# Datacenter Simulation Methodologies Case Studies

# Tamara Silbergleit Lehman, Qiuyun Wang, Seyed Majid Zahedi and Benjamin C. Lee



This work is supported by NSF grants CCF-1149252, CCF-1337215, and STARnet, a Semiconductor Research Corporation Program, sponsored by MARCO and DARPA.

Time	Торіс
09:00 - 09:30	Introduction
09:30 - 10:30	Setting up MARSSx86 and DRAMSim2
10:30 - 11:00	Break
11:00 - 12:00	Spark simulation
12:00 - 13:00	Lunch
13:00 - 13:30	Spark continued
13:30 - 14:30	GraphLab simulation
14:30 - 15:00	Break
15:00 - 16:15	Web search simulation
16:15 - 17:00	Case studies



# Big data demands big computing, yet we face challenges...

- Architecture Design
- Systems Management
- Research Coordination



# Heterogeneity

- Tailors hardware to software, reducing energy
- Complicates resource allocation and scheduling
- Introduces risk

# Sharing

- Divides hardware over software, amortizing energy
- Complicates task placement and co-location
- Introduces risk



#### Heterogeneity for Efficiency

Heterogeneous datacenters deploy mix of server- and mobile-class hardware



#### Heterogeneity for Efficiency

Heterogeneous datacenters deploy mix of server- and mobile-class hardware

#### Heterogeneity and Markets

Agents bid for heterogeneous hardware in a market that maximizes welfare



#### Heterogeneity for Efficiency

Heterogeneous datacenters deploy mix of server- and mobile-class hardware

#### Heterogeneity and Markets

Agents bid for heterogeneous hardware in a market that maximizes welfare

#### Sharing and Game Theory

Agents share multiprocessors with game-theoretic fairness guarantees



# Datacenter Design and Management

#### Heterogeneity for Efficiency

Heterogeneous datacenters deploy mix of server- and mobile-class hardware

- "Web search using mobile cores" [ISCA'10]
- "Towards energy-proportional datacenter memory with mobile DRAM" [ISCA'12]



# Mobile versus Server Processors

### • Simpler Datapath

- Issue fewer instructions per cycle
- Speculate less often

### • Smaller Caches

- Provide less capacity
- Provide lower associativity

# Slower Clock

• Reduce processor frequency



# Applications in Transition



#### **Conventional Enterprise**

- Independent requests
- Memory-, I/O-intensive
- Ex: web or file server

### **Emerging Datacenter**

- Inference, analytics
- Compute-intensive
- Ex: neural network





- Distribute web pages among index servers
- Distribute queries among index servers
- Rank indexed pages with neural network



Bing.com [Microsoft]



- Joules per second ::  $\downarrow$  **10**× on Atom versus Xeon
- Queries per second ::  $\downarrow$  **2**×
- Queries per Joule ::  $\uparrow$  **5**×

# Case for Processor Heterogeneity



#### **Mobile Core Efficiency**

• Queries per Joule  $\uparrow$  **5**×

#### Mobile Core Latency

- 10% queries exceed cut-off
- Complex queries suffer

#### Heterogeneity

- Small cores for simple queries
- Big cores for complex queries



#### Reddi et al., "Web search using mobile cores" [ISCA'10]

# Memory Architecture and Applications



### **Conventional Enterprise**

- High bandwidth
- Ex: transaction processing

### **Emerging Datacenter**

- Low bandwidth (< 6% DDR3 peak)
- Ex: search [Microsoft], memcached [Facebook]



# Memory Capacity vs Bandwidth

# • Online Services

- Use < 6% bandwidth, 65-97% capacity
- Ex: Microsoft mail, map-reduce, search [Kansal+]

# • Memory Caching

- 75% of Facebook data in memory
- Ex: memcached, RAMCloud [Ousterhout+]

# • Capacity-Bandwidth Bundles

- Server with 4 sockets, 8 channels
- Ex: 32GB capacity, >100GB/s bandwidth



# Mobile-class Memory

# • Operating Parameters

- Lower active current (130mA vs 250mA)
- Lower standby current (20mA vs 70mA)

### • Low-power Interfaces

- No delay-locked loops, on-die termination
- Lower bus frequency (400 vs 800MHz)
- Lower peak bandwidth (6.4 vs 12.8GBps)



#### LP-DDR2 vs DDR3 [Micron]

# Source of Disproportionality



#### **Activity Example**

• 16% DDR3 peak

### Energy per Bit

- Large power overheads
- High cost per bit



#### "Calculating memory system power for DDR3" [Micron]



#### **Mobile Memory Efficiency**

• Bits / Joule  $\uparrow$  5×

#### Mobile Memory Bandwidth

• Peak B/W ↓ 0.5×

### Heterogeneity

- LPDDR for search, memcached
- DDR for databases, HPC



#### Heterogeneity and Markets

Agents bid for heterogeneous hardware in a market that maximizes welfare

- "Navigating heterogeneous processors with market mechanisms" [HPCA'13]
- "Strategies for anticipating risk in heterogeneous system design" [HPCA'14]



	RHEL	SLES	Windows	Windows with SQ	L Standard Windows with SC	2L Web
Region:	US East (N. Virginia)		×			
		VCPU	ECU	Memory (GiB)	Instance Storage (GB)	Linux/UNIX Usage
eneral P	urpose -	Current Gen	eration			
n3.medik	m	1	3	3.75	1 x 4 SSD	\$0.113 per Hour
n3Jarge		2	6.5	7.5	1 x 32 SSD	\$0.225 per Hour
n3.xlarge		4	13	15	2 x 40 SSD	\$0.450 per Hour
n3.2xlarg	10	8	25	30	2 × 80 SSD	\$0.900 per Hour
eneral P	urpose - I	Previous Ger	neration			
n1.small		1	1	1.7	1 x 160	\$0.060 per Hour
n1.medik	ım	1	2	3.75	1 x 410	\$0.120 per Hour
n1Jarge		2	4	7.5	2 × 420	\$0.240 per Hour
n1.xlarge		4	8	15	4 x 420	\$0.480 per Hour
ompute	Optimize	d - Current G	ieneration			
:3.large		2	7	3.75	2 x 16 SSD	\$0.150 per Hour
:3.xlarge		4	14	7.5	2 x 40 SSD	\$0.300 per Hour
:3.2xlarg	,	8	28	15	2 x 80 SSD	\$0.600 per Hour
:3.4xlarg		16	55	30	2 x 160 SSD	\$1.200 per Hour
:3.8xlarg	•	32	108	60	2 x 320 SSD	\$2.400 per Hour
ompute	Optimize	d - Previous	Generation			
:1.mediu	m	2	5	1.7	1 x 350	\$0.145 per Hour
1.xlarge		8	20	7	4 x 420	\$0.580 per Hour
c2.8xlar	10	32	88	60.5	4 x 840	\$2,400 per Hour

#### Systems are heterogeneous

- Virtual machines are sized
- Physical machines are diverse

### Heterogeneity is exposed

- Users assess machine price
- Users select machine type

### Burden is prohibitive

 Users must understand hardware-software interactions



Elastic Compute Cloud (EC2) [Amazon]

Risk: the possibility that something bad will happen

# **Understand Heterogeneity and Risk**

- What types of hardware?
- How many of each type?
- What allocation to users?

# Mitigate Risk with Market Allocation

- Ensure service quality
- Hide hardware complexity
- Trade-off performance and power



# Market Mechanism

- User specifies value for performance
- Market shields user from heterogeneity







#### User Provides...

- Task stream
- Service-level agreement

### Proxy Provides...

- µarchitectural insight
- Performance profiles
- Bids for hardware



Guevara et al. "Navigating heterogeneous processors with market mechanisms" [HPCA'13] Wu and Lee, "Inferred models for dynamic and sparse hardware-software spaces" [MICRO'12]

# Visualizing Heterogeneity (2 Processor Types)



- Ellipses represent hardware types
- Points are combinations of processor types
- Colors show QoS violations



# Further Heterogeneity (4 Processor Types)



- Best configuration is heterogeneous
- QoS violations fall  $16\% \rightarrow 2\%$
- Trade-offs motivate design for manageability



Guevara et al. "Navigating heterogeneous processors with market mechanisms" [HPCA'13] Guevara et al. "Strategies for anticipating risk in heterogeneous datacenter design" [HPCA'14]

#### Sharing and Game Theory

Agents share multiprocessors with game-theoretic fairness guarantees

"REF: Resource elasticity fairness with sharing incentives for multiprocessors" [ASPLOS'14]





### **Big Servers**

- Hardware is under-utilized
- Sharing amortized power

#### Heterogeneous Users

- Tasks are diverse
- Users are complementary
- Users prefer flexibility

### **Sharing Challenges**

- Allocate multiple resources
- Ensure fairness





- Alice and Bob are working on research papers
- Each has \$10K to buy computers
- Alice and Bob have different types of tasks
- Alice and Bob have different paper deadlines



# Strategic Behavior

- Alice and Bob are strategic
- Which is better?
  - Small, separate clusters
  - Large, shared cluster

- Suppose Alice and Bob share
  - Is allocation fair?
  - Is lying beneficial?





#### Image: [www.websavers.org]

### Users must share

• Overlooks strategic behavior

### Fairness policy is equal slowdown

• Fails to encourage envious users to share

### Heuristic mechanisms enforce equal slowdown

• Fail to give provable guarantees



# "If an allocation is both equitable and Pareto efficient, ... it is fair." [Varian, Journal of Economic Theory (1974)]



# • Sharing Incentives

Users perform no worse than under equal division



# • Sharing Incentives

Users perform no worse than under equal division

### • Envy-Free

No user envies another's allocation



## • Sharing Incentives

Users perform no worse than under equal division

### Envy-Free

No user envies another's allocation

### • Pareto-Efficient

No other allocation improves utility without harming others



## • Sharing Incentives

Users perform no worse than under equal division

### Envy-Free

No user envies another's allocation

### • Pareto-Efficient

No other allocation improves utility without harming others

### • Strategy-Proof

No user benefits from lying



# Cobb-Douglas Utility

$$\mathbf{u}(\mathbf{x}) = \prod_{\mathsf{r}=1}^{\mathsf{R}} \mathbf{x}_{\mathsf{r}}^{lpha_{\mathsf{r}}}$$

- u utility (e.g., performance)
- $\mathbf{x}_{\mathbf{r}}$  allocation for resource  $\mathbf{r}$  (e.g., cache size)
- $\alpha_{\mathbf{r}}$  elasticity for resource  $\mathbf{r}$
- Cobb-Douglas fits preferences in computer architecture
- Exponents model diminishing marginal returns
- Products model substitution effects



Zahedi et al. "REF: Resource elasticity fairness with sharing incentives for multiprocessors" [ASPLOS'14]

$$\mathsf{u}_1 = \mathsf{x}_1^{0.6} \mathsf{y}_1^{0.4} \qquad \mathsf{u}_2 = \mathsf{x}_2^{0.2} \mathsf{y}_2^{0.8}$$

 $\begin{array}{lll} u_1, u_2 & \mbox{performance} \\ x_1, x_2 & \mbox{allocated memory bandwidth} \\ y_1, y_2 & \mbox{allocated cache size} \end{array}$ 



# Possible Allocations

- 2 users
- 12MB cache
- 24GB/s bandwidth





# Envy-Free (EF) Allocations

- Identify EF allocations for each user
- $u_1(A_1) \ge u_1(A_2)$
- $u_2(A_2) \ge u_2(A_1)$





• No other allocation improves utility without harming others





Fairness = Envy-freeness + Pareto efficiency



Memory Bandwidth

### Many possible fair allocations!



Zahedi et al. "REF: Resource elasticity fairness with sharing incentives for multiprocessors" [ASPLOS'14]



### Guarantees desiderata

- Sharing incentives
- Envy-freeness
- Pareto efficiency
- Strategy-proofness





# Off-line profiling

• Synthetic benchmarks

# **Off-line simulations**

• Various hardware

# Machine learning

•  $\alpha = 0.5$ , then update





•  $\mathbf{u} = \prod_{\mathbf{r}=1}^{\mathbf{R}} \mathbf{x}_{\mathbf{r}}^{\alpha_{\mathbf{r}}}$ 

• 
$$\log(\mathbf{u}) = \sum_{\mathbf{r}=1}^{\mathbf{R}} \alpha_{\mathbf{r}} \log(\mathbf{x}_{\mathbf{r}})$$

- Use linear regression to find  $\alpha_{\rm r}$ 



# Cobb-Douglas Accuracy



Ferret Sim. + Ferret Est.

- Utility is instructions per cycle
- Resources are cache size, memory bandwidth •



Zahedi et al. "REF: Resource elasticity fairness with sharing incentives for multiprocessors" [ASPLOS'14]



• Compare users' elasticities on same scale

• 
$$u = x^{0.2}y^{0.3} \rightarrow u = x^{0.4}y^{0.6}$$





$$\begin{aligned} \mathbf{u}_1 &= \mathbf{x}_1^{0.6} \mathbf{y}_1^{0.4} & \mathbf{u}_2 &= \mathbf{x}_2^{0.2} \mathbf{y}_2^{0.8} \\ \mathbf{x}_1 &= \left(\frac{0.6}{0.6+0.2}\right) \times 24 = 18 \text{GB/s} \\ \mathbf{x}_2 &= \left(\frac{0.2}{0.6+0.2}\right) \times 24 = 6 \text{GB/s} \end{aligned}$$



# Equal Slowdown versus REF



- Equal slow-down provides neither SI nor EF
- Canneal receives < half of cache, memory



# Equal Slowdown versus REF



- Equal slow-down provides neither SI nor EF
- Canneal receives < half of cache, memory



- Resource elasticity fairness provides both SI and EF
- Canneal receives more cache, less memory



Zahedi et al. "REF: Resource elasticity fairness with sharing incentives for multiprocessors" [ASPLOS'14]

# Performance versus Fairness



- Measure weighted instruction throughput
- REF incurs < 10% penalty



Zahedi et al. "REF: Resource elasticity fairness with sharing incentives for multiprocessors" [ASPLOS'14]

# Datacenter Design and Management

#### Heterogeneity for Efficiency

Heterogeneous datacenters deploy mix of server- and mobile-class hardware

- "Web search using mobile cores" [ISCA'10]
- "Towards energy-proportional datacenter memory with mobile DRAM" [ISCA'12]

#### Heterogeneity and Markets

Agents bid for heterogeneous hardware in a market that maximizes welfare

- "Navigating heterogeneous processors with market mechanisms" [HPCA'13]
- "Strategies for anticipating risk in heterogeneous system design" [HPCA'14]

#### Sharing and Game Theory

Agents share multiprocessors with game-theoretic fairness guarantees

"REF: Resource elasticity fairness with sharing incentives for multiprocessors" [ASPLOS'14]



# Datacenter Design and Management

#### Heterogeneity for Efficiency

Heterogeneous datacenters deploy mix of server- and mobile-class hardware

- Processors hardware counters for CPI stack
- Memories simulator for cache, bandwidth activity

#### Heterogeneity and Markets

Agents bid for heterogeneous hardware in a market that maximizes welfare

- Processors simulator for core performance
- Server Racks queueing models (e.g., M/M/1)

#### Sharing and Game Theory

Agents share multiprocessors with game-theoretic fairness guarantees

Memories – simulator for cache, bandwidth utility

