and leaf-spine Data Center architectures [6]. However, the on-chip and chip-to-chip photonic technologies are still far away from commercialization, despite the fact that various photonic Network-on-Chip (NoC) architectural concepts have already proposed [7].

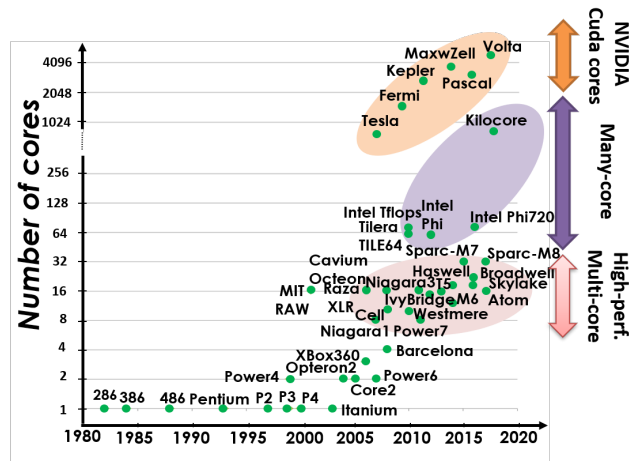


Fig. 1. Evolution from single- to many-core computing architectures.

In parallel, computing has also experienced some radical advances by turning from simple dual- and quad-core layouts into a highly heterogeneous environment both at chip- and system-level. As shown in Fig. 1, General-Purpose Graphic Processing Units (GP-GPUs) [8] can host more than 4000 CUDA cores on the same die, offering, however, only a 2 Gflop per core processing power. Processing power per core increases in manycore architectures, where up to 1000 cores can be employed [9]. However, when high-performance cores are required as in the case of Chip Multiprocessor (CMP) configurations [10] only a number of up to 32 cores can fit on the same die. The ideal scenario towards boosting processing power would of course imply a die that employs as many cores as a GPU does, but with core capabilities similar to the high-performance cores available in CMPs.

The number of high-performance cores performing as a single computational entity can scale to higher values only through multi-socket designs with 4 or maximum 8 interconnected sockets. The most recent top-class Intel Xeon 8-socket board yields a total number of up to 224 cores [11], requiring, of course, the use of high-bandwidth off-chip inter-socket interconnects. Going one step beyond the multisocket scheme, disintegration of processor dies has been coined in the recent years as a way to form macrochips that will synergize a high amount of high-performance cores, usually exploiting optical inter-die links [12]. This versatile environment at chip-scale suggests a diverse set of requirements that has to be met by optics, depending on the application. However, it creates also a new opportunity to rethink the role of optics in on- and off-chip computing, building upon the proven capabilities of optical hardware towards strengthening the compute architecture/technology co-design perspective.

In this paper, we attempt to investigate the new perspectives for optics in computing, reviewing the high-priority challenges faced currently by the computing industry and evaluating the credentials of state-of-the-art photonics to address them successfully. We provide a review of the work on photonic NoCs, highlighting the bottlenecks towards their materialization. Building on the state-of-art pNoC implementations [13-

33], we conclude to a solid case for employing integrated photonics in inter-chip multi-socket and disintegrated layouts rather than in Network-on-Chip (NoC) implementations, proposing at the same time a flat-topology chip-to-chip multi-socket interconnect technology. We demonstrate experimental results for 40 Gb/s multi-socket boards (MSBs) operation, showing the potential to scale to >8-socket designs boosting the number of directly interconnected high-performance cores. Combined with the Hipo λ os Optical Packet Switch (OPS) that has been recently shown to support sub- μ sec latencies [34], an optically-enabled rack-scale 256-socket disaggregated setting using a number of 32 interconnected optical 8-socket MSBs, could be implemented, forming in this way a powerful disaggregated rack-scale computing scheme.

The paper is organized as follows: Section II outlines the main challenges faced today in the computing landscape, providing also an overview of the research on pNoC architectures, concluding to their main limitations. Section III argues for the employment of optics in MSBs and provides experimental results on a 40Gb/s flat-topology 8-node chip-to-chip (C2C) layout, using O-band integrated photonic transceiver and routing circuitry. Finally, Section IV concludes the paper.

2 Overview of the PNoC architectures

In order to define and refine the role of optics in the current computing landscape, it is critical to identify the main challenges currently experienced by the computing industry along the complete hierarchy from on-chip through multi-socket chip-to-chip computational modules. Fig. 2 provides an illustrative overview of the main bandwidth, latency and energy needs for different on-chip and off-chip interconnect layers and data transfer operations in a $20 \times 20 \text{mm}^2$ processor chip fabricated by a 28nm Integrated Circuit (IC) CMOS technology.

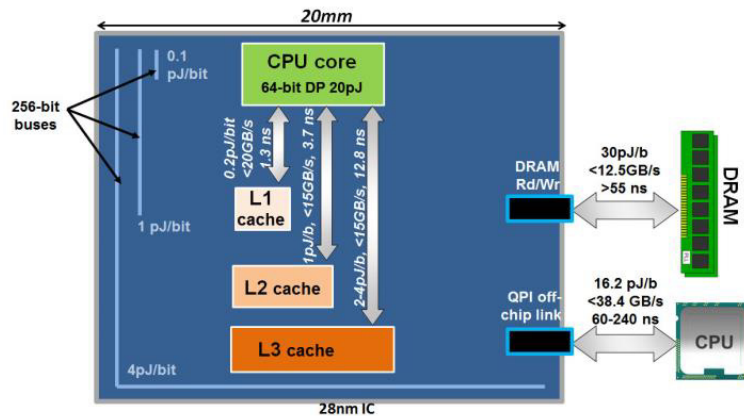


Fig. 2. Energy, bandwidth and latency requirements at different on-chip and off-chip communication needs. The size of every cache memory is bigger for larger capacity caches and their distance from the core is higher as the cache hierarchy increases.

A digital processing operation performed by the core consumes only 20pJ/bit, but sending data across the chip requires 0.1pJ/bit for a 1mm long electrical link, 1pJ/bit for a 10mm link and goes up to 4pJ/bit for a link length of 40mm. When going off-chip in order to access DRAM, a high amount of 30pJ/bit is consumed, while a chip-to-chip interconnect link like QPI requires 16.2pJ/bit. Accessing L1 cache requires 0.2pJ/bit, while L2 and L3 access requires 1 and 2-4pJ/bit, respectively. Memory bandwidth reduces with increasing memory hierarchy, with L1 memory bandwidth approaching 20GB/sec and gradually decreasing when going to L2 and L3 access until an upper limit of 12.5GB/sec in the case of DRAM access. Latency follows the inverse path, starting from a high >55nsec value when fetching from DRAM and gradually reducing with increased memory hierarchy, with L1 access latency being around 1.3nsec. Having this overview, the main challenges today are formed around:

i) Interconnect energy consumption: A modern CPU consumes around 1.7nJ per floating-point operation [35-36], being 85x higher than the 20pJ per floating point required for reaching the Exascale milestone within the gross 20MW power envelope. Current architectures rely to a large degree on data movement, with electronic interconnects forming the main energy consuming factor in both on- and off-die setups [36]. With the energy of a reasonable standard-cell-based, double-precision fused-multiply add (DFMA) being only ~20 pJ, it clearly reveals that fetching operands is much more energy-consuming than computing on them [35-36].

ii) Memory bandwidth at an affordable energy envelope: The turn of computing into strongly heterogeneous and parallel settings have transformed memory throughput into a key factor for increasing processing power [35], with the most efficient way for improvement still being the use of wider memory buses and hierarchical caching. However, the highest memory bandwidth per core in modern multicore processors can hardly reach 20 GB/sec [37], with L1 cache latency values still being >1nsec.

iii) Die area physical constraints: The need to avoid the latency and energy burden of DRAM access has enforced a rich on-chip L1, L2 and L3 cache hierarchy that typically occupies >40% of the chip real-estate [38], suggesting that almost half of the die area is devoted to memory and interconnects instead of processing functions.

iv) Cache coherency-induced multi- and broadcasting traffic patterns: The need for cache coherency at intra-chip multi- and manycore setups, as well as at inter-chip multisocket systems, yields communication patterns with strong multi- and broadcast characteristics, that have to be satisfied at a low- latency low-energy profile by the interconnect and network-on-chip infrastructure. Multibus ring topologies form a widely adopted multicast-enabling NoC architecture in current modern multi-core processors [39], but still the cache coherency control messages may often account for more than 30% of the total available bandwidth, which may reach even 65% in multi-socket settings [40].

The first attempts to exploit photonics for overcoming the on-chip bandwidth, energy and latency bottlenecks mainly inspired by the rapidly growing field of silicon photonics [2]. A number of breakthrough computing architectures relying on pNoC were demonstrated, proposing and utilizing novel silicon photonic transceiver and switching schemes. The pioneering work on photonic Torus [7] was followed by performance and energy advances in pNoC-enabled many-core designs, addressing even

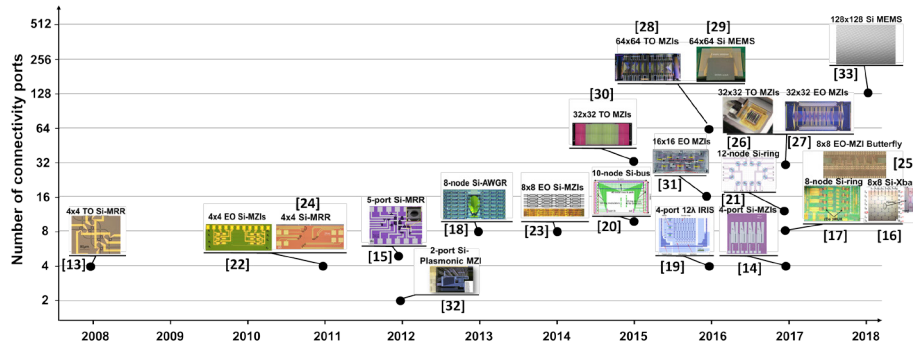


Fig. 3. Evolution of photonic Network-on-Chip and on-chip photonic switches.

cache-coherency needs [41]. All this shaped a promising roadmap for the many-core computing architectures [7], [42-47]. At the same time, it revealed the requirements to be met by silicon photonics towards materializing their on-chip employment in practical NoC layouts: transceiver line-rates between 1-40 Gb/s and optoelectronic conversion energies between a few tens to a few hundreds of fJ/bit were considered in the vast majority of pNoC schemes [7], [42-47]. Driven by these efforts, photonic integration technology achieved the performance metrics required by pNoC architectures with silicon photonic modulators and SiGe Photo-diodes (PDs) operating at data rates up to 56Gb/s exhibiting an energy efficiency less than a few tens of fJ/bit [48].

Fig. 3 summarizes the most important pNoC and on-chip switches up to now [13-33]. Silicon switches have witnessed a remarkable progress yielding high-port connectivity arrangements with a variety of underlying physical mechanisms like the thermo-optic (TO), electro-optic (EO) and opto-mechanical effects [49], allowing for 32×32 EO Mach-Zehnder Interferometric (MZI)-based layouts [27], 64×64 TO MZI designs [28] and up to 128×128 Microelectromechanical switches (MEMS) [33].

All these demonstrations indicate that integrated photonics can now indeed offer the line-rate, energy, footprint and connectivity credentials required by pNoC-enabled manycore computing architectures. However, the realization of a manycore machine that employs a pNoC layer seems to be still an elusive target, with the main reason being easily revealed when inspecting the non-performance-related co-integration and integration level details of a pNoC-enabled computational setting. Manycore architectures necessitate the on-die integration of a few thousands of photonic structures [7], residing either on 3D integration schemes [50] or on monolithically co-integrated electronic and photonic structures, with transistors and optics being almost at the same layer [51]. However, 3D integration has still not managed to fulfil the great expectations that were raised and is still struggling to overcome a number of significant challenges [52]. On the other hand, monolithic integration has recently accomplished some staggering achievements reporting on real workload execution over an opto-electronic die with optical core-memory interconnection [53]. Nevertheless, this technology has still a long-way to go until reaching the complexity and functionality level required by a many-core pNoC design.

With almost the complete Photonic Integrated Circuit (PIC) technology toolkit being today available as discrete photonic chips, computing can reap the benefits of optics by employing photonics for off-die communication in i) multi-socket and ii) disintegrated layouts. Both schemes can yield a high number of directly interconnected high-performance cores, unleashing solutions that cannot be met by electronics. At the same time, this approach is fully inline with the 2.5D integration scheme that employs discrete photonic and electronic chips on the same silicon interposer and has made tremendous progress in the recent years [54]. To this end, the employment of off-die communications via discrete photonic chips can form a viable near-term roadmap for the exploitation of photons in computational settings.

3 Optics for multi-socket boards

MSB systems rely currently on electrically interconnected sockets and can be classified in two categories:

i) “glueless” configurations, where point-to-point (P2P) interconnects like Intel’s QPI [55] can offer high-speed, low-latency, any-to-any C2C communication for a number of 4 or 8 sockets. A 4-socket setup can yield a cache-coherent layout with directly interconnected sockets and latency values that range between 60-240nsec. Scaling to 8-socket designs can only be met through dual-hop links, degrading latency performance but still comprising a very powerful cache-coherent computational setting: Intel’s Xeon E7-8800 v4 was the first processor supporting 8-socket configurations and was by that time advertised as being suitable to “dominate the world” [56]. Fig. 4(a) depicts a 4-socket (4S) and 8-socket (8S) layout, respectively, along with their respective interconnects. A typical interconnect like Intel’s QPI operates at a 9.6 Gb/s line-rate and consumes 16.2 pJ/bit, while the total bandwidth communicated by every socket towards all three possible directions is 38.4 GB/s, i.e. 307.2 Gb/s [57].

ii) “glued” configurations, where scaling beyond 8-socket layouts is accomplished by exploiting active switch-based setups, such as PCI-Express switches, in order to interconnect multiple 4- or 8-socket QPI “islands”[57].

With latency and bandwidth comprising the main performance criteria in releasing powerful MSB configurations, “glueless” layouts offer a clear latency-advantage over the “glued” counterparts avoiding by default the use of any intermediate switch. Photonics can have a critical role in transforming “glued” into “glueless” architectures even when the number of interconnected sockets is higher than 8, enabling single-hop configurations, with Fig. 4(b) illustrating how the basic flat-topology can be accomplished for the case of an 8-Socket layout. This has been initially conceived and proposed by UC Davis in their pioneering work on Flat-Topology computing architectures [58] via Arrayed Waveguide Grating Router (AWGR) interconnects, utilizing low-latency, non-blocking and all-to-all optical connectivity credentials enabled by their cyclic-routing wavelength properties. UC Davis demonstrated via gem5 simulations the significant execution time and energy savings accomplished over the electronic baseline [58], revealing also additional benefits when employing bit-parallel

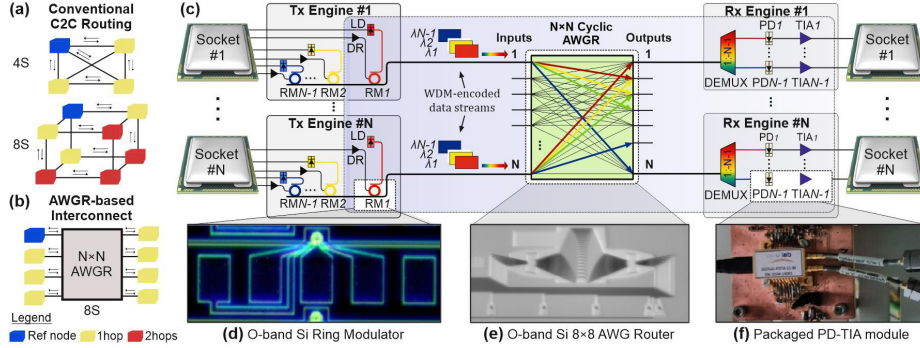


Fig. 4. (a) C2C routing in current electronic 4S and 8S MSBs, (b) Flat-topology 8S layout using AWGR-based routing, (c) proposed $N \times N$ AWGR-based optical C2C interconnect for MSB connectivity. Photonic integrated circuits employed as the basic building blocks in the 40Gb/s experimental demonstration: (d) Ring Modulator, (e) 8×8 cyclic-frequency AWGR and (f) PD-TIA module. (blue-highlighted areas: part of the architecture demonstrated experimentally, white-highlighted areas: basic building blocks used for the demonstration).

transmission and flexible bandwidth-allocation techniques. Experimental demonstrations of AWGR-based interconnection for compute node architectures were, however, constrained so far in the C-band regime, limiting their compatibility with electro-optic Printed Circuit Board (PCB) technology that typically offers a low waveguide loss figure at the O-band [59]. As such, AWGR-based experimental compute node interconnect findings were reported so far only in pNoC architectural approaches, using a rather small line-rate operation of 0.3 Gb/s [60].

The European H2020 project ICT-STREAMS is currently attempting to deploy the necessary silicon photonic and electro-optical PCB technology toolkit for realizing the AWGR-based MSB interconnect benefits in the O-band and at data rates up to 50Gb/s [61]. It aims to exploit wavelength division multiplexing (WDM) Silicon photonics transceiver technology at the chip edge as the socket interface and a board-pluggable O-band silicon-based AWGR as the passive routing element, as shown in a generic N -socket architecture depicted in Fig. 4(c). Each socket is electrically connected to a WDM-enabled Tx optical engine equipped with $N-1$ laser diodes (LD), each one operating at a different wavelength. Every LD feeds a different Ring Modulator (RM) to imprint the electrical data sent from the socket to each one of the $N-1$ wavelengths, so that the Tx engine comprises finally $N-1$ RMs along with their respective RM drivers (DR). All RMs are implemented on the same optical bus to produce the WDM-encoded data stream of each socket. The data generated by each socket enters the input port of the AWGR and is forwarded to the respective destination output that is dictated by the carrier wavelength and the cyclic-frequency routing properties of the AWGR [58]. In this way, every socket can forward data to any of the remaining 7 sockets by simply modulating its electrical data onto a different wavelength via the respective RM, allowing direct single-hop communication between all sockets through passive routing. At every Rx engine, the incoming WDM-encoded data stream gets demultiplexed with a $1:(N-1)$ optical demultiplexer (DEMUX), so

that every wavelength is received by a PD. Each PD is connected to a transimpedance amplifier (TIA) that provides the socket with the respective electrical signaling.

The AWGR-based interconnect scheme requires a higher number of transceivers compared to any intermediate switch solution, but this is exactly the feature that allows to combine WDM with AWGR's cyclic frequency characteristics towards enabling single-hop communication and retaining the lowest possible latency. Utilizing an 8×8 AWGR, the optically-enabled MSB can allow single-hop all-to-all interconnection between 8 sockets, while scaling the AWGR to 16×16 layouts can yield single-hop communication even between 16 sockets, effectively turning current "glued" into "glueless" designs. The ICT-STREAMS on-board MSB aims to incorporate 50GHz single-mode O-band electro-optical PCBs [62], relying on the adiabatic coupling approach between silicon and polymer waveguides [63] for low-loss interfacing of the Silicon-Photonics (Si-Pho) transceiver and AWGR chips with the EO-PCB.

Next, the first 40Gb/s experimental results of demonstration with the fiber-interconnected integrated photonic building blocks is presented, extending the recently presented operation of the 8-socket architecture at 25 Gb/s [64]. The main integrated transmitter, receiver and routing building blocks that were used, comprise three discrete chips, i.e. a Si-based RM [48], a Si-based 8×8 AWGR routing platform [65] and a co-packaged PD-TIA [66], which are depicted in Fig. 4(d), (e) and (f), respectively. The silicon O-band carrier-depletion micro-ring modulator is an all-pass ring resonator fabricated on imec's active platform with demonstrated 50 Gb/s modulation capabilities [48]. The RM can be combined with a recently developed low-power driver [67], leading to an energy efficiency of 1 pJ/bit at 40 Gb/s. For the routing platform, the demonstration relied on an O-band integrated silicon photonic 8×8 AWGR device [65] with 10 nm-channel spacing, a maximum channel loss non-uniformity of 3.5 dB and a channel crosstalk of 11 dB. Finally, the Rx engine employed a co-packaged uni-traveling InGaAs-InP PIN photodiode (PD) connected with a low-power TIA implemented in $0.13 \mu\text{m}$ SiGe BiCMOS [66]. The PD-TIA energy efficiency for operation at 40 Gb/s is 3.95 pJ/bit.

The energy efficiency of the proposed 40 Gb/s chip-to-chip (C2C) photonic link is estimated at 5.95 pJ/bit, assuming a 10% wall-plug efficiency for the external laser. This indicates that the proposed architecture has the credentials to lead to 63.3% reduction in energy compared to the 16.2 pJ/bit link energy efficiency of Intel QPI [57]. Fig. 5 (a)-(h) show the eye diagrams of the signal at the 8 outputs of the AWGR corresponding to the 8 routing scenarios for all possible input-output port combina-

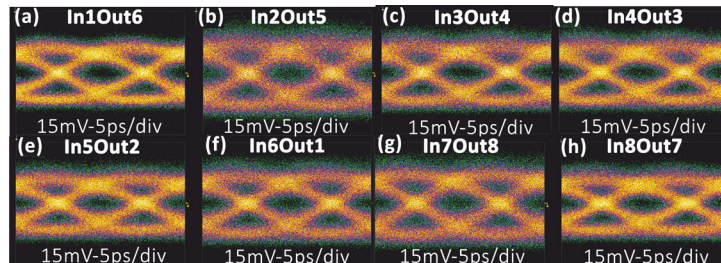


Fig. 5. Eye diagrams (a)-(h) after routing via the respective In#iOut#j I/O ports of the AWGR.

tion, indicating clear eye openings and successful routing at 40 Gb/s with ER values of 4.38 ± 0.31 dB and AM values of 2.3 ± 0.3 dB, respectively. The RM was electrically driven with a peak-to-peak voltage of 2.6 V_{pp}, while the applied reverse DC bias voltage was -2.5 V. The optical power of the CW signal injected at the RM input was 8 dBm, with the modulated data signal obtained at the RM output having an average optical power level of -6.3 dBm.

Going a step further, the proposed optically-enabled MSBs can be beneficially employed in rack-scale disaggregated systems when equipped with an additional transceiver lane for dealing with the off-board traffic and are combined with the recently demonstrated Hipo λ os high-port switch architecture [34]. Recently in [68], it was shown that rack-scale disaggregation among a 256-node system can be successfully accomplished for a variety of communication patterns with an ultra-low mean latency value of < 335 nsec for 10 Gb/s data rates. The disaggregated architecture are expected to improve drastically when scaling Hipo λ os data-rates to 40Gb/s, making this compatible with the 40Gb/s silicon photonic transmitter reported in this paper.

4 Conclusion

We reviewed the pNoC-enabled manycore architectures proposed over the last decade. After analyzing the co-integration aspects as the main limitation for the realization of pNoC-based computing, we have defined a new role for photonics in the landscape of computing related to off-die communication. We discussed how optics can yield single-hop low-latency multsocket boards for even more than 4 interconnected sockets, demonstrating experimental results for 40Gb/s C2C interconnection in a 8-node setup via integrated photonic transmitter and routing circuits. Combining 8-socket optical boards with a Hipo λ os optical packet switch shown in [34], photonics can yield a powerful 256-node compute disaggregated system with latency below the sub- μ s threshold considered for memory disaggregation environments.

Acknowledgment: This work is supported by the H2020 projects ICT-STREAMS (688172) and L3MATRIX (688544).

References

1. K. Bergman: Photonic Networks for Intra-Chip, Inter-Chip, and Box-to-Box Interconnects in High Performance Computing. In: Eur. Conf. on Optical Comm. (ECOC), Cannes, France, 2006.
2. M. Lipson: Guiding, modulating, and emitting light on Silicon-challenges and opportunities. *J. Lightw. Techn.* 23 (12), 4222-4238 (2005).
3. Intel SiP 100G PSM4 Optical Tx, <https://www.intel.com/content/www/us/en/architecture-and-technology/silicon-photonics/optical-transceiver-100g-psm4-qsf28-brief.html>, last accessed 2019/04/12.
4. Luxtera 2x100G-PSM4 OptoPHY Product Family, <http://www.luxtera.com/embedded-optics/>, last accessed 2019/04/12.

5. N G. Zervas, H. Yuan, A. Saljoghei, Q. Chen and V. Mishra: Optically disaggregated data centers with minimal remote memory latency: Technologies, architectures, and resource allocation. *J. Opt. Comm. and Netw.*, 10 (2), A270-A285 (2018)
6. M. Bielski, et al.: dReDBox: Materializing a Full-stack Rack-scale System Prototype of a Next-Generation Disaggregated Datacenter. In: 2018 Design, Automation & Test Conference & Exhibition (DATE), 2018.
7. A. Shacham, K. Bergman and L. Carloni: Photonic Networks-on-Chip for Future Generations of Chip Multiprocessors. *Trans. on Computers* 57 (9), 1246-1260 (2008).
8. J. Kider, NVIDIA Fermi architecture
http://www.seas.upenn.edu/~cis565/Lectures2011/Lecture16_Fermi.pdf, last accessed 2019/04/12.
9. B. Bohnenstiehl, et. al.: KiloCore: A 32-nm 1000-Processor Computational Array", *IEEE Journal of Solid-State Circuits* 52 (4), 891-902 (2017).
10. Intel Xeon Platinum 8180 Processor. <https://ark.intel.com/products/120496>, last accessed 2019/04/12.
11. Supermicro Super Server 7089P-TR4T.
www.supermicro.com/products/system/7U/7089/SYS-7089P-TR4T.cfm, last accessed 2019/04/12.
12. K. Raj, et al.: "Macrochip" computer systems enabled by silicon photonic interconnects. In: *Optoelectronic Interconnects and Component Integration IX*, 2010.
13. N. Sherwood-Droz, H. Wang, L. Chen, B. Lee, A. Biberman, K. Bergman and M. Lipson-Optical 4x4 hitless silicon router for optical networks-on-chip (NoC). *Optics Express* 16(20), 15915 (2008).
14. H. Jia, et. al.: Four-port Optical Switch for Fat-tree Photonic Network-on-Chip. *Journal of Lightw. Techn.* 35(15), 3237–3241 (2017).
15. Lin Yang, et. al.: Optical routers with ultra-low power consumption for photonic networks-on-chip. In *Proc: Conf. on Lasers and Electro-Optics (CLEO)*, San Jose, CA, 2012.
16. G. Fan, R. Orobtcouk, B. Han, Y. Li and H. Li: 8 x 8 wavelength router of optical network on chip. *Optics Express* 25 (20), 23677 (2017).
17. C. Zhang, S. Zhang, J. Peters and J. Bowers: 8 × 8 × 40 Gbps fully integrated silicon photonic network on chip. *Optica* 3(7), 785 (2016).
18. R. Yu, et. al.: A scalable silicon photonic chip-scale optical switch for high performance computing systems. *Optics Express* 21 (26), 32655 (2013).
19. F. Testa, et. al.: Design and Implementation of an Integrated Reconfigurable Silicon Photonics Switch Matrix in IRIS Project. *J. of Sel. Topics in Quant. Electr.*, 22 (6), 155-168 (2016).
20. P. Dong, et. al.: Reconfigurable 100 Gb/s Silicon Photonic Network-on-Chip. In *Proc: Optical Fiber Communication Conference (OFC)*, 2014.
21. F. Gambini, et. al.: Experimental demonstration of a 24-port packaged multi-microring network-on-chip in silicon photonic platform", *Optics Express* 25(18), 22004 (2017).
22. M. Yang, et. al.: Non-Blocking 4x4 Electro-Optic Silicon Switch for On-Chip Photonic Networks. *Optics Express*, vol. 19, no. 1, p. 47, 2010.
23. B. Lee, et. al.: Monolithic Silicon Integration of Scaled Photonic Switch Fabrics, CMOS Logic, and Device Driver Circuits. *Journal of Lightwave Techn.* 32(4), 743-751(2014).
24. T. Hu, et. al. Wavelength-selective 4×4 nonblocking silicon optical router for networks-on-chip. *Optics Letters* 36(23), 4710 (2011).
25. N. Dupuis, et. al. Nanosecond-scale Mach-Zehnder-based CMOS Photonic Switch Fabrics. *Journal of Lightwave Technology*, 1-1 (2016).

26. P. Dumais, et. al.: Silicon Photonic Switch Subsystem With 900 Monolithically Integrated Calibration Photodiodes and 64-Fiber Package. *J of Lightw. Techn.* 36(2), 233-238 (2018).
27. L. Qiao, W. Tang and T. Chu: 32×32 silicon electro-optic switch with built-in monitors and balanced-status units. *Scientific Reports* 7(1), 2017.
28. L. Qiao, W. Tang and T. Chu: Ultra-large-scale silicon optical switches. In *Proc: 2016 IEEE Int. Conference on Group IV Photonics (GFP)*, Shanghai, 2016.
29. T. J. Seok: 64×64 Low-loss and broadband digital silicon photonic MEMS switches. In *Proc: European Conference on Optical Communication (ECOC)*, Valencia, 2015.
30. K. Tanizawa, K. Suzuki, M. Toyama, M. Ohtsuka, N. Yokoyama, K. Matsumaro, M. Seki, K. Koshino, et. al.: Ultra-compact 32×32 strictly-non-blocking Si-wire optical switch with fan-out LGA interposer", *Optics Express* 23(13), 17599 (2015).
31. L. Lu, et. al.: 16×16 non-blocking silicon optical switch based on electro-optic Mach-Zehnder interferometers. *Optics Express* 24(9), 9295 (2016).
32. S. Papaioannou, et. al.: Active plasmonics in WDM traffic switching applications. *Scientific Reports* 2(1), 2012.
33. K. Kwon, et. al.: 128×128 Silicon Photonic MEMS Switch with Scalable Row/Column Addressing. In *Proc: Conference on Lasers and Electro-Optics*, 2018.
34. N. Terzenidis, M. Moralis-Pegios, G. Mourgias-Alexandris, K. Vysokinos, N. Pleros: High-port low-latency optical switch architecture with optical feed-forward buffering for 256-node disaggregated data centers. *Op. Ex.* 26, 8756-8766 (2018).
35. S. Parker: The Evolution of GPU Accelerated Computing, In *Proc: Extreme Scale Computing*, IL, USA, July 29, 2013.
36. B. Dally: Challenges for Future Computing Systems. In *Proc.: HiPEAC 2015*, Amsterdam, NL, 2015.
37. S. Saini et al: Performance Evaluation of the Intel Sandy Bridge Based NASA Pleiades Using Scientific and Engineering Applications. *NAS Technical Report: NAS-2015-05*.
38. S. Borkar, A.A. Chien: The future of microprocessors. *Commun.ACM* 54(5), 67-77 (2011).
39. M. Kumashikar, S. Bendi, S. Nimmagadda, A. Deka and A. Agarwal: 14nm Broadwell Xeon® processor family: Design methodologies and optimizations. In *Proc.: 2017 IEEE Asian Solid-State Circuits Conference (A-SSCC)*, 2017.
40. Bull SAS. An efficient server architecture for the virtualization of business-critical applications. white paper 2012. https://docuri.com/download/bullion-efficient-server-architecture-for-virtualization_59c1dc51f581710b28689168_pdf, last accessed 2019/04/12.
41. G. Kurian, et. al.: ATAC. In: *International Conference on Parallel Architectures and Compilation techniques - PACT '10*, 2010.
42. C. Chen and A. Joshi: Runtime Management of Laser Power in Silicon-Photonic Multibus NoC Architecture. *J. of Sel. Top. in Quant. Electr.* 19 (2), 3700713-3700713 (2013).
43. Z. Li, A. Qouneh, M. Joshi, W. Zhang, X. Fu and T. Li: Aurora: A Cross-Layer Solution for Thermally Resilient Photonic Network-on-Chip. *Trans. on VLSI Syst.* 23 (1), 170-183 (2015).
44. S. Bahirat and S. Pasricha: METEOR. *ACM Transactions on Embedded Computing Systems* 13 (3), 1-33 (2014).
45. X. Wang, H. Gu, Y. Yang, K. Wang and Q. Hao: RPNOC: A Ring-Based Packet-Switched Optical Network-on-Chip. *Phot. Techn. Lett.* 27 (4), 423-426 (2015).
46. H. Gu, K. Chen, Y. Yang, Z. Chen and B. Zhang: MRONOC: A Low Latency and Energy Efficient on Chip Optical Interconnect Architecture. *Phot. Journal* 9 (1), 1-12 (2017).

47. S. Werner, J. Navaridas and M. Luján: Efficient sharing of optical resources in low-power optical networks-on-chip. *J. of Opt. Comm. and Netw.* 9 (5), 364-374 (2017).
48. M. Pantouvaki, et al.: Active Components for 50 Gb/s NRZ-OOK Optical Interconnects in a Silicon Photonics Platform. *J. Lightw. Techn.* 35 (4), 631-638 (2017).
49. B. Lee: Silicon Photonic Switching: Technology and Architecture. In: 2017 European Conference on Optical Communication, 2017.
50. S. J. B. Yoo, B. Guan and R. Scott: Heterogeneous 2D/3D photonic integrated microsystems. *Microsyst. & Nanoeng.* 2 (1), 2016.
51. C. Sun, et al.: Single-chip microprocessor that communicates directly using light. *Nature* 528 (7583), 534-538 (2015).
52. C. Li, et al.: Chip Scale 12-Channel 10 Gb/s Optical Transmitter and Receiver Subassemblies Based on Wet Etched Silicon Interposer. *J. Lightw. Techn.* 35 (15), 2017.
53. R. C. Sun, et. al.: Single-chip microprocessor that communicates directly using light. *Nature* 528 (7583), 534-538 (2015).
54. Xiaowu Zhang, et. al: Heterogeneous 2.5D integration on through silicon interposer. *Applied Physics Reviews* 2(2), 021308 (2015).
55. Intel. An Introduction to the Intel QuickPath Interconnect. <https://www.intel.com/content/www/us/en/io/quickpath-technology/quick-path-interconnect-introduction-paper.html>, last accessed 2019/04/12.
56. Intel, Intel® Xeon® Processor E7-8800/4800/2800 Families, <https://www.intel.com/content/www/us/en/processors/xeon/xeon-e7-8800-4800-2800-families-vol-2-datasheet.html>, last accessed 2019/04/12.
57. R. Maddox, G. Singh and R. Safranek, Weaving high performance multiprocessor fabric. Hillsboro, Intel Press, 2009.
58. P. Grani, R. Proietti, S. Cheung, S. J. B. Yoo: Flat-Topology High-Throughput Compute Node With AWGR-Based Optical-Interconnects. *J. Lightw. Techn.* 34(12), 2016.
59. A. Sugama, K. Kawaguchi, M. Nishizawa, H. Muranaka and Y. Arakawa: Development of high-density single-mode polymer waveguides with low crosstalk for chip-to-chip optical interconnection. *Optics Express* 21 (20), 24231 (2013).
60. R. Yu, et. al.: A scalable silicon photonic chip-scale optical switch for high performance computing systems. *Optics Express* 21(26), 32655 (2013).
61. G. T. Kanellos and N. Pleros: WDM mid-board optics for chip-to-chip wavelength routing interconnects in the H2020 ICT-STREAMS. In: SPIE, Febr. 2017.
62. T. Lamprecht, et. al.: EOCB-Platform for Integrated Photonic Chips Direct-on-Board Assembly within Tb/s Applications. In Proc: 68th Electronic Components and Technology Conference (ECTE), pp. 854–858 (2018).
63. R. Dangel et al.: Polymer Waveguides Enabling Scalable Low-Loss Adiabatic Optical Coupling for Silicon Photonics. *J. of Select. Topics in Quant. Electr.* 24(4),1-11(2018).
64. M. Moralis-Pegios et al: Chip-to-Chip Interconnect for 8-socket direct connectivity using 25Gb/s O-band integrated transceiver and routing circuits. In: ECOC, Rome, Italy, 2018.
65. S. Pitris, et. al.: Silicon photonic 8×8 cyclic Arrayed Waveguide Grating Router for O-band on-chip communication. *Opt. Express* 26(5), 6276-6284 (2018).
66. B. Moeneclaey et al.: A 40-Gb/s Transimpedance Amplifier for Optical Links. *IEEE Photonics Technology Letters* 27(13), 1375-1378 (2015).
67. H. Ramon et al.: Low-Power 56Gb/s NRZ Microring Modulator Driver in 28nm FDSOI CMOS. *IEEE Photonics Technology Letters* 30(5), 467-470 (2018).
68. N. Terzenidis, et. al.: Dual-layer Locality-Aware Optical Interconnection Architecture for Latency-Critical Resource Disaggregation Environments. Accepted in: *Int. Conf. on Opt. Netw. Design and Modeling (ONDM)*, May 2019.