

ortheastern

CSR: Small: Collaborative Research: Leveraging Intra-chip/Inter-chip Silicon-Photonic Networks for Designing Next-Generation Accelerators Amir Kavyan Ziabari, Saiful Mojumder, Yifan Sun, José L. Abellán, John Kim, Ajay Joshi (PI), and David Kaeli (PI)

- and data-level parallelism
- throughput









UMH: A Hardware-Based Unified Memory Hierarchy [3]

Unified Memory

- Supported since CUDA 6.0, OpenCL 2.0 and HSA
- Improves programmability
- Severe runtime inefficiencies

Unified Memory Hierarchy (UMH)

Extends the memory hierarchy so that the GPU memory • acts as a last-level cache for system memory

• No redundant copies Block state Remote access only on miss Μ 0 GPU CPU GPU GPU GPU **GPU** Chip Data Cache request Line Stacked Stacked Stacked Stacked DRAM DRAM * * (SMDs SMDs SMDs LLC SMD SMD SMD SMD Performance From/to HMD HMD Host Memory Interconnect speedup Kernel Execution Memory Copy BSch DCT DH FW FWT Hist MM

Summary

- A Hybrid Silicon Photonic crossbar provides significantly improved performance, at the cost of slightly higher power consumption. • UMH improves the performance in a multi-GPU system by eliminating redundant memcpy operations. A combination of hybrid Silicon Photonic NoC and UMH can further improve multi-GPU performance.

Benchmarks

- system.

[1] Ziabari, A.-K., et al. "Leveraging Silicon-Photonic NOC for Designing Scalable GPUs." Proc. of the 29th ACM International Conference on Supercomputing. ACM, 2015, pp. 273-282. [2] Ziabari, A.-K., et al. "Asymmetric NoC Architectures for GPU Systems." Proc. of the 9th ACM International Symposium on Networks-on-Chip. (NOCS'15), 2015, pp. 1-8.

[3] Ziabari, A.-K., et al. "UMH: A Hardware-Based Unified Memory Hierarchy for Systems with Multiple Discrete GPUs." ACM Transactions on Architecture and Code Optimization (TACO) 13.4 (2016): 35, pp. 1-25. [4] Sun, Y., et al. "Hetero-mark, a Benchmark Suite for CPU-GPU Collaborative Computing." Proc. of the IEEE International Symposium on Workload Characterization." 2016, pp. 1-10.

SMD + HMD

NMOESI cache coherency protocol is used



 Allows for GPU data sharing • Enables CPU coherency

• Adds the N-state for non-coherent non-exclusive accesses 72MB overhead for 1GB data (to support coherence)

Processor action			Incoming request		
load	store	n-store	Eviction	Read request	Write request
hit	write request → M	hit	writeback → I	-	send data →∎
hit	hit	hit	writeback →	send data → O	send data →
hit	write request → M	hit	writeback → I	send data	send data →∎
hit	hit → M	hit → N	→	send data → S	send data →∎
hit	write request → M	hit → N	\rightarrow	-	\rightarrow
read request $\rightarrow \mathbf{S}$ or $\rightarrow \mathbf{E}$	write request → M	read request → N	-	-	-

• UMH improves performance by 1.9x on a 4-GPU system, as compared to memcpy

Going from 1-GPU to 4-GPUs, UMH achieves a 2.3x speedup, while memcpy only achieves a 1.7x



Future Work

• The performance of a multi-GPU system can be improved further by providing a truly shared memory

 The execution time of heterogeneous workloads [4], including machine learning algorithms, can be significantly reduced by introducing shared memory and a high bandwidth interconnection network.