Contents lists available at ScienceDirect

# Optical Switching and Networking

# An optically-enabled chip–multiprocessor architecture using a single-level shared optical cache memory

P. Maniotis [a,b,*], S. Gitzenis [b], L. Tassiulas [b,c], N. Pleros [a,b]

[a] Department of Informatics, Aristotle University of Thessaloniki, Thessaloniki, Greece
[b] Information Technologies Institute, Center for Research and Technology Hellas, Thessaloniki, Greece
[c] Department of Electrical and Computer Engineering, University of Thessaly, Volos, Greece

## ABSTRACT

We present an optical bus-based chip–multiprocessor architecture where the processing cores share an optical single-level cache implemented in a separate chip next to the Central-Processing-Unit (CPU) die. The interconnection system is realized through Wavelength-Division-Multiplexed optical interfaces connecting the shared cache with the cores and the Main-Memory via spatial-multiplexed waveguides. Evaluating the proposed approach, we realize system-level simulations of a wide-range parallel work-loads using Gem5. Optical cache architecture is compared against the conventional one that uses dedicated on-chip Level-1 electronic caches and a shared Level-2 cache. Results show significant Level-1 miss rate reduction of up to 96% for certain cases; on average, a performance speed-up of 19.4% or cache capacity requirements reduction of ∼63% is attained. Combined with high-bandwidth CPU-Dynamic Random Access Memory (DRAM) bus solutions based on optical interconnects, the proposed design is a promising architecture bridging the gap between high-speed optically connected CPU-DRAM schemes and high-speed optical memory technologies.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

It has been more than twenty years ago since the speed mismatch between the Central Processing Unit (CPU) and Main Memory (MM), commonly referred to as the "Memory Wall", was identified as one of the main barriers against increase in the computer performance [1]. Solutions such as deployment of large on chip cache memories, the widening of the CPU-MM buses and prefetching from the MM have been devised to ease the limited off-chip bandwidth and the high MM's response latency imposed by the constraints of the electronic technology [2]. Higher spatial multiplexing degrees through wider buses allow for more efficient and simultaneous multi-bit data transfer within a single cycle. On the other hand, trading bandwidth for reduced average delay and buffering data close to CPU in anticipation of future requests through prefetching reduced on average the access delay stalling processing. However, the introduction of modern Chip Multi-Processor (CMP) configurations has further aggravated the bottleneck between the CPU and MM and led to larger two- or even three-level cache memory hierarchies that take up almost 40% of the total chip energy consumption [3] and more that 40% of chip real estate [4–6]. Taking into account the distance- and speed-dependent energy dissipation and the low bandwidth density associated with electronic technology, novel approaches are required against the Memory Wall.

Promising emerging solutions emerge from the optical interconnects and photonic integration technology fields thanks to their proven high-speed data transfer abilities. The introduction of these technologies in the interconnection system between the memory and processing elements is expected to relieve computing from some energy-demanding and slow electronics. Focusing on the CPU–MM interconnection system, the main effort has shifted to the replacement of the electronic busses with optical wires. With the current technology, fetching 256 bits of operand from the MM module consumes more than 16 nJ [7] and requires four stages of transmission (assuming a 64-bit wide bus at an operation speed just above 1 GHz). Bringing photonics into the game by replacing the electrical busses with optical wiring solutions, either over a Silicon-on-Insulator (SOI) platform or over an Optical Printed Circuit Board (OPCB), is expected to (1) reduce energy consumption down to 1 mW/Gbps [8], (2) raise operation speed to several tens of GHz and at the same time, (3) dispense with the traditional issue of Resistance-Capacitance (RC)-induced delay of the electrical wiring. This roadmap is rapidly gaining interest with several works demonstrating the benefits of switching from

* Corresponding author at: Department of Informatics, Aristotle University of Thessaloniki, Thessaloniki, Greece.
    E-mail address: ppmaniot@csd.auth.gr (P. Maniotis).

electronic to optical CPU-MM buses [9–13] and introducing novel fully functional all-optical interfaces for Dynamic Random Access Memory (DRAM) integration [9]. However, all these enhancements cannot mitigate the need for memory caching as CMP dies will continue to struggle against finding an optimum balance in the processor, cache and interconnect circuitry considerations.

Going a step further with the emerging technologies the field of optics can also lead to novel solutions in the data buffering domain, such as cache memories. Although the lack of electric charge places photons at disadvantage when coming to storage, a variety of optical flip-flop and Random Access Memory (RAM) cell technologies have appeared for storing information. These technologies exploit the Set-Reset flip-flop architectural layout reducing at the same time the access delay. Representative all-optical flip-flop technologies include coupled Semiconductor Optical Amplifiers (SOAs) [14], III-V-on-SOI microdisk lasers [15], polarization bistable Vertical Cavity Surface Emitting Lasers (VCSELs) [16] and SOA-based Mach-Zehnder Interferometers (SOA-MZIs) [17].

Proceeding moreover to multi-bit storage, Photonic Crystal (PhC) nanocavities technology has already demonstrated more than 100-bit integrated storage capacity with significant benefits in terms of speed, energy consumption and footprint [18,19]. Extending the elementary flip-flop operation, the first optical Static RAM (SRAM) cell allowed for fully functional random access read/write operation at 5 Gbps [20]. In this configuration the cell deploys two SOA access gates and a SOA-MZI-based flip-flop, and can theoretically operate at speeds of up to 40 Gbps [21]. Next generation SRAM cells introduced significant improvements in terms of active elements and energy consumption reduction through the introduction of wavelength diversity in the incoming signals [22].

Expanding our view from single elementary memory cells to complete optical RAM architectures, [23] has highlighted the benefits of employing Wavelength Division Multiplexing (WDM)-formatted data and address fields in RAM peripheral circuits such as row [24] and column decoders [24,25]. Taking advantage of all these advances, we recently presented a complete and fully functional optical cache memory architecture that successfully performs both read and write operations directly in the optical domain [26]. In [26] we moved forward with designing an all-optical cache memory that combines all the optical subsystems, such as read/write selection modules, row and column decoders, 2D RAM banks and tag comparison circuits. Physical layer simulation scenarios carried out with the commercially available VPI Photonics simulation suite [27] indicate error-free operation at speeds up to 16 GHz for both direct [26] and 2-way associative [28] cache mapping schemes. However, its system-scale performance in CMP configurations has been so far evaluated only for the *bodytrack* and *blackscholes* benchmarks [29] from the PARSEC benchmark suite [30], providing only a limited amount of information about its advantageous architectural perspectives in true CMP systems. In order to analyze its application potential in real CMP settings, it is among the prerequisites to follow the well-known practice of validating the proposed scheme with a broad range of workloads [31].

This paper extends our previous work by presenting a detailed optical bus-based CMP architecture where all-optical Level-1 instruction (L1i) and Level-1 data (L1d) caches are shared among the processing cores and the MM. In this work we focus on the interconnection and system-level architecture performance advantages of the shared all-optical cache memory scheme using the physical-layer cache design of [26] as the basic building block. Both L1i and L1d caches are placed off-chip, sparing precious chip area from the die in favor of processing elements. On the other hand, the choice of a shared cache unit has been taken on the basis that the cycle of the optical cache memories is a fraction of the cycle of the electronic cores, making thus possible to serve

multiple concurrent core requests without stalling the execution.

The optical bus-based CMP architecture's system-scale performance is addressed for 12 parallel workloads, using the PARSEC benchmark suite [30] on top of the Gem5 simulator [32]. The simulation findings suggest that the shared optical cache architecture can improve substantially the Level-1 (L1) cache miss rate (up to 96% reduction in certain cases), and either speed-up execution (by 19.4% on average) or slash the required cache capacity (by ∼63% on average). Following the proposed CMP architecture, the connection of the cache memory modules with both the (cache-free) CMP dies and the DRAM elements can be realized completely in the optical domain, relieving processor dies from interconnection and caching modules.

The rest of this paper is organized as follows: Section 2 describes the detailed physical layer architecture of the bus-based CMP architecture, Section 3 presents the system-scale simulation results, Section 4 recapitulates the findings of this work and finally Section 5 concludes the paper.

## 2. Optical-bus-based CMP architecture with optical cache memories

Fig. 1(a) presents a typical example of a modern CMP with multi-level electronic caches and an indicative number of eight processing cores. Specifically, the standard approach is to put dedicated L1d and L1i caches at each core that run at the same speed with the core in order to maintain stall-free core operation assuming cache hits. L1d and L1i caches independently buffer the instruction and data fetch and store operations towards doubling the cache bandwidth and reducing interference between the instruction and data streams. Reducing the interference between the instruction and data streams can improve the overall system's performance and is the standard approach followed by most current CMP systems [33]. Going down through the memory hierarchy, a second unified Level-2 (L2) cache stores both instructions and data and depending on the number of cores and the target application, Level-3 (L3) caches may be eventually also employed and shared among the processing cores. Last, the MM connects to the CPU chip with a spatially multiplexed electrical bus. Although L2 and L3 are slower than L1, they are much faster to access than MM and, typically much larger in size than L1, diminishing thus the penalty of an L1 miss.

In contrast to the conventional paradigm of Fig. 1(a), we study the impact of using optical cache memory technology as single-level shared cache between the processing cores and the MM. The optical cache memory technology has relied on the hardware design initially presented in [26,28] and being capable of operating at 16 GHz line-rate speeds. In the proposed CMP architecture, the shared L1 cache is an optical cache memory technology, connected to CPU and MM via spatial multiplexed optical waveguides. The direct sharing of the cache among the cores does not necessarily stall the core operation as the optical cache operates at significant higher speeds compared to the electronic cores, managing to serve multiple concurrent requests from many cores during each electronic core cycle.

Fig. 1(b) depicts the layout of the proposed optical-bus-based CMP architecture comprising three discrete subsystems: (i) the cache-free CMP chip (8 cores are shown as in Fig. 1(a)), (ii) the optical cache chip with separate L1i and L1d caches lying next to the CMP chip, and (iii) the MM module. The interconnection system between the three subsystems consists of three optical buses with proper WDM optical interfaces at the edge of the CPU cores and the MM. Note that optical to electronic conversion is not required at the cache memory's sides as the optical cache memory operates completely in the optical domain. Section 2.1 presents all