

# The Case for Photonic Memory Disaggregation in Datacenters

John Shalf<sup>1</sup>, George Michelogiannakis<sup>1</sup>, Brian Austin<sup>1</sup>, Taylor Groves<sup>1</sup>, Manya Ghobadi<sup>2</sup>, Larry Dennison<sup>3</sup>, Tom Gray<sup>3</sup>, Yiwen Shen<sup>4</sup>, Min Yee Teh<sup>4</sup>, Madeleine Glick<sup>4</sup>, Keren Bergman<sup>4</sup>

<sup>1</sup>Lawrence Berkeley National Laboratory, 1 Cyclotron Road, Berkeley, California, 94720

<sup>2</sup>MIT CSAIL, <sup>3</sup>NVIDIA Corporation, <sup>4</sup>Columbia University

Corresponding author: jshalf@lbl.gov

**Abstract:** Datacenter and HPC workloads have diverse memory capacity requirements, but typical node architectures accommodate a superset of those requirements. High bandwidth density photonic link technologies enable efficient memory disaggregation – potentially reducing wasted memory capacity by 8-10x.

## 1. Introduction

HPC and commercial datacenters purchase large quantities of identical server nodes to simplify system management and leverage volume discounted prices. However, the workloads running on these nodes have extremely diverse requirements, motivating operators to equip each node with enough resources that can accommodate a superset of application requirements even though this approach tends to waste a lot of resources. For example, a recent memory utilization study at the NERSC HPC center shows that 50% of the workloads use < 20% of available node memory. Recent advances in high-speed, high bandwidth-density photonic link technologies, combined with broadband optical circuit switching, enable resource disaggregation across node boundaries at the rack scale, where the connectivity between typical node-level resources such as memory, storage, CPUs, and GPUs could be configured at runtime. Rack-scale resource disaggregation can provide substantial reductions in acquisition cost and energy consumption by reducing the amount of underutilized memory capacity by 8-10x.

### 1.1. Resource Disaggregation in the Datacenter

Datacenter workloads show a large diversity in their resource demands: Training algorithms for deep machine learning stress compute and interconnect elements, in-memory databases stress integrated Non-Volatile Memory (NVM) storage bandwidth, and data-intensive analytics workloads stress memory capacity and bandwidth. For this reason, datacenter operators aspire to move towards fully “disaggregated rack” architectures able to flexibly and dynamically allocate resources such as memory, storage, and compute in response to the mixture of tasks assigned to each cluster. Many contemporary resource disaggregation solutions are built upon Ethernet-based fabrics with electrical packet switches and 100 Gbps Ethernet NICs. However, the cost, power consumption, and latency of conventional Ethernet fabrics are severe inhibitors to efficient resource sharing. In particular, disaggregating high-performance memory over system-wide distances poses an exceptional challenge because of the extremely high data rates required [5]. The emergence of ultra-high-bandwidth-density photonic link technologies driven by efficient comb-laser sources open up the opportunity to break out of the package to enable more efficient rack-scale resource disaggregation.

### 1.2. The Challenges of Package Escape Bandwidth

To meet continued bandwidth demands, memory-intensive compute applications, such as GPUs, have largely moved to 2.5D integrated “in-package” memory technologies, which rely on higher areal density for electrical signals rather than higher signaling rates. Each signal connection runs at a comparatively low signaling rate (sub 5 GHz per IO connection), but with a much smaller pitch between the signal connections. The Heterogeneous Integration Roadmap (HIR) projects that five generations of exponential improvements in bandwidth density can be realized from known packaging technology roadmaps [2] with minimal SERDES increase. We posit that such wide-and-slow approach has a natural synergy with silicon photonic ring resonator technologies where the optimal efficiency is obtained from using more channels that operate at a slower-rate per channel rather than a small number of high-signal-rate SERDES rate channels. The emergence of efficient solid-state comb-laser sources [3] will take this technology to the next level, where energy efficient (1 picojoule/bit) short-reach (< 7 meters) high-bandwidth (terabyte/second) optical links can become practical and cost-effective. In particular, the bandwidth density of this emerging link technology will enable package escape bandwidths exceeding 1 Terabyte/second that can support high memory bandwidths, and enable scalable disaggregation strategies for datacenters in the future.

## 2. PINE: Photonics for Disaggregated Datacenters

The *Photonic Integrated Networked Energy efficient datacenter (PINE)* architecture [2] takes a systems-perspective on photonics integration to build resource disaggregation into the datacenter at a fundamental level. The PINE datacenter vision builds on three innovative pillars, shown in Figure 1:

1. A new generation of optical links specifically optimized for energy efficiency.
2. A concept of embedded, high bandwidth density photonic interconnectivity between various types of multi-chip modules (MCMs) in a unique interposer platform.
3. A re-designed system architecture with adaptive connectivity using bandwidth steering.

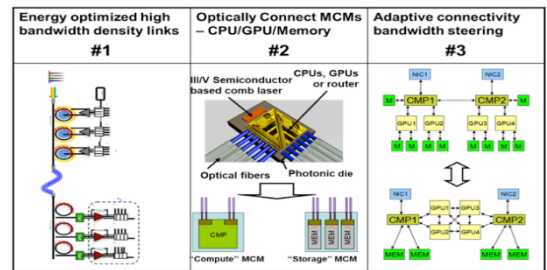


Figure 1: Three pillars of the ARPae PINE.

The PINE architectural design disaggregates key elements of the traditional datacenter and reorganizes them around a reconfigurable optical network fabric. Within the PINE datacenter, nodes are interconnected with ultra-low power consuming photonic links uniquely co-integrated into multi-chip modules (MCMs) in our interposer platform.

PINE disaggregation mechanisms assigns datacenter resources to workloads so that **only the required amount of computation power, memory capacity, and interconnectivity bandwidth are made available over the needed time period.** This efficient usage of resources reduces the vast amounts of wasted energy consumption of current datacenters.

## 3. Impact on Datacenter Operation

This section illustrates the opportunities for rack-scale disaggregation with PINE architecture by focusing on the need for memory disaggregation in today's clusters. The rule of thumb for memory per node over many decades has been to purchase one byte of memory capacity per peak floating-point operation per second of compute performance (one byte-per-flop in shorthand). However, this mythical ratio has rarely been tested in practice. Figure 2 shows the cumulative distribution function of memory utilization from the Cori supercomputing system at NERSC collected from over 9,000 nodes for a period of 48 months.

Although each server is equipped with 128 GB of memory, the figure shows that only 15% of jobs require nearly all of the 128 Gigabytes of memory in the node. Over 50% of CPU hours go to jobs that consume 25 Gigabytes or less.

The PINE photonic architecture offers the flexibility and bandwidth density needed to carry memory traffic at full rate, and offers the reach necessary to access as little or as much memory as is required for the job running on the node. By using broadband optical circuit switches to patch in the desired amount of memory (right-sizing) at full in-package bandwidths, memory disaggregation reduces the total amount of memory capacity purchased system-wide. Since memory consumes nearly 25-50% of total system cost and total system power consumption [4], the opportunity to reduce capital acquisition costs and operating costs for power are substantial.

Memory pressure at NERSC, 2018

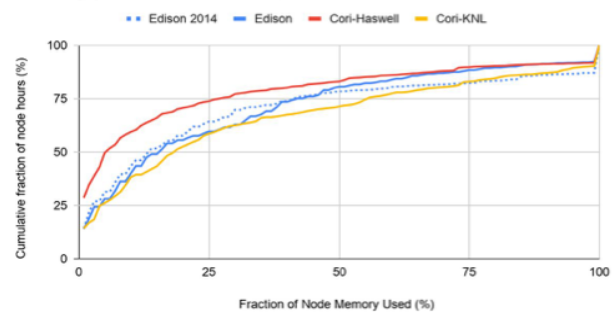


Figure 2: Memory utilization on NERSC systems.

## 4. References

- [1] 2019 Heterogeneous Integration Roadmap (HIR), <https://eps.ieee.org/technology/heterogeneous-integration-roadmap/2019-edition.html>
- [2] K. Bergman et. al., "PINE: An Energy Efficient Flexibly Interconnected Photonic Data Center Architecture for Extreme Scalability" in *2018 IEEE Optical Interconnects Conference (OI)*, June 4-6, 2018.
- [3] M. Bahadori, K. Bergman: Low-Power Optical Interconnects based on Resonant Silicon Photonic Devices: Recent Advances and Challenges. *ACM Great Lakes Symposium on VLSI 2018*: 305-310, 2018.
- [4] S. Ghose et al. 2018. What Your DRAM Power Models Are Not Telling You: Lessons from a Detailed Experimental Study. In *Proc. ACM Measurement and Analysis of Computing Systems*, Vol. 2, No. 3,). ACM, 2018.
- [5] Daniel Brunina, Caroline P. Lai, Ajay S. Garg, and Keren Bergman, "Building Data Centers With Optically Connected Memory," *Journal of Optical Communications and Networking*, Vol. 3, No. 8, 2011.