



SIMR: Single Instruction Multiple Request Processing for Energy-Efficient Data Center Microservices

<u>Mahmoud Khairy*</u>, Ahmad Alawneh, Aaron Barnes, and Timothy G. Rogers Purdue University

*Currently at AMD Research

Home page: <u>https://mkhairy.github.io/</u> Contact: <u>abdallm@purdue.edu</u>

MICRO 2022

10/3/2022

Datacenter Power Breakdown



Datacenter Power Breakdown (from Google)

25-45% of datacenter power is consumed in CPU's instruction supply (frontend & OoO)

Barroso, Luiz André, and Urs Hölzle. "The datacenter as a computer: An introduction to the design of warehouse-scale machines." Synthesis lectures on computer architecture. 2018 Haj-Yihia, Jawad, et al. "Fine-grain power breakdown of modern out-of-order cores and its implications on skylake-based systems." ACM TACO 2016 Notes: in the TACO paper, Execution includes ALU+Reg+OoO. In the fig above, we exclude the OoO and add it to the frontend. Caches power include dynamic L1/L2/L3 cache power.. The numbers are collected with McPAT

CPU Power Breakdown



Key Observation #1: Single Program Multiple Data (SPMD) are abundant in the datacenters

Server Workloads on GPU's

- Key Idea: Exploit SPMD by batching requests and run them on <u>GPU's Single Instruction Multiple Thread</u> (SIMT) or <u>CPU's SIMD</u>
- Advantage: Significant energy efficiency (throughput/watts) vs multi-threaded CPU

• Drawbacks:

- (1) Hindering programmability (C++/PHP vs CUDA/OpenCL)
- (2) Limited system calls support
- (3) High service latency (10-6000x)
 - GPUs tradeoff single threaded optimizations (OoO, speculative execution, etc.) in favor of excessive multithreading
 - In SIMD, relying on branch predicates & fine grain context

<u>Recall</u>: GPUs and SIMDs were designed to execute data parallel portion (i.e., loops) not the entire application

Rhythm: Harnessing Data Parallel Hardware for Server Workloads

Sandeep R Agrawal Duke University sandeep@cs.duke.edu

John Tran NVIDIA johntran@nvidia.com Valentin Pistol Duke University pistol@cs.duke.edu

David Tarjan * NVIDIA

Rhythm, ASPLOS 2014

Jun Pang Duke University pangjun@cs.duke.edu

Alvin R Lebeck Duke University alvy@cs.duke.edu

MemcachedGPU: Scaling-up Scale-out Key-value Stores

Tayler H. Hetherington The University of British Columbia taylerh@ece.ubc.ca Mike O'Connor NVIDIA & UT-Austin moconnor@nvidia.com Tor M. Aamodt The University of British Columbia aamodt@ece.ubc.ca

MemcachedGPU, SoCC 2015

ispc: A SPMD Compiler for High-Performance CPU Programming

Matt Pharr Intel Corporation matt.pharr@intel.com William R. Mark Intel Corporation william.r.mark@intel.com

ispc, InPar 2012

"Slower but energy-efficient wimpy cores only win for general data center workloads if their singlecore speed is reasonably close to that of mid-range brawny cores"

> Up to 2x slower latency can be tolerated by data center providers

Barroso, Luiz André, and Urs Hölzle. "The datacenter as a computer: An introduction to the design of warehouse-scale machines." Synthesis lectures on computer architecture. 2018 Hölzle, Urs. "Brawny cores still beat wimpy cores, most of the time." IEEE MICRO 2010



Urs Hölzle Google SVP

Off-Chip BW Scaling



Key Observation #2: There is available headroom to increase on-chip throughput (thread count) in the foreseeable future.

How to increase on-chip throughput of CPU?

- Direction#1 (industry standard): Add more Chiplets + Cores + SMT 1
- Direction#2 (this work): Move to SIMT
 - More energy efficient (throughput/watts)
 - Cost-effective (throughput/area)
 - Better scalability



"Let's bring SIMT efficiency to the CPU world!"

SIMT Efficiency

CPU Multi-Core with Simultaneous Multi-Threading



SIMR System Overview



Client Requests (HTTP/RPC calls)

Batch Similar Requests (e.g. per API)



RPU Core

CPU vs GPU vs RPU

Metric	CPU	GPU		
Core model	000	In-Order		
Programming	General-Purpose	CUDA/OpenCL		
ISA	x86/ARM	HSAIL/PTX		
System Calls Support	Yes	No		
Thread grain	Coarse grain	Fine grain		
Threads per core	Low (1-8)	Massive (2K)		
Thread model	SMT	SIMT		
Consistency	Variant	Weak+NMCA*		
Interconnect	Mesh/Ring	Crossbar		

*NMCA: non-multi copy atomicity

Ren, Xiaowei, et al. "HMG: Extending cache coherence protocols across modern hierarchical multi-GPU systems." HPCA 2020



The RPU takes advantage of the latency optimizations and programmability of the CPU

& SIMT efficiency and memory model scalability of the GPU

RPU's Challenges

- Control Divergence
 - <u>Challenge</u>: Control divergence with high latency path
 - Solution: Optimized batching & System-level batch split
- Memory Divergence
 - <u>Challenge</u>: Cache/TLB contention & bank conflicts
 - <u>Solution</u>: Batch tuning, stack/memory coalescing and SIMR-aware memory allocation
- Larger execution units & cache resources
 - <u>Challenge</u>: Higher instruction execution & L1 hit latency
 - Solution: Exploit low IPC, less generated traffic and employ sub-batching interleaving



ncy employ sub-batching interleaving

RPU's Challenges

- Control Divergence
 - <u>Challenge</u>: Control divergence with high latency path
 - Solution: Optimized batching & System-level batch split
- Men
 - C Read more details in the paper on how we address these challenges
 - <u>Solution</u>: Batch tuning, stack/memory coalescing and SIMR-aware memory allocation
- Larger execution units & cache resources
 - <u>Challenge</u>: Higher instruction execution & L1 hit latency
 - Solution: Exploit low IPC, less generated traffic and employ sub-batching interleaving



ncy employ sub-batching interleaving

SIMT Control Efficiency



Notes: (1) Batch Size = 32 & #batches=75, (2) System Calls are not traced, (3) SIMT Eff = scalar-instructions / (batch-instructions * batch-size), (4) fine-grain locking are assumed. Other assumptions are included in the paper.



Efficiency and Service Latency Results (Simulation)

CPU(SMT-1) CPU(SMT-8) RPU(SIMT-32)



Higher Is better

Efficiency and Service Latency Results (Simulation)

■ CPU(SMT-1) ■ CPU(SMT-8) ■ RPU(SIMT-32)



U(SIMT-32) Higher Is better

16

Summary

- *Request Similarity* is abundant in the data center.
- We start with OoO CPU design and augment it with SIMT execution to maximize chip utilization and exploit the similarity.
- We co-design the software stack to support *batching* and awareness of SIMT execution.

SIMT efficiency is high in the open-source microservices we study.

We are very interested in evaluating SIMT control efficiency in proprietary production microservices.





 μ Suite: A Benchmark Suite for Microservices

Google facebook

Thank You! Q&A?

Instruction level parallelism (ILP) & Thread level parallelism (TLP)

Data level parallelism (DLP)

Request level parallelism (RLP)







Backup Slides

SIMT-friendly Microservices



Monolithic Service

Key Observation#3: Microservices reduce the per-thread cache requirement and minimize control-flow variations between concurrent threads

Aicroservices architecture +Smaller cache footprint +Less divergent

Batching Optimization

From Google's Production DL Inference

Production				MLPerf 0.7				
DNN	ms	batch	DNN	ms	batch	DNN	ms	batch
MLP0	7	200	RNN0	60	8	Resnet50	15	16
MLP1	20	168	RNN1	10	32	SSD	100	4
CNN0	10	8	BERT0	5	128	GNMT	250	16
CNN1	32	32	BERT1	10	64			

Memcached servers



Table 5. Latency limit in ms and batch size picked for TPUv4i.

DL Inference Batching

Network Batching

Key Observation#4: Modern data centers already rely on request batching heavily

Jouppi, Norman P., et al. "Ten Lessons From Three Generations Shaped Google's TPUv4i: Industrial Product." 2021 ISCA https://memcached.org/blog/nvm-multidisk/

Meisner, David, and Thomas F. Wenisch. "Dreamweaver: architectural support for deep sleep." ASPLOS 2012

Batching for deep sleep

Energy efficiency



Single Thread Latency

Latency & Energy-Efficiency Tradeoff

Energy efficiency



Single Thread Laten



HW/SW Stack

Webservice (C++, PHP, ...) ARM/x86 compiler HTTP server Runtime/libs (pthread, cstdlib, ..) OS (Process, VM, I/Os)

Multi Core CPU

CUDA compiler CUDA compiler Nvidia Triton HTTP server CUDA runtime/libs (cudalib, tensorRT, ..) OS (I/Os management) CUDA driver (VM/thread management) GPU Hardware

CPU SW Stack

GPU SW Stack

→ For RPU, we keep the SW programming interface as in the CPU
→ Some VM&process management system calls are reimplemented in the RPU driver to be batch-aware



RPU SW Stack

RPU HW



Energy Efficiency of CPU vs RPU (Analytical Model)



 \rightarrow an anticipated 2-10x energy efficiency gain can be achieved with RPU vs CPU

CPU Dynamic Energy Breakdown





Experimental Setup

Dynamic Instrumentation



Khairy, Mahmoud, et al. "Accel-Sim: An extensible simulation framework for validated GPU modeling." ISCA 2020 Zhang, Yangi, Yu Gan, and Christina Delimitrou. "uqSim: Scalable and Validated Simulation of Cloud Microservices." ISPASS 2019 Alawneh, Ahmad, et al. "A SIMT Analyzer for Multi-Threaded CPU Applications." ISPASS 2022 Sriraman, Akshitha, and Thomas F. Wenisch. "µ suite: a benchmark suite for microservices." IISWC 2018 Gan, Yu, et al. "An open-source benchmark suite for microservices and their hardware-software implications for cloud & edge systems." ASPLOS 2019 Li, Sheng, et al. "McPAT: An integrated power, area, and timing modeling framework for multicore and manycore architectures." MICRO 2009

Batching Opportunity for Facebook Services

- To amortize batching overhead, you either need:

 - (1) High service latency, with low traffic so service latency will amortize batching **OR** • (2) High traffic, with low service latency so high traffic will amortize batching **OR** • (3) High traffic and high service latency (ideal case)
- Let's take a look at Facebook in-production services:

μservice	Throughput (QPS)	Req. latency	Insn./query
Web	O (100)	O (ms)	O (10 ⁶)
Feed1	O (1000)	O (ms)	$O(10^9)$
Feed2	O (10)	O (s)	O (10 ⁹)
Ads1	O (10)	O (ms)	O (10 ⁹)
Ads2	O (100)	O (ms)	$O(10^9)$
Cache1	O (100K)	Ο (μs)	$O(10^3)$
Cache2	O (100K)	Ο (μs)	$O(10^3)$

Note: I was not able to calculate the exact batching overhead as the exact numbers are not shown and SLA (P99 latency) is not specified.

Sriraman, Akshitha, Abhishek Dhanotia, and Thomas F. Wenisch. "Softsku: Optimizing server architectures for microservice diversity@ scale." ASPLOS 2019



Batching Opportunity for Google Services

- (1) From Google in-production ML inference services:
 - Batching is widely used for DL inference with size = 8-20 reqs based on traffic and latency

Production				MLPerf 0.7				
DNN	ms	batch	DNN	ms	batch	DNN	ms	batch
MLP0	7	200	RNN0	60	8	Resnet50	15	16
MLP1	20	168	RNN1	10	32	SSD	100	4
CNN0	10	8	BERT0	5	128	GNMT	250	16
CNN1	32	32	BERT1	10	64			

Table 5. Latency limit in ms and batch size picked for TPUv4i.

• (2) Further, Google search service has a high service latency (~10s ms) and high traffic (~100K QPS), so they are a good candidate for batching

Jouppi, Norman P., et al. "Ten Lessons From Three Generations Shaped Google's TPUv4i: Industrial Product." 2021 ISCA

Quoted: "Clearly, datacenter applications limit latency, not batch size. Future DSAs should take advantage of larger batch sizes"

Thank You! Q&A?



Hardware



Accelerators