

# Roofline analysis and profiling of AI / HPC apps

# Guide open systems for AI

Georgios Michelogiannakis, Raghu Shankar, and Boon Chong Ang

OCP AI/ML Infrastructure Workshop

June 2025

# Chiplets Modularity is a promising solution to make HPC&AI specialization accessible



Science

### Application Analysis (Profiling) Guides Chiplet Design

Profiling helps us pick a collection of chiplets (newly-developed or available from the community) for a set of applications

This way, we can serve the mix of applications found in open-science AI & HPC systems and maintain modularity

Collection of chiplets that maximize metric of interest



Application characterization



# **Application Requirements: Literature Search**

# Compute to Data IO intensity ranges in orders of magnitude – Opportunity for modular & tuned systems optimizing KPIs



- GPU nodes AMD EPYC 7763 "Milan", 256GB DDR4, 4x Nvidia A100, 40GB HBM2,
- 135,000 jobs, 650 distinct apps during Jul 2024 (avg. ~200 jobs/app)
- Majority of jobs out of machine balance
  - Jobs ~0.1 TF/s
  - Systems: ~10 TF/s
  - Source: "System-Wide Roofline Profiling a Case Study on NERSC's Perlmutter Supercomputer," Perf Modeling, Benchmarking, and Simulation (PMBS) Nov 2024, Brian Austin, Dhruva Kulkarni, et. al. LBNL

AI – Compute to Memory (Data) ranges in orders of magnitude Opportunity for modular & tuned systems optimizing KPIs in open systems for AI



Memory Bound

Compute Bound

•\* Arithmetic intensity varies with model size, model, sequence length, data set size, # parameters, #batches, layer

· Adapted from sources:

"HBM Roadmap ver 1.7 Workshop", KAIST Teralab, June 2025, Terabyte Interconnection and Package Lab "Full stack optimization of Transformer Inference: A Survey", Sehoon Kim, et. al, UC Berkeley, arXiv, Feb 2023

21

### Fugaku: HPC jobs characterizations shows wide range Opportunity: Modularity & IP built for AI leverage for HPC



- 2.1 M jobs
  - Memory bound 1.6M
  - Compute Bound 0.5M

#### HPC ranges more to the left .. 0.001 to 1

 Source: "MCBound: An Online Framework to Characterize and Classify Memory/Compute-bound HPC Jobs", Antici, et. al., SC24, Nov 2024

#### Opportunity for modular & tuned systems optimizing KPIs in open systems for AI



Memory Bound

Compute Bound

# Profiling Results Collected by HPC/AI OCP Workstream Members

### Particle Tracking For High Energy Physics (CPU Version)



A Combinatorial Explosion in Computational Complexity With Increased Luminosity for Future Experiments

### Impact of Parallelism to Memory



### Canneal: Simulated Annealing Approximates Global Optimum

### Memory bandwidth

Runtime 3s

#### Uncore Events ≽ 🖺

Uncore Event Short Name / Uncore Event Type Uncore Event Count

UNC_IMC_DRAM_DATA_READS	2,139,123,811
UNC_IMC_DRAM_DATA_READS	2,139,123,811
UNC_IMC_DRAM_DATA_WRITES	265,424,926
UNC_IMC_DRAM_DATA_WRITES	265,424,926

#### Bandwidth Utilization Histogram h

\*N/A is applied to non-summable metrics.

Explore bandwidth utilization over time using the histogram and identify memory objects or functions with maximum contribution to the high bandwidth utilization.

Bandwidth Domain: DRAM, GB/sec ~

#### 😔 Bandwidth Utilization Histogram 🎼

This histogram displays the wall time the bandwidth was utilized by certain value. Use sliders at the bottom of the histogram to define thresholds for Low, Medium and High utilization levels. You can use these bandwidth utilization types in the Bottom-up view t utilization type. To learn bandwidth capabilities, refer to your system specifications or run appropriate benchmarks to measure them; for example, Intel Memory Latency Checker can provide maximum achievable DRAM and Interconnect bandwidth.



#### Microarchitecture Usage<sup>®</sup>: 27.3% ► of Pipeline Slots ≿ $\odot$

Slots Slots

	Retiring <sup>(2)</sup> :	27.3%	of Pipeline Slots
	Front-End Bound <sup>®</sup> :	30.6% 🏲	of Pipeline Slots
	Bad Speculation <sup>(2)</sup> :	7.7%	of Pipeline Slots
ତ	Back-End Bound <sup>③</sup> :	34.4%	of Pipeline Slot
	☑ Memory Bound <sup>③</sup> :	21.4% 🖻	of Pipeline Slot
	S L1 Bound <sup></sup>	16.9% 🏲	of Clockticks
	OTLB Overhead <sup></sup> 𝔅	100.0% 🕅	of Clockticks
	Load STLB Hit <sup></sup> :	100.0% 🏲	of Clockticks
	Load STLB Miss <sup>@</sup> :	8.2%	of Clockticks
	Loads Blocked by Store Forwarding <sup>(2)</sup> :	0.2%	of Clockticks
	Lock Latency <sup>(2)</sup> :	1.3%	of Clockticks
	Split Loads ②:	0.0%	of Clockticks
	4K Aliasing <sup>③</sup> :	0.1%	of Clockticks
	FB Full <sup>(2)</sup> :	11.3% 🖻	of Clockticks
	L2 Bound <sup>(2)</sup> :	0.9%	of Clockticks
	S L3 Bound <sup></sup>	5.4%	of Clockticks
	Contested Accesses 2:	0.3%	of Clockticks
	Data Sharing <sup>(2)</sup> :	0.1%	of Clockticks
	L3 Latency <sup>(2)</sup> :	7.0% 🖻	of Clockticks
	SQ Full :	0.0% 🏲	of Clockticks
	Solution Sector Sec	18.8% 🛤	of Clockticks
	Memory Bandwidth <sup>(2)</sup> :	23.0% 🖻	of Clockticks
	Memory Latency <sup>(2)</sup> :	20.1% 🏲	of Clockticks
	Store Bound <sup>®</sup> :	0.6%	of Clockticks
	Core Bound <sup>(2)</sup> :	13.0% 🏲	of Pipeline Slots

#### Memory Bound<sup>®</sup>: 21.4% ► of Pipeline Slots $\odot$

	Cache Bound <sup>(2)</sup> :	23.3% 🏲	of Clockticks
$\odot$	DRAM Bound <sup>®</sup> :	18.8% 🏲	of Clockticks
	Average DRAM Bandwidth, GB/s:	4.169	



36971.496

# What's Next

# Request for profiling and prototyping resources

Results critical for recommending next-gen specialized modular chiplet systems

- 1. Compute/GPU resource for profiling from cloud service providers
  - 2-socket x86 servers nodes + 2x GPUs 6 months
  - 256 GB DIMM +
  - 2x Nvidia H100, 80 GB, HBM3 or equivalent
  - SSD 300GB +
  - Ubuntu 22.04
- 2. Need trial licenses for virtual prototyping:
  - Synopsys (platform architect, zebu for FPGA prototyping, design compiler, etc)
  - Cadence
  - Siemens EDA
- 3. Keysight CloudBuilder, Keysight Al Data Center Builder,
  - Software suite designed to emulate real-world AI workloads, enabling validation and optimization of AI infrastructure components like networks, hosts, and accelerators
- 4. Grow team of contributors!