Server Architecture for the Post-Moore Era

Babak Falsafi Director, EcoCloud ecocloud.ch

ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE



Data Economics

П

П

П

Π

Π





1 terabyte =

1000 gigabytes

1TB

1 EB = ¹ billion gigabytes or **250 billion DVDs**



Big data is projected to grow into a market by 2017, up from \$10.2 BILLION in 2013

П

П П

All of the world's digital data equals about 900 exabytes, of which is created by individuals

1 petabyte = 1000 terabytes

1PB

l exabyte = 1000 petabytes

of the world's

data by 2020



1ZB

for more than

is nearly 2 times as large] FB = as the web archive at the **US Library of Congress**

ŒÐ

Internet-of-Things (IoT): Data in Flight





20 Billion Connected Devices



4 Zettabytes of Data, 10% of Digital Universe



Source: IDC Worldwide and Regional IoT forecast, EMC Digital Universe with Research and Analysis by IDC

Data Shaping All Science & Technology



Science entering 4th paradigm

Analytics using IT on

- Instrument data
- Simulation data
- Sensor data
- Human data
- ...

Complements theory, empirical science & simulation



The FOURTH PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

TONY HEY STEWART TANKIEY AND KRISTIN TOLL

Data-centric science key for innovation-based economies!





Perspective on Scaling

"Invent



Daily IT growth in 2014 = AII of AWS in 2004!

Modern Datacenters are Warehouse-Scale Computers



- Millions of interconnected home-brewed servers
- Centralization helps exploit economies of scale
- Network fabric provides micro-second connectivity
- At physical limits
- Need sources for
 - Electricity
 - Network
 - Cooling



20MW, 20x Football Field \$3 billion

Warning! Datacenters are not Supercomputers

- an EPFL research center
- Run heterogeneous data services at massive scale
- Driven for commercial use
- Fundamentally different design, operation, reliability, TCO
 - Density 10-25KW/rack as compared to 25-90KW/rack
 - Tier 3 (~2 hrs/downtime) vs.Tier I (upto I day/downtime)
 -and lots more

Datacenters are the IT utility plants of the future







Supercomputing

Cloud Computing

Cloud Taking Over Enterprise





But..., silicon is running out of steam!



Moore's law dying





Manycore Accelerators



With voltages leveling:

- Parallelism has emerged as the only silver bullet
- Use simpler cores
 - Prius instead of Audi R8
- Restructure software
- Each core →

fewer joules/op

Concentional Server Concentional Server CPU (e.g., Xeon) CPU (e.g., Xeon) CPU (e.g., Xeon)



Modern Manycor





Data parallelHigher memory b/w

Super simple cores

- Shared front end
- IOx slower clocks

Great for dense parallel computation



Parallelism Alone Can't Help ecocloue

ISA opportunities

- Integration
 - Less energy moving
 - Closer to memory
- Specialization
 - Customize work
 - Less work/computation
- Approximation
 - Adjust precision





Center to bring efficiency to data

- I 8 faculty, 50 researchers
- \$6M/year external funds

Mission:

- Energy-efficient data-centric IT
- From algorithms to infrastructure
- Maximizing value for data





ecocloud.ch



CREDIT SUISSE

Our Vision: Holistic Optimization of Datacenters



From algorithms to infrastructure:

- Cross-layer integration & specialization
- Introspection & resource provisioning

Open technologies!







Memory-Centric Servers Near-Memory Processing

Summary



Scale-Out Datacenters



Vast data sharded across servers

Memory-resident workloads

- Necessary for performance
- Major TCO burden

Put memory at the center

- Design system around memory
- Optimize for data services



Servers driven by the DRAM market!



Server Benchmarking with CloudSuite 3.0 (cloudsuite.ch)





Building block for Google PerfKit, EEMBC Big Data!

Services are Stuck in Memory





- On-chip memory overprovisioned
- Instruction supply is bottlenecked

Manycore Accelerator for Data Serving



CAVIUM

Case for Workload Optimized Processors For Next Generation Data Center & Cloud

Gopal Hegde VP/GM, Data Center Processing Group

Cavium Thunder X

- Based on SOP @ EPFL
- Designed to serve data
- Optimized code supply
- Trade off SRAM for cores
- Runs stock software
- 10x faster than Xeon for CloudSuite

NOC-Out: NoC for Server Processors



Exactly the **opposite** of current NoCs

- Cache coherent
- But, designed for core-to-cache traffic
- Not core-to-core!

LLC network:

Flattened Butterfly (FB) topology

Request & Reply networks:

- Tree topology
- Limited connectivity for efficiency

FB's performance at 1/10th cost





DRAM Cache with Storage-Class Memory



SCM extends DRAM capacity as memory

DRAM cache + SCM:

- Provides high b/w access to hot data
- Mitigates the read/write disparity
 - Data read/written in pages (bulk)

DRAM cache + battery:Helps mitigate persistent ordering stalls

Scale-Out NUMA: Rack-Scale Memory Pooling





Pool memory over a light fabric:

- Balance load skew in data serving
- Mitigate partition skew in analytics

soNUMA:

- Socket-integrated network interface (e.g., Sonoma)
- Protected global memory read/write + synch





Custom Computing [FPGA's vs. GPU's in Data centers, IEEE Micro'17]



Reconfigurable

- Best for spatial computing
- Not caching/reuse
- Parallel, dataflow
- IOx slower clocks
- Better for sparse arithmetic

Microsoft, Amazon & Intel



FPGA's in Servers





Latest version:

- One FPGA per blade
- Sits on the network
- Backend connected to CPU/NI
- Originally to accelerate Bing, Azure
- Now ML service called BrainWave
- Intel's HARP: tighter integration

Microsoft Unveils Catapult to Accelerate Bing! [EcoCloud Annual Event, June 5th, 2014]



Memory-Centric Accelerators Abound



- Linear algebra for ML/NN
- IOx over GPU
- ML as a service

Oracle's RAPID:

- Accelerator for analytics in SQL
- Data movement engine in hardware
- Custom message passing cores
- Up to 15x better perf/Watt over Xeon



Walkers: Accelerating Data Management



- Pointer-based data structures (e.g., hash table, B-tree)
- Parallel lookups require traversing chains
- Decouple chains in co-designed hw/sw



15x better performance/Watt over Xeon

Walkers in Software [VLDB'16]



Use insights to help Xeon

- Decouple hash & walk in software
- Create & manage queues in wraparound code

2.3x speedup on Xeon

- Unclogs dependences in microarchitecture
- Maximizes memory level parallelism
- Under consideration by SAP HANA [VLDB'18]

The Specialization Funnel



Specialized

- GPU/ThunderX
- DBToaster
- IX Kernel
- Tensorflow

ASIC

- Crypto/Bitcoin
- Network logic

General Purpose

- Intel CPU
- Oracle Database
- Linux
- Java/C

Specialize as algorithms mature Domain-specific languages to platforms



Modern apps/services are statistical Analog input, analog output

Key:

Much redundancy in data/arithmeticOutput quality not accuracy or error

Exploit inProcessing, communication, storage





) verview

Memory-Centric Servers

Near-Memory Processing

Summary



What happens on servers?



Huge datasets reside in memory

- Fetch data
- Perform minimal computation
- Repeat over dataset



Data-centric services revolve around data movement



[Dally, SC'14 Panel]

Energy is dominated by data movement!





CPU-DRAM bandwidth (DIMM Channel) → 24 GB/s
 External interfaces either low in b/w or power hungry

But, internal DRAM bandwidth (Chip) → 128 GB/s



Parallelism is limited by connectivity

NMP comes to rescue



Near-memory processing (NMP):

- A layer of logic placed closer to DRAM
- Die-stacked or on interposer
- Helps exploit parallelism

Reduces data movement



Redesign algorithms, SW & HW to realize NMP potential

Why not compute inside DRAM? eco

Idea emerged in the 90's

- IRAM/PIM
- Logic & Memory on the same die

Did not make it

- Lowers DRAM density
- Increases DRAM costs
- DRAM is highly cost-sensitive



an EPFL research center

Must maintain DRAM cost advantages

NMP Commandments [IEEE Micro issue on Big Data' I 6]



Not (CPU) business as usual

- I. DRAM favors streaming vs. random access
 - CPU's leverage reuse & locality in cache hierarchy
- 2. DRAM favors wide (slow) cores vs. many (fast) cores
 - Stream-level parallelism to match DRAM b/w
- 3. Memory must maintain semantics relative to CPU
 Shared address space + coherence between NMP & CPU

Must co-design algorithm/HW for NMP!

Why not random access?

Internal DRAM structure dictates

- Activating a IKB row of data
- Dominates access latency & energy

Treat DRAM as a block-oriented device

- Stream data
- Maximize bandwidth & efficiency

an EPFL research center

DRAM row



Example:

- For DRAM with 128 GB/s internal bandwidth
- Optimal (parallel) random access only captures ~8 GB/s
- Requires 5x more power

Must use algorithms that favor sequential access!

The Mondrian Data Engine [ISCA'17]



- SIMD cores + data streaming
 - I 024-bit SIMD @ I GHz
 - No caches







Algorithm/hardware co-design maximize near-memory performance





Iterates over a pair of tables Finds the matching keys in two tables Major operation in data management

Q: SELECT ... FROM R, S WHERE R.Key = S.Key



CPU-centric (Hash) Join



Performed in two phases: Partition & Probe

- I. Partition tables based on keys
- 2. Probe joins partitions
 - Optimized for random accesses to cached data



Access patterns in hash Join ecocloud



: Random access (local or remote)

Join operation on Mondrian



Revisiting Sort join [ASBD'14]:

- Sort join (O(nlogn)) vs. Hash Join (O(n))
- Sort tables and then merge join

Perform more work

But, sort and merge use streaming access

Trade algorithm complexity for streaming memory accesses

Comparing access patterns



Phases	Hash	Sort
I. Partitioning		<mark>()</mark>
2. Build / Sort		
3. Probe / Merge	$\overline{\boldsymbol{\bigotimes}}$	\odot

Random access (local or remote)
Streaming access (remote)
Streaming access (local)

Performance





- Algorithm alone gets ~ I0x [ASBD'I5]
- Algorithm/hardware co-design gets 50x





Trends for data & online services:

- Data growing fast
- Online services are in-memory
- Memory is a big fraction of TCO

Post-Moore server designs:

- Opportunities abound
- Processors, accelerators, memory, network, system

Integration + Specialization + Approximation





For more information please visit us at **ecocloud.ch**





BEFERRE