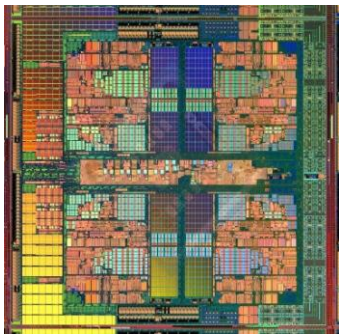


Computer Architecture: Multi-Core Evolution and Design

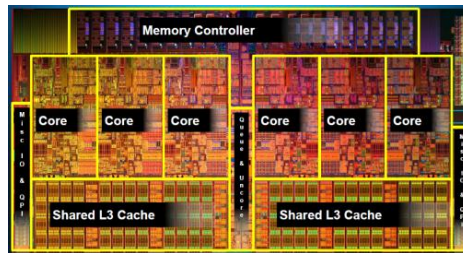
Prof. Onur Mutlu
Carnegie Mellon University

Multiple Cores on Chip

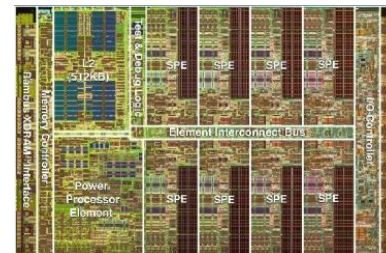
- Simpler and lower power than a single large core
- Large scale parallelism on chip



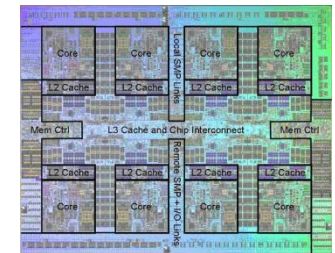
AMD Barcelona
4 cores



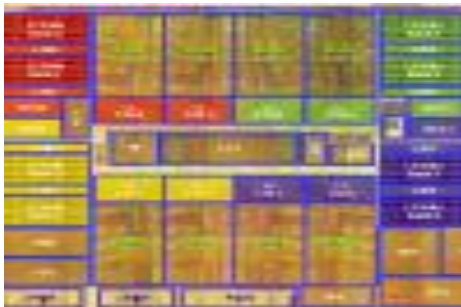
Intel Core i7
8 cores



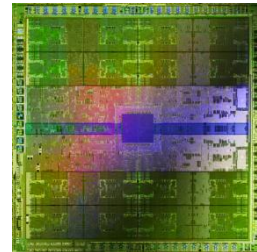
IBM Cell BE
8+1 cores



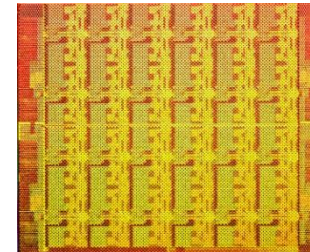
IBM POWER7
8 cores



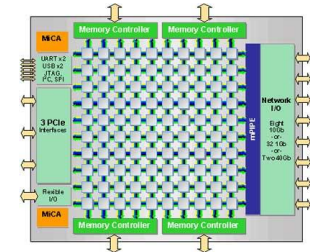
Sun Niagara II
8 cores



Nvidia Fermi
448 "cores"



Intel SCC
48 cores, networked



Tiler TILE Gx
100 cores, networked

With Multiple Cores on Chip

- What we want:
 - N times the performance with N times the cores when we parallelize an application on N cores

- What we get:
 - Amdahl's Law (serial bottleneck)
 - Bottlenecks in the parallel portion

Caveats of Parallelism

■ Amdahl's Law

- f: Parallelizable fraction of a program
- N: Number of processors

$$\text{Speedup} = \frac{1}{1 - f + \frac{f}{N}}$$

- Amdahl, “Validity of the single processor approach to achieving large scale computing capabilities,” AFIPS 1967.
- **Maximum speedup limited by serial portion: Serial bottleneck**
- **Parallel portion is usually not perfectly parallel**
 - **Synchronization** overhead (e.g., updates to shared data)
 - **Load imbalance** overhead (imperfect parallelization)
 - **Resource sharing** overhead (contention among N processors)

The Problem: Serialized Code Sections

- Many parallel programs cannot be parallelized completely
- Causes of serialized code sections
 - Sequential portions (Amdahl's "serial part")
 - Critical sections
 - Barriers
 - Limiter stages in pipelined programs
- Serialized code sections
 - Reduce performance
 - Limit scalability
 - Waste energy

Example from MySQL

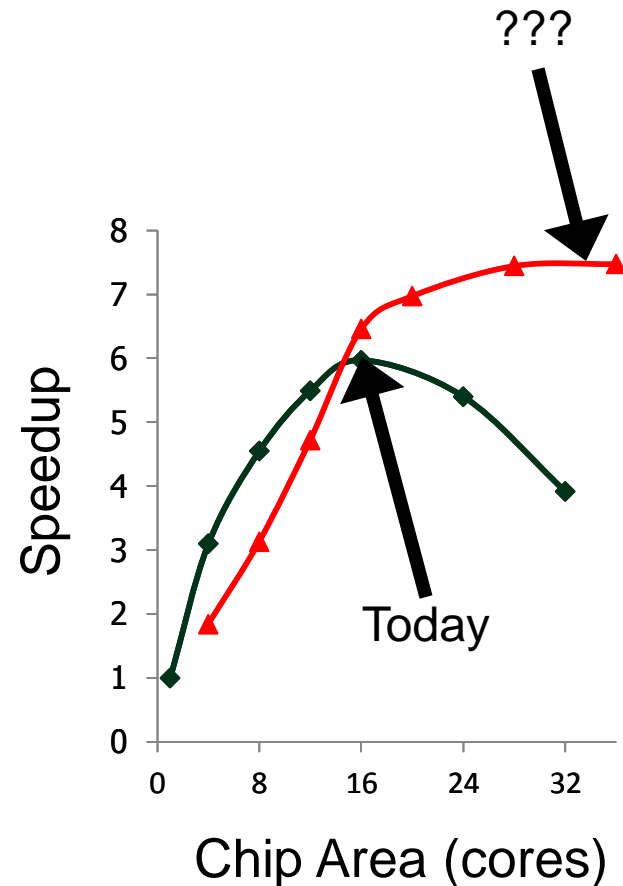
**Critical
Section**

Access Open Tables Cache

Open database tables

Perform the operations
....

Parallel



Demands in Different Code Sections

- What we want:
- In a serialized code section → one powerful “large” core
- In a parallel code section → many wimpy “small” cores
- These two conflict with each other:
 - If you have a single powerful core, you cannot have many cores
 - A small core is much more energy and area efficient than a large core

“Large” vs. “Small” Cores

Large Core

- *Out-of-order*
- *Wide fetch e.g. 4-wide*
- *Deeper pipeline*
- *Aggressive branch predictor (e.g. hybrid)*
- *Multiple functional units*
- *Trace cache*
- *Memory dependence speculation*

Small Core

- *In-order*
- *Narrow Fetch e.g. 2-wide*
- *Shallow pipeline*
- *Simple branch predictor (e.g. Gshare)*
- *Few functional units*

Large Cores are power inefficient:
e.g., 2x performance for 4x area (power)

Large vs. Small Cores

- Grochowski et al., “Best of both Latency and Throughput,” ICCD 2004.

	Large core	Small core
Microarchitecture	Out-of-order, 128-256 entry ROB	In-order
Width	3-4	1
Pipeline depth	20-30	5
Normalized performance	5-8x	1x
Normalized power	20-50x	1x
Normalized energy/instruction	4-6x	1x

Meet Small Cores: Piranha Chip Multiprocessor

- Barroso et al., “Piranha: A Scalable Architecture Based on Single-Chip Multiprocessing,” ISCA 2000.
 - An early example of a symmetric multi-core processor
 - Large-scale server based on CMP nodes
 - Designed for commercial workloads
 - Read:
 - Barroso et al., “Memory System Characterization of Commercial Workloads,” ISCA 1998.
 - Ranganathan et al., “Performance of Database Workloads on Shared-Memory Systems with Out-of-Order Processors,” ASPLOS 1998.
-

Commercial Workload Characteristics

- Memory system is the main bottleneck
 - Very high CPI
 - Execution time dominated by memory stall times
 - Instruction stalls as important as data stalls
 - Fast/large L2 caches are critical
- Very poor Instruction Level Parallelism (ILP) with existing techniques
 - Frequent hard-to-predict branches
 - Large L1 miss ratios
 - Small gains from wide-issue out-of-order techniques
- No need for floating point and multimedia units

Piranha Processing Node

CPU

Alpha core:
1-issue, in-order,
500MHz

Next few slides from

Luiz Barroso's ISCA 2000 presentation of

Piranha: A Scalable Architecture
Based on Single-Chip Multiprocessing



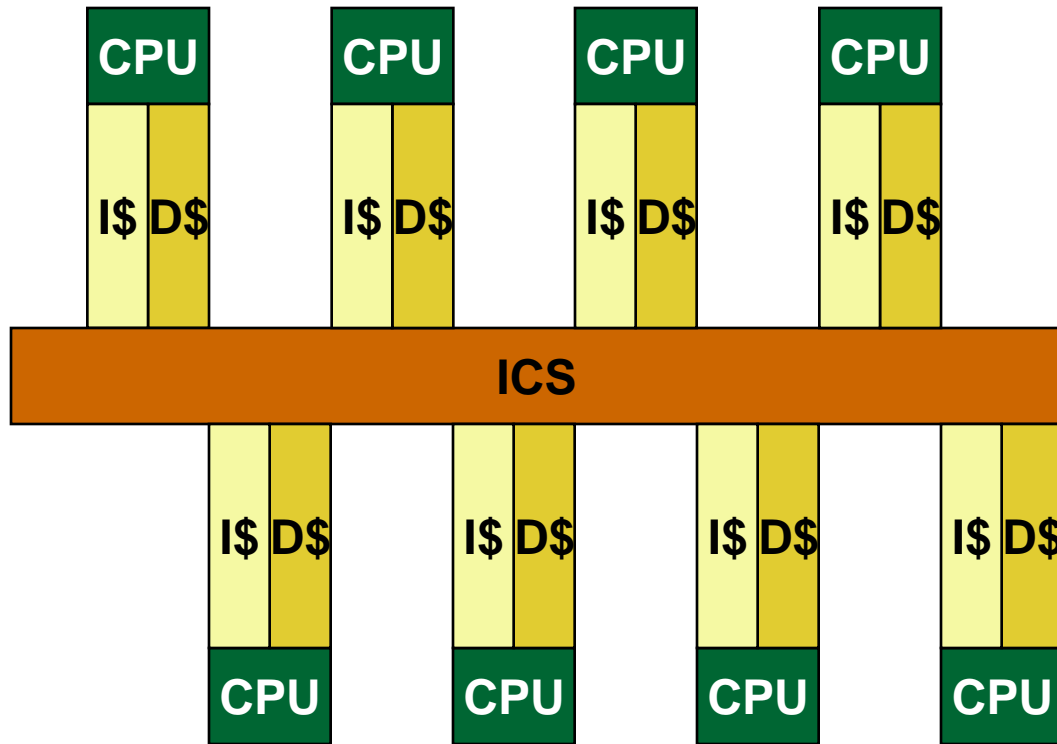
Piranha Processing Node



Alpha core:
1-issue, in-order,
500MHz
L1 caches:
I&D, 64KB, 2-way



Piranha Processing Node



Alpha core:

1-issue, in-order,
500MHz

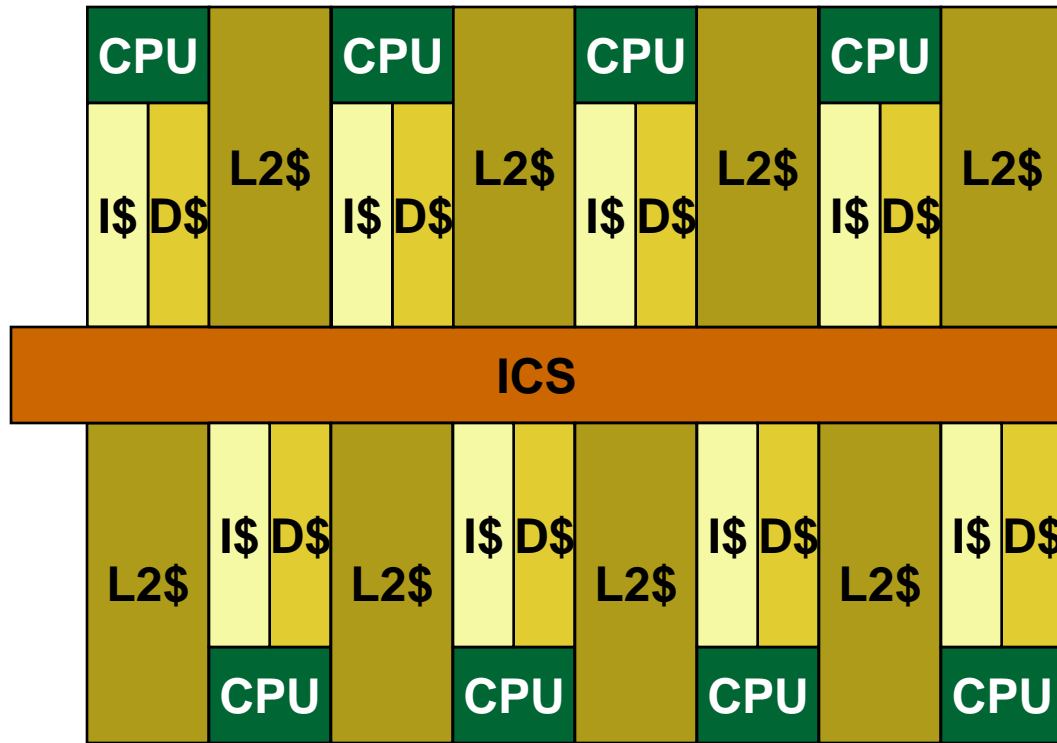
L1 caches:

I&D, 64KB, 2-way

Intra-chip switch (ICS)
32GB/sec, 1-cycle
delay



Piranha Processing Node



Alpha core:

1-issue, in-order,
500MHz

L1 caches:

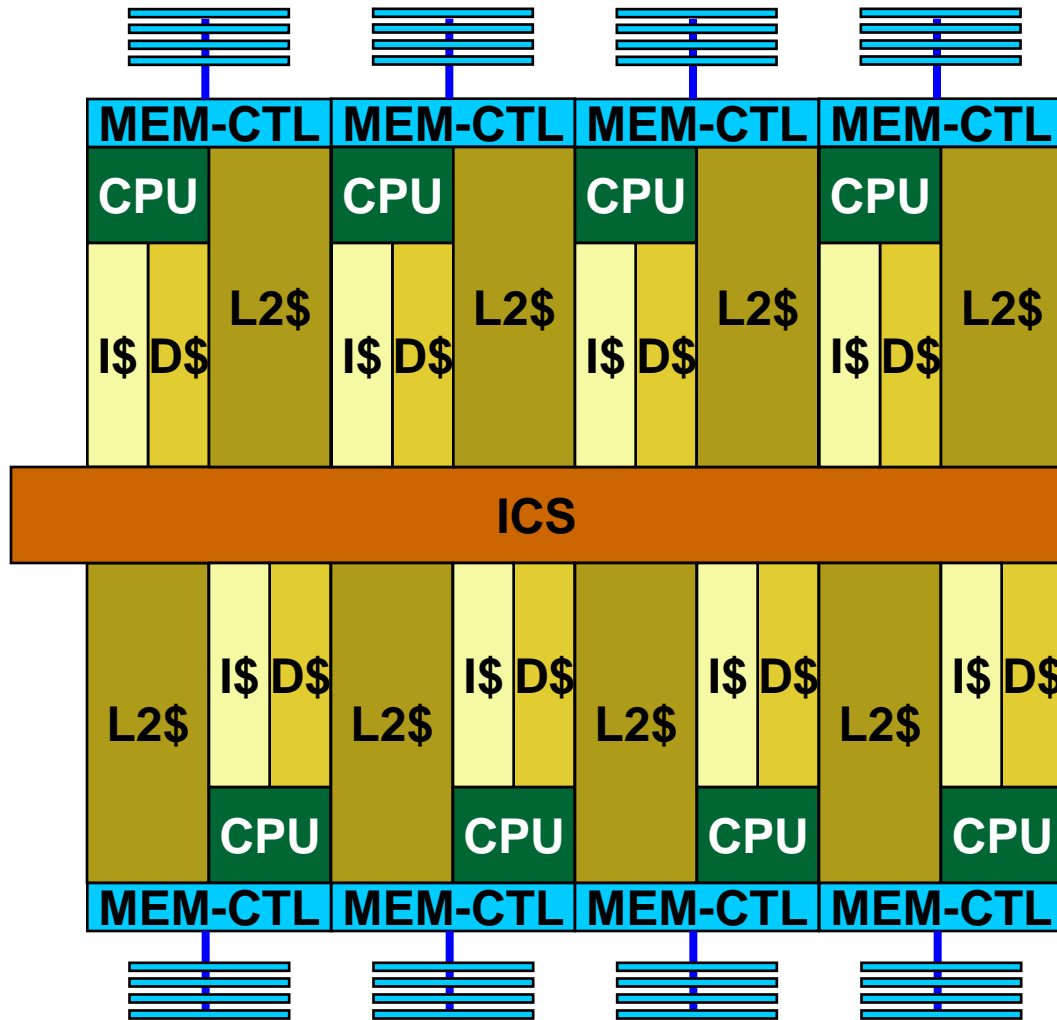
I&D, 64KB, 2-way
Intra-chip switch (ICS)
32GB/sec, 1-cycle
delay

L2 cache:

shared, 1MB, 8-way



Piranha Processing Node

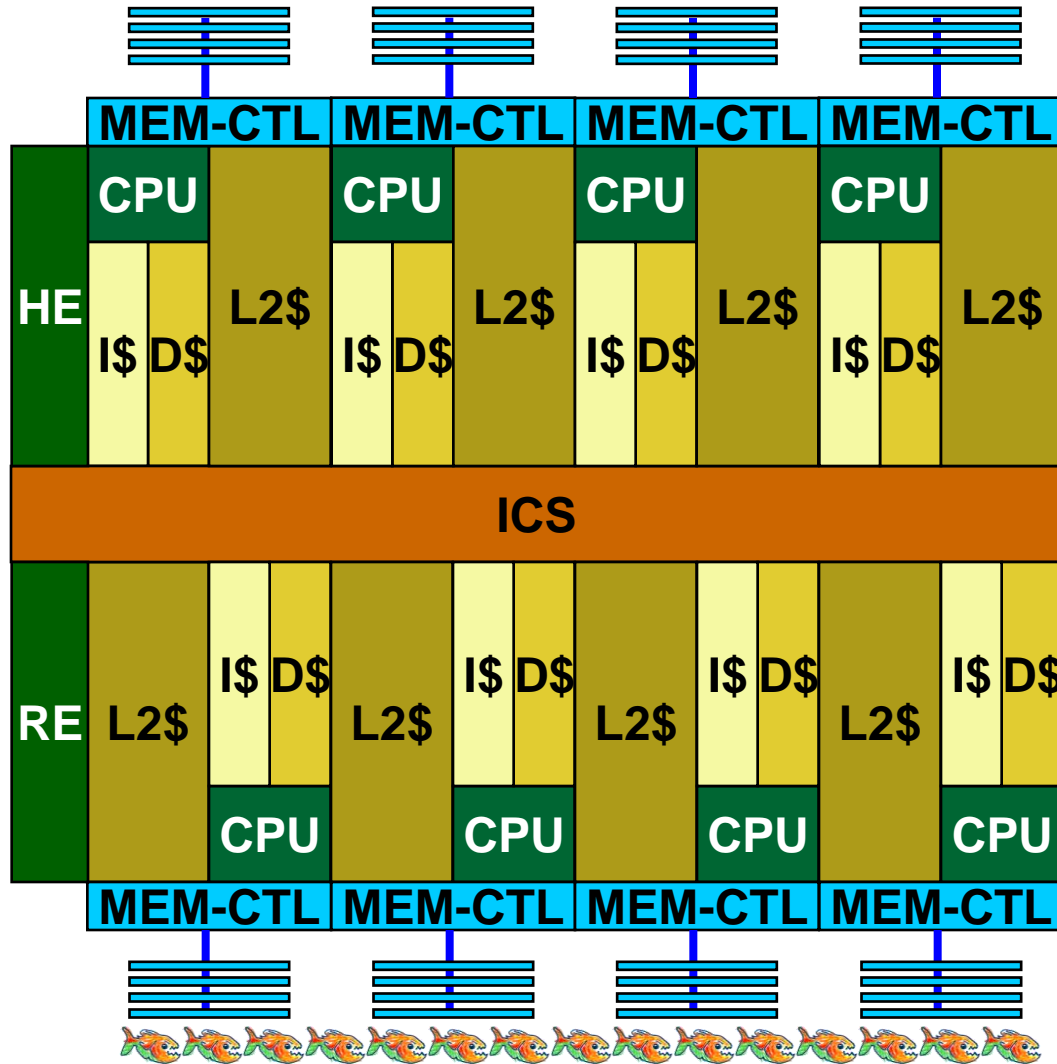


Alpha core:
1-issue, in-order,
500MHz
L1 caches:
I&D, 64KB, 2-way
Intra-chip switch (ICS)
32GB/sec, 1-cycle
delay
L2 cache:
shared, 1MB, 8-way
**Memory Controller
(MC)**
RDRAM, 12.8GB/sec

8 banks
@1.6GB/sec



Piranha Processing Node



Alpha core:

1-issue, in-order,
500MHz

L1 caches:

I&D, 64KB, 2-way
Intra-chip switch (ICS)
32GB/sec, 1-cycle
delay

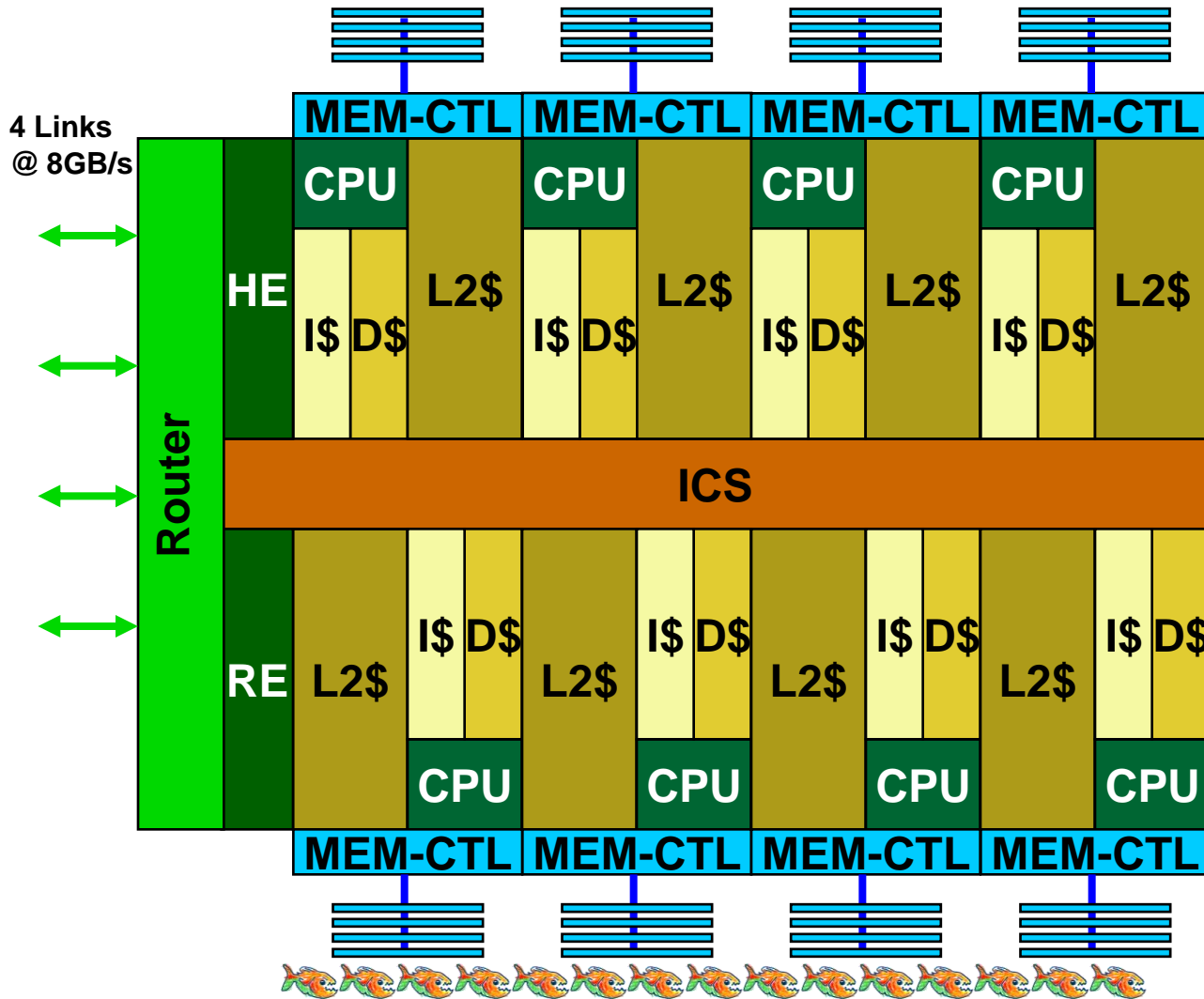
L2 cache:

shared, 1MB, 8-way
Memory Controller (MC)
RDRAM, 12.8GB/sec

Protocol Engines (HE & RE)

μ prog., 1K μ instr.,
even/odd interleaving

Piranha Processing Node



Alpha core:

1-issue, in-order,
500MHz

L1 caches:

I&D, 64KB, 2-way
Intra-chip switch (ICS)
32GB/sec, 1-cycle
delay

L2 cache:

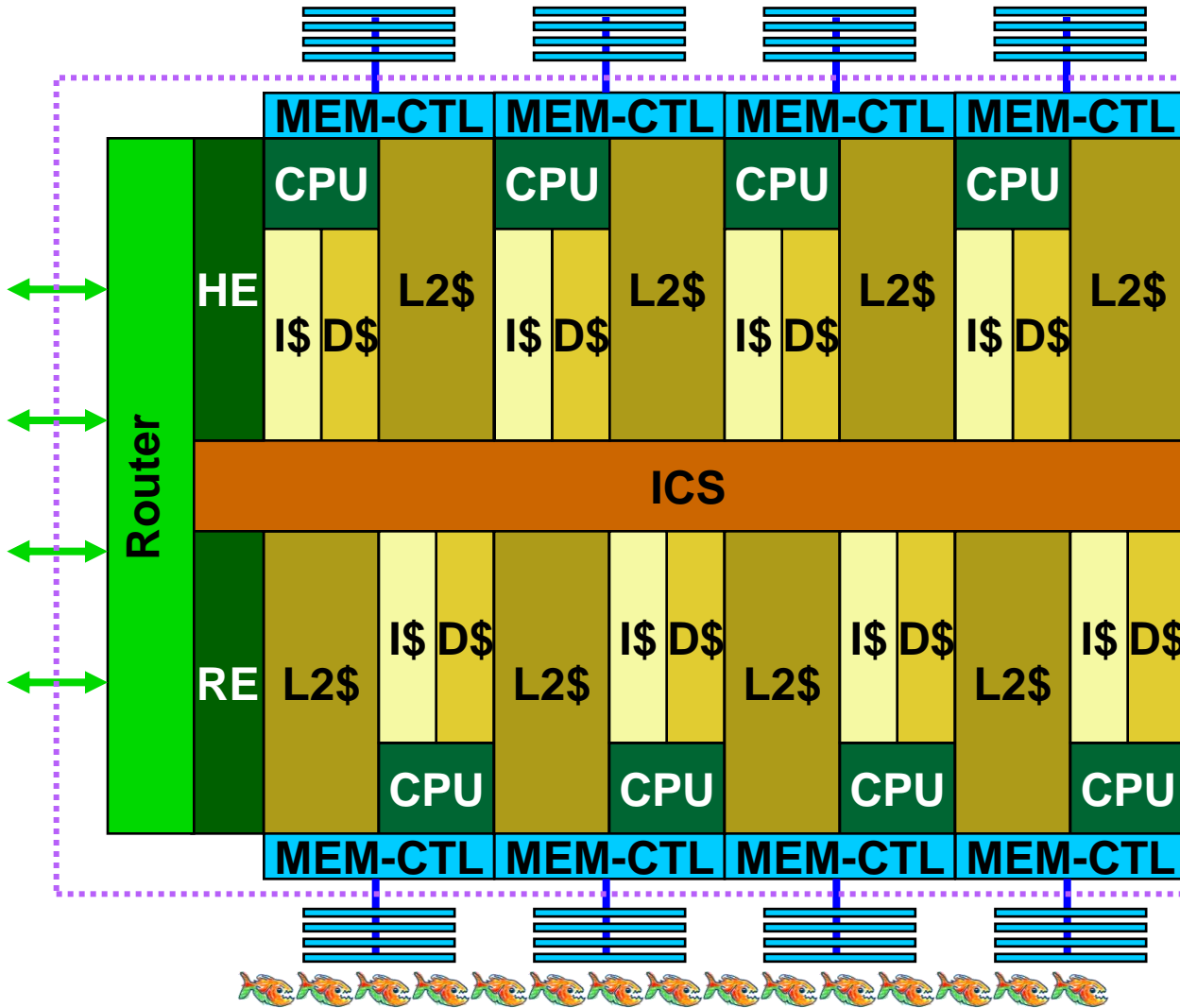
shared, 1MB, 8-way
Memory Controller (MC)
RDRAM, 12.8GB/sec
Protocol Engines (HE &
RE):

μ prog., 1K μ instr.,
even/odd interleaving

System Interconnect:

4-port Xbar router
topology independent
32GB/sec total
bandwidth

Piranha Processing Node



Alpha core:
 1-issue, in-order,
 500MHz

L1 caches:
 I&D, 64KB, 2-way

Intra-chip switch (ICS)
 32GB/sec, 1-cycle
 delay

L2 cache:
 shared, 1MB, 8-way

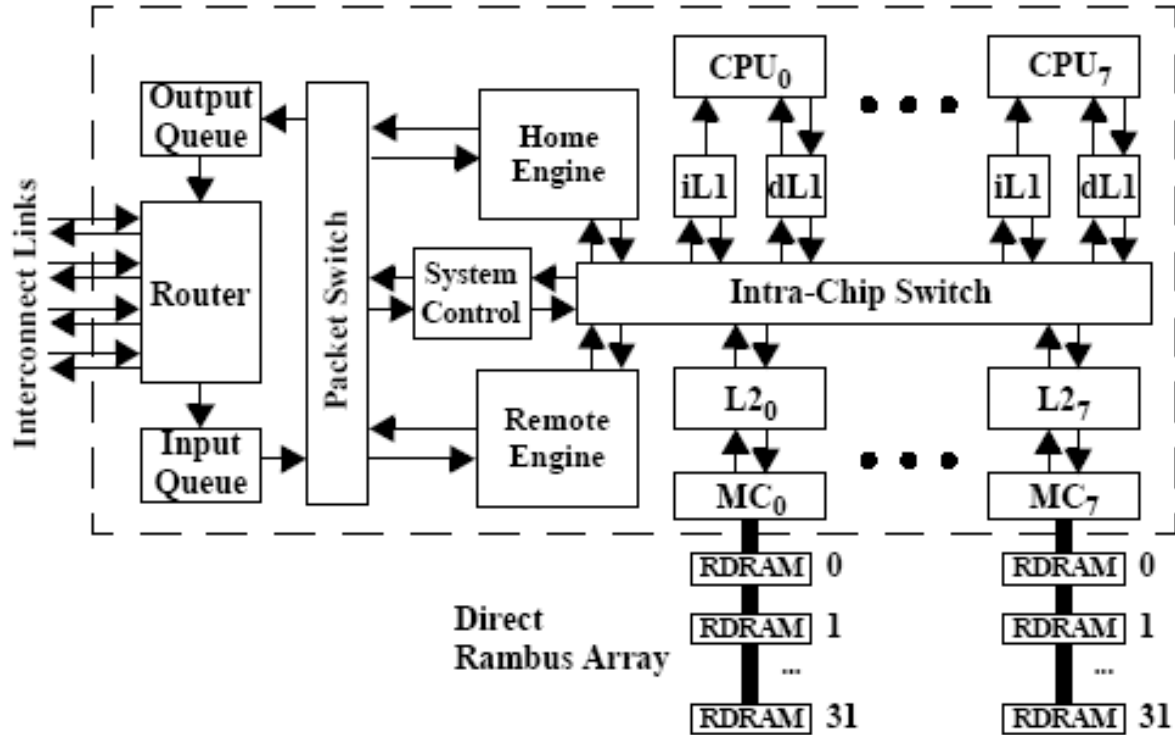
Memory Controller (MC)
 RDRAM, 12.8GB/sec

Protocol Engines (HE &
 RE):
 µprog., 1K µinstr.,
 even/odd interleaving

System Interconnect:
 4-port Xbar router
 topology independent
 bandwidth

Single Chip

Piranha Processing Node



Inter-Node Coherence Protocol Engine

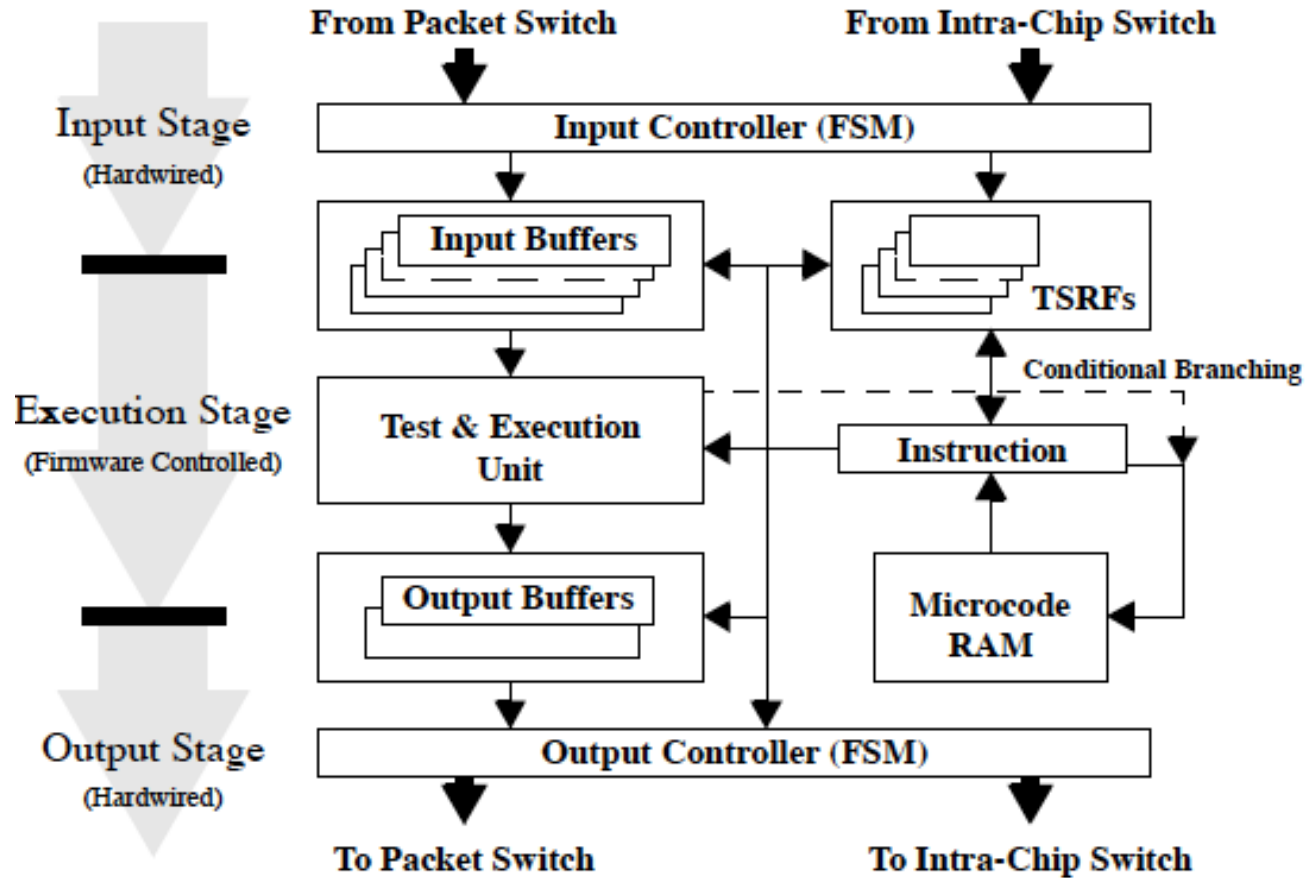


Figure 4. Block diagram of a protocol engine.

Piranha System

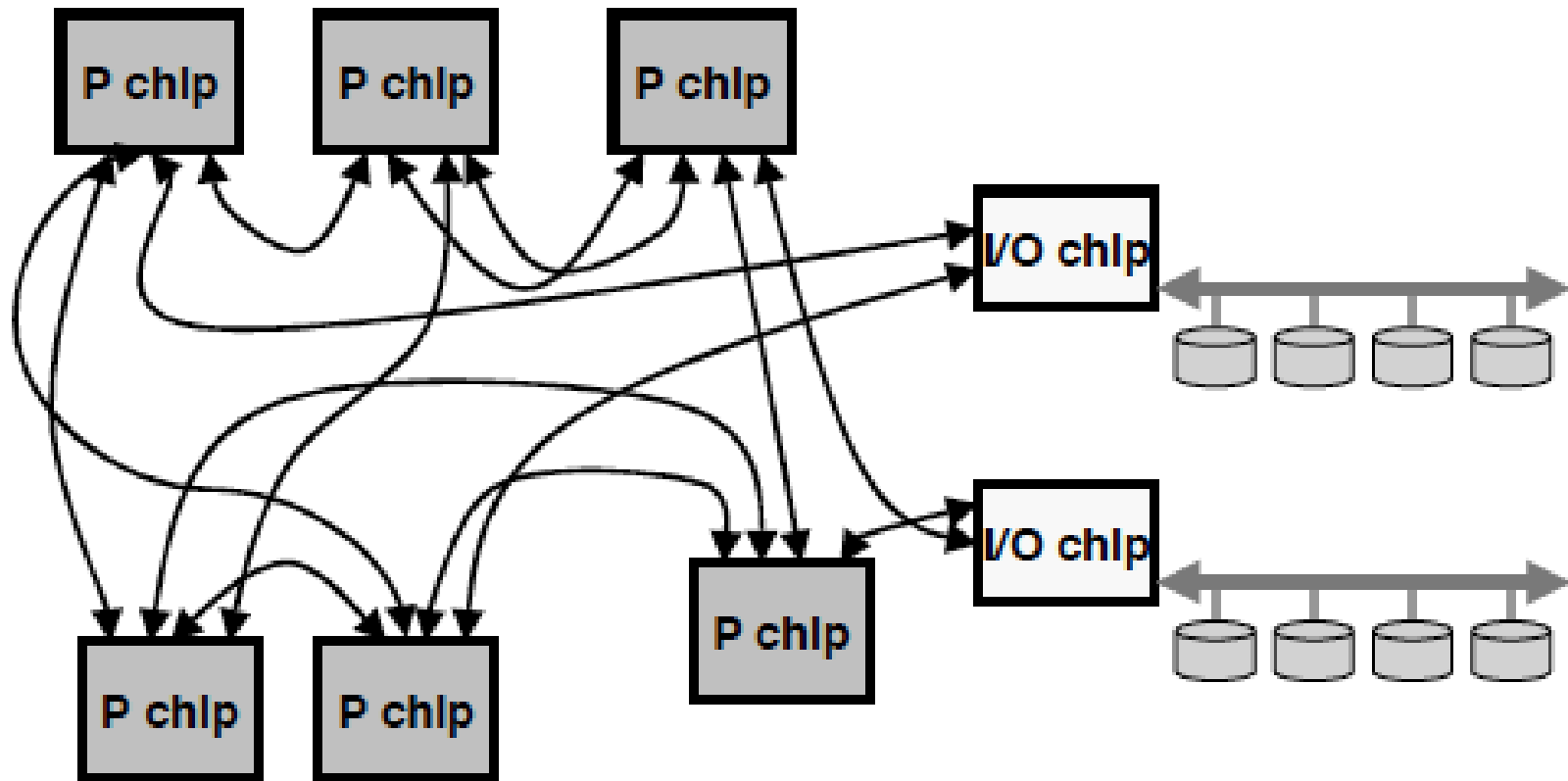


Figure 3. Example configuration for a Piranha system with six processing (8 CPUs each) and two I/O chips.

Piranha I/O Node

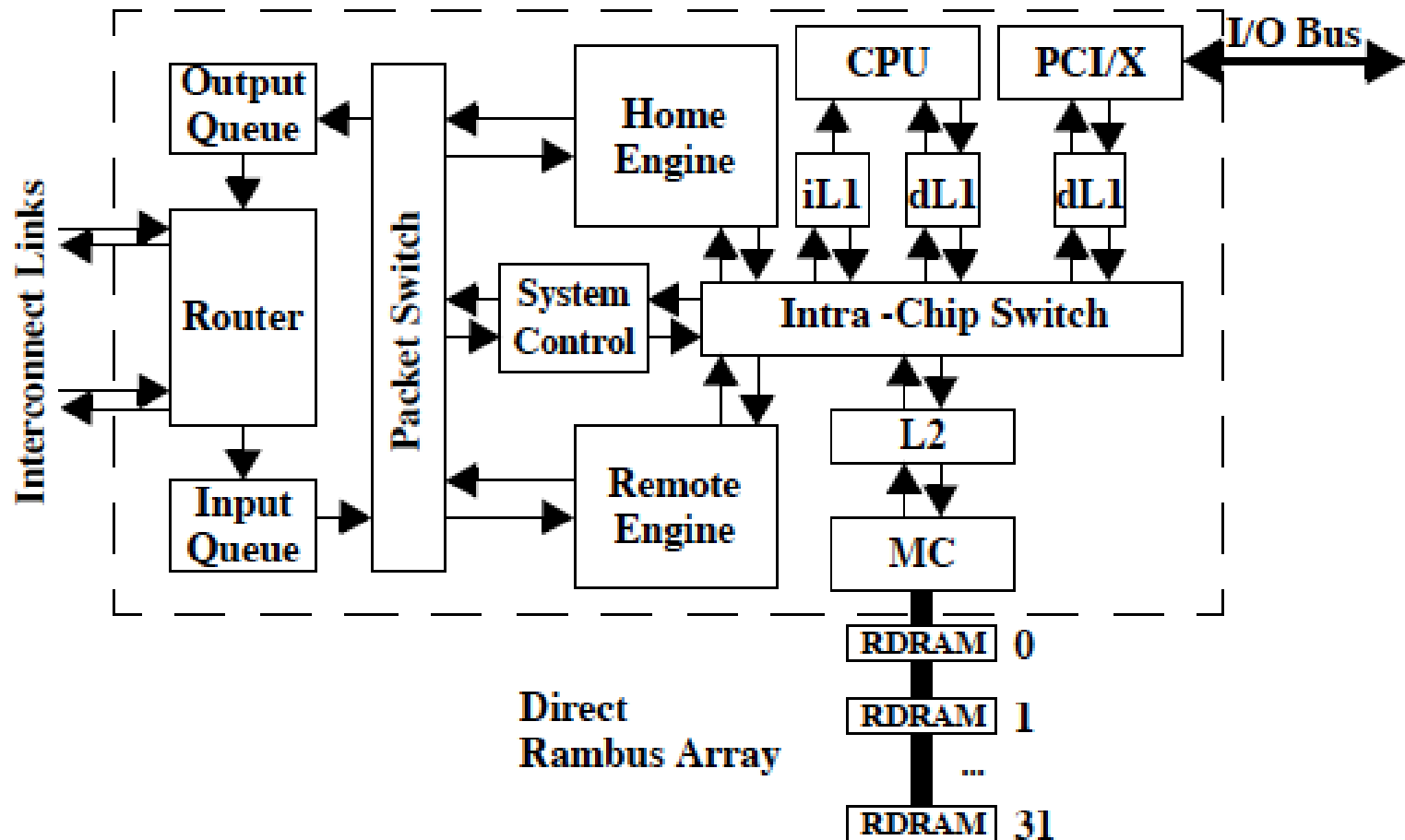
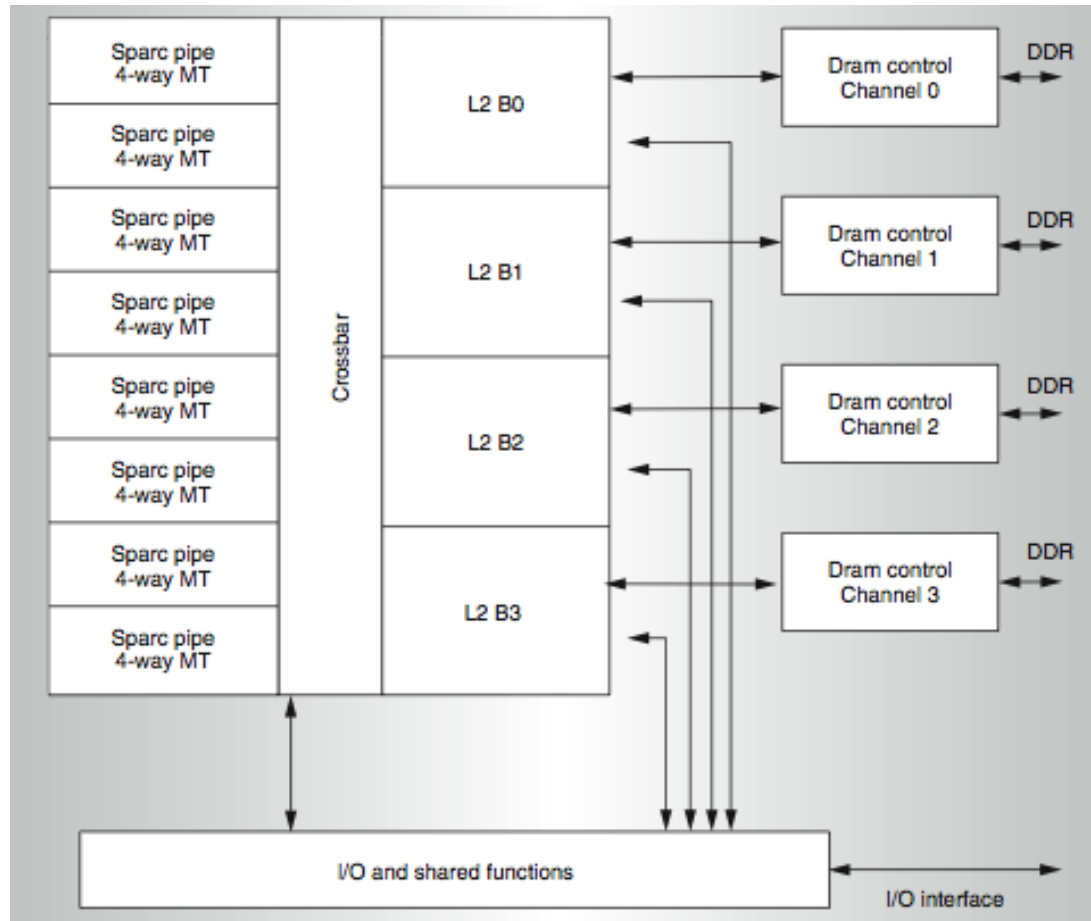


Figure 2. Block diagram of a single-chip Piranha I/O node.

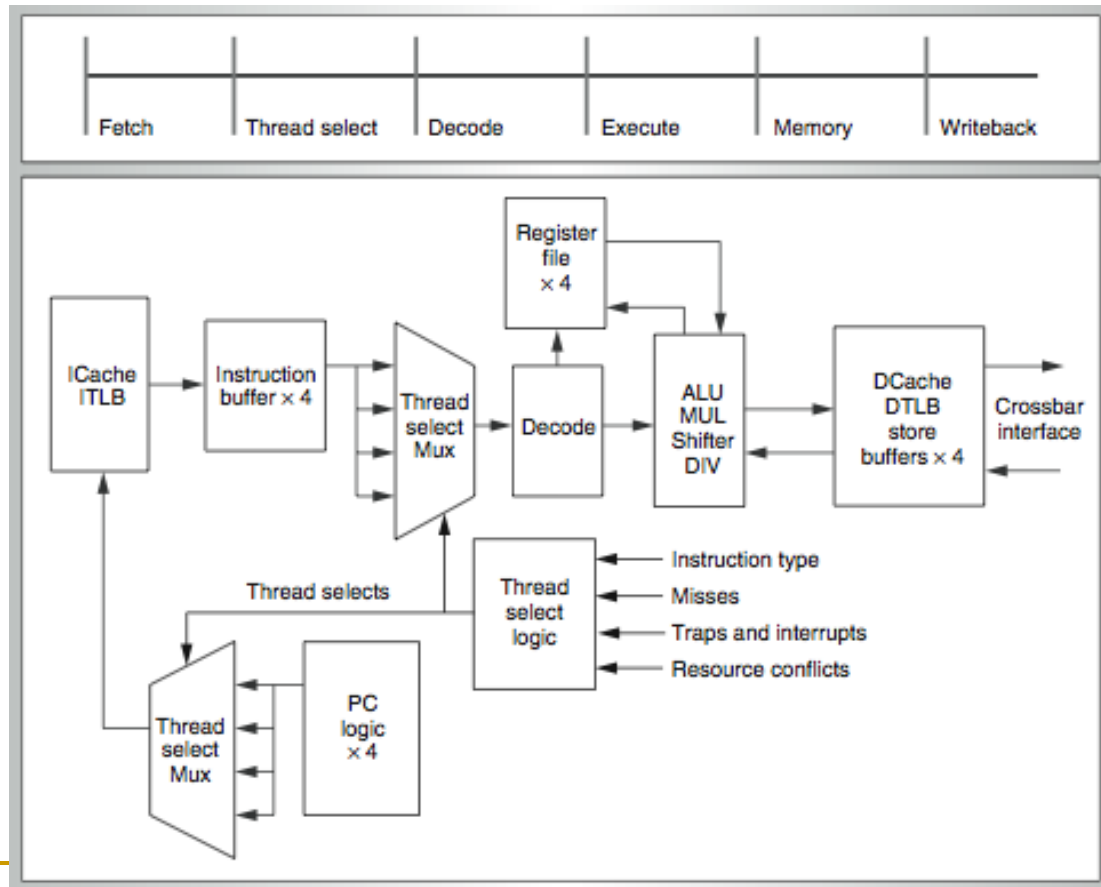
Meet Small: Sun Niagara (UltraSPARC T1)

- Kongetira et al., “Niagara: A 32-Way Multithreaded SPARC Processor,” IEEE Micro 2005.



Niagara Core

- 4-way fine-grain multithreaded, 6-stage, dual-issue in-order
- Round robin thread selection (unless cache miss)
- Shared FP unit among cores

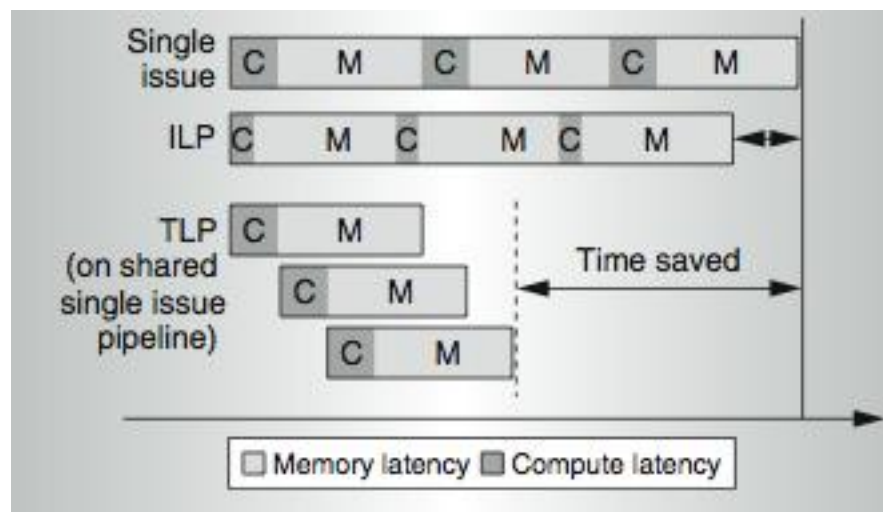


Niagara Design Point

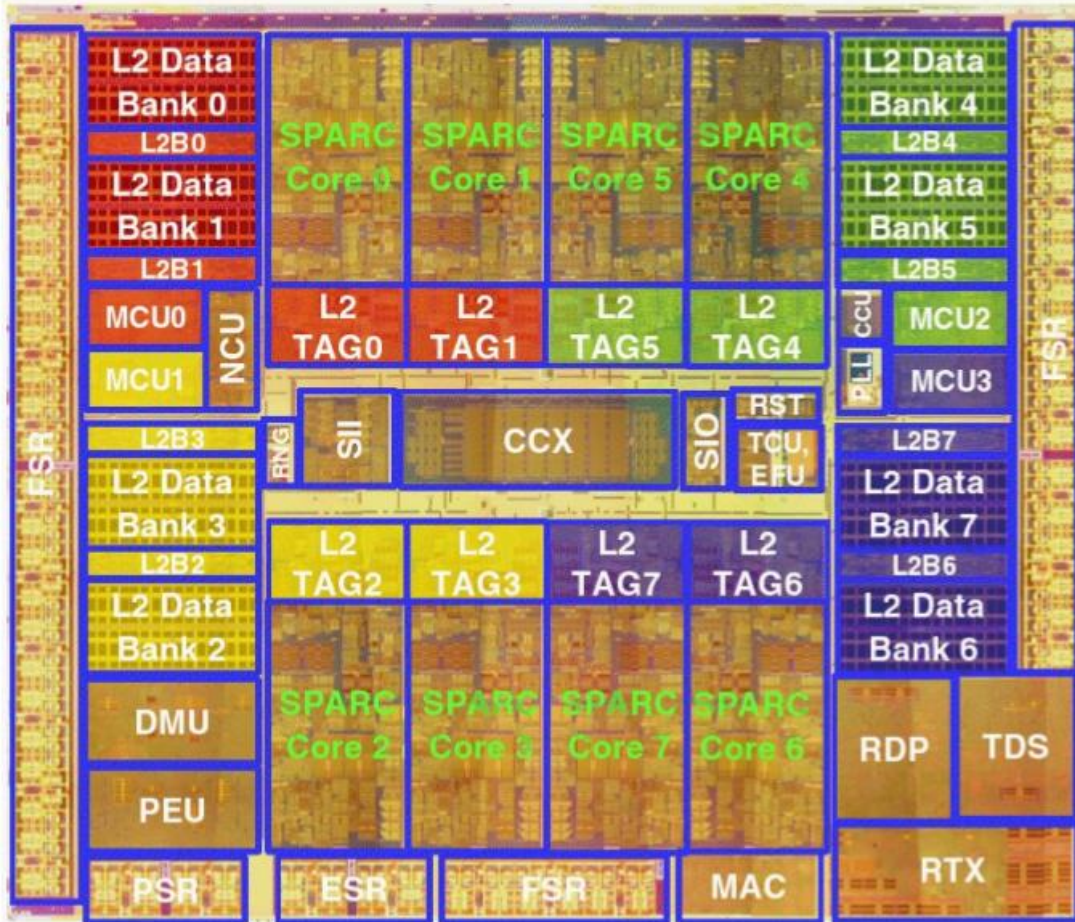
- Designed for commercial applications

Table 1. Commercial server applications.

Benchmark	Application category	Instruction-level parallelism	Thread-level parallelism	Working set	Data sharing
Web99	Web server	Low	High	Large	Low
JBB	Java application server	Low	High	Large	Medium
TPC-C	Transaction processing	Low	High	Large	High
SAP-2T	Enterprise resource planning	Medium	High	Medium	Medium
SAP-3T	Enterprise resource planning	Low	High	Large	High
TPC-H	Decision support system	High	High	Large	Medium



Meet Small: Sun Niagara II (UltraSPARC T2)

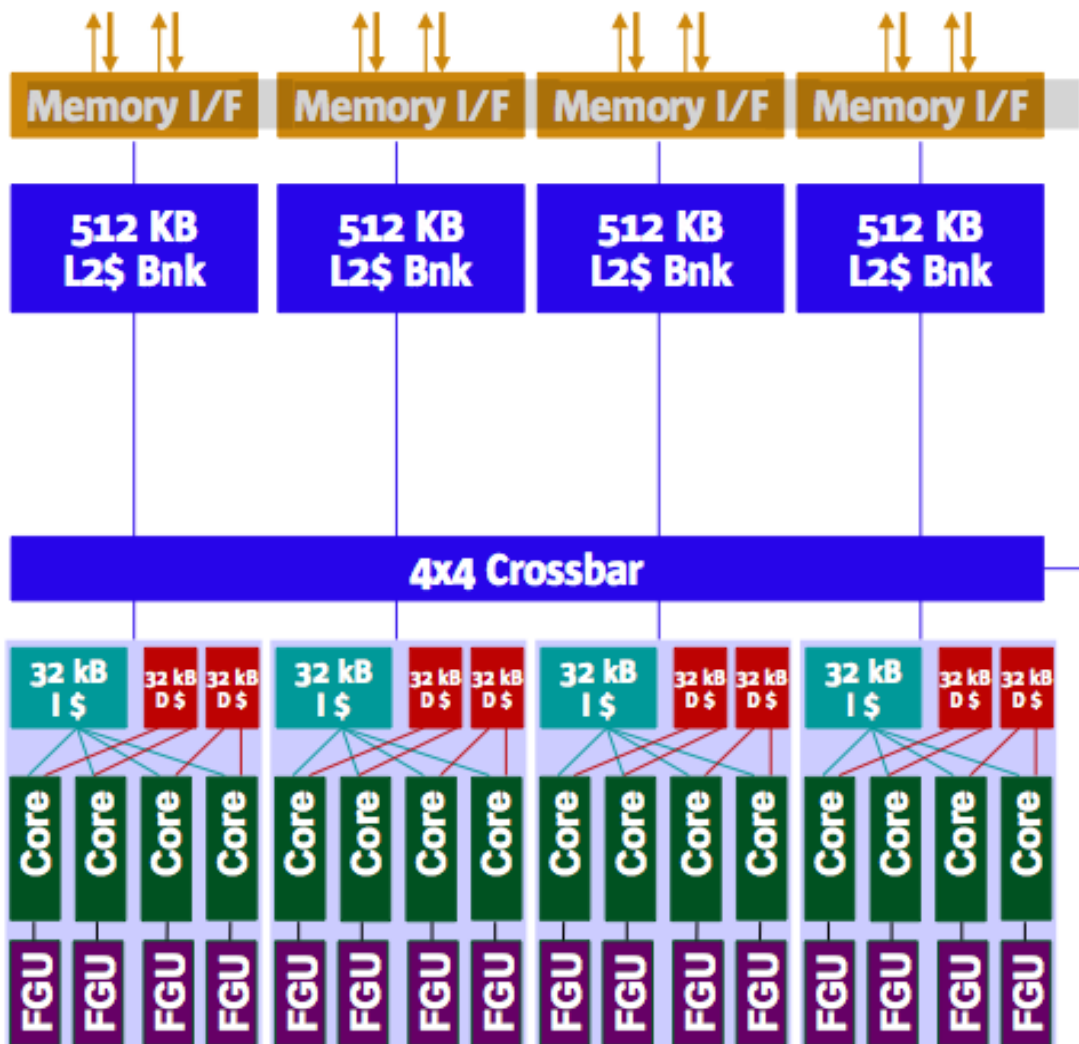


- 8 SPARC cores, 8 threads/core. 8 stages. 16 KB I\$ per Core. 8 KB D\$ per Core. FP, Graphics, Crypto, units per Core.
- 4 MB Shared L2, 8 banks, 16-way set associative.
- 4 dual-channel FBDIMM memory controllers.
- X8 PCI-Express @ 2.5 Gb/s.
- Two 10G Ethernet ports @ 3.125 Gb/s.

Meet Small, but Larger: Sun ROCK

- Chaudhry et al., “Rock: A High-Performance Sparc CMT Processor,” IEEE Micro, 2009.
- Chaudhry et al., “Simultaneous Speculative Threading: A Novel Pipeline Architecture Implemented in Sun's ROCK Processor,” ISCA 2009
- Goals:
 - Maximize throughput when threads are available
 - Boost single-thread performance when threads are not available and on cache misses
- Ideas:
 - Runahead on a cache miss → ahead thread executes miss-independent instructions, behind thread executes dependent instructions
 - Branch prediction (gshare)

Sun ROCK



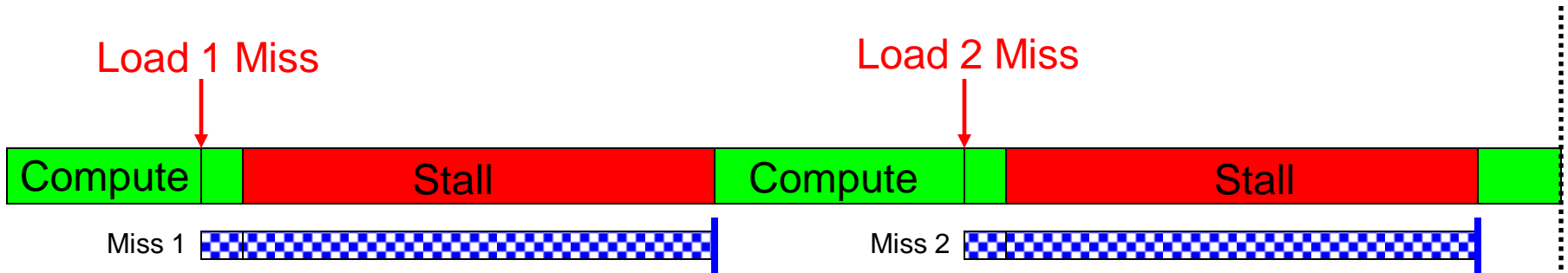
- 16 cores, 2 threads per core (fewer threads than Niagara 2)
- 4 cores share a 32KB instruction cache
- 2 cores share a 32KB data cache
- 2MB L2 cache (smaller than Niagara 2)

Runahead Execution (I)

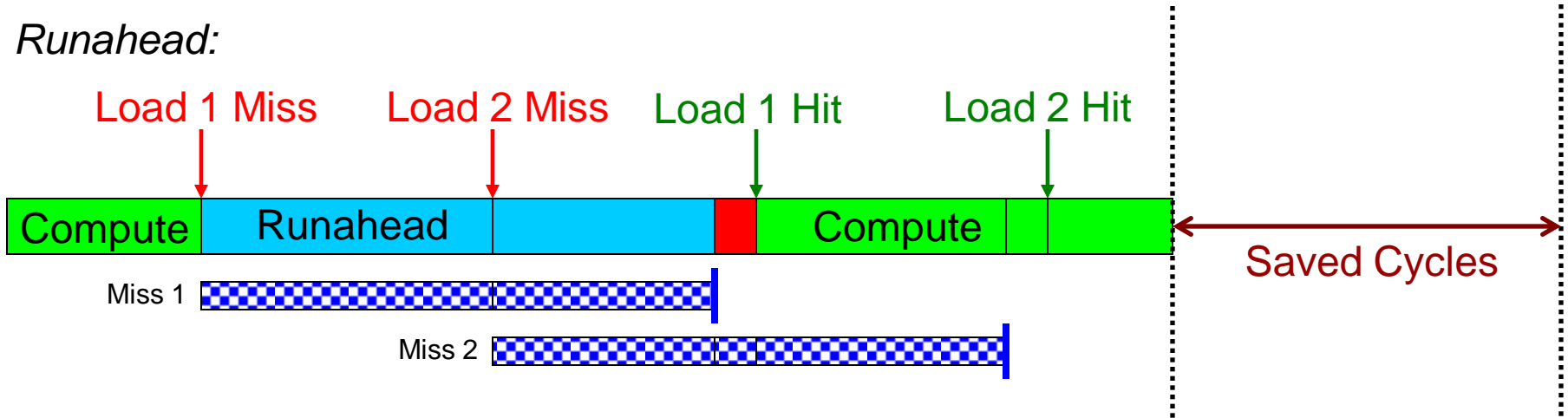
- A simple pre-execution method for prefetching purposes
- Mutlu et al., “Runahead Execution: An Alternative to Very Large Instruction Windows for Out-of-order Processors,” HPCA 2003, IEEE Micro 2003.
- When the oldest instruction is a long-latency cache miss:
 - Checkpoint architectural state and enter runahead mode
- In runahead mode:
 - **Speculatively pre-execute instructions**
 - **The purpose of pre-execution is to generate prefetches**
 - L2-miss dependent instructions are marked INV and dropped
- Runahead mode ends when the original miss returns
 - Checkpoint is restored and normal execution resumes

Runahead Execution (II)

Small Window:



Runahead:



Runahead Execution (III)

■ Advantages

- + Very **accurate** prefetches for data/instructions (all cache levels)
 - + Follows the program path
- + **Simple to implement**, most of the hardware is already built in

■ Disadvantages

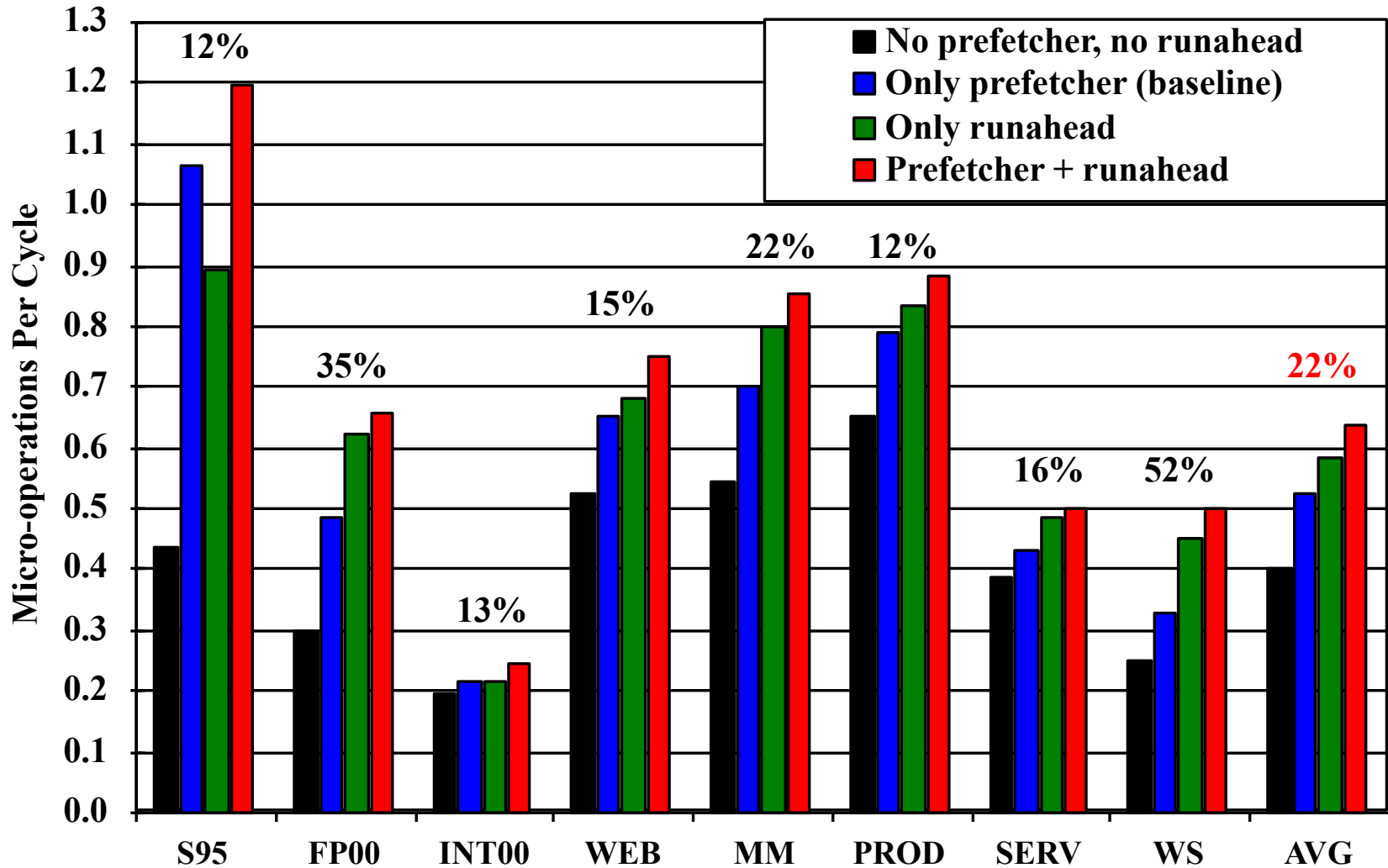
- **Extra executed instructions**

■ Limitations

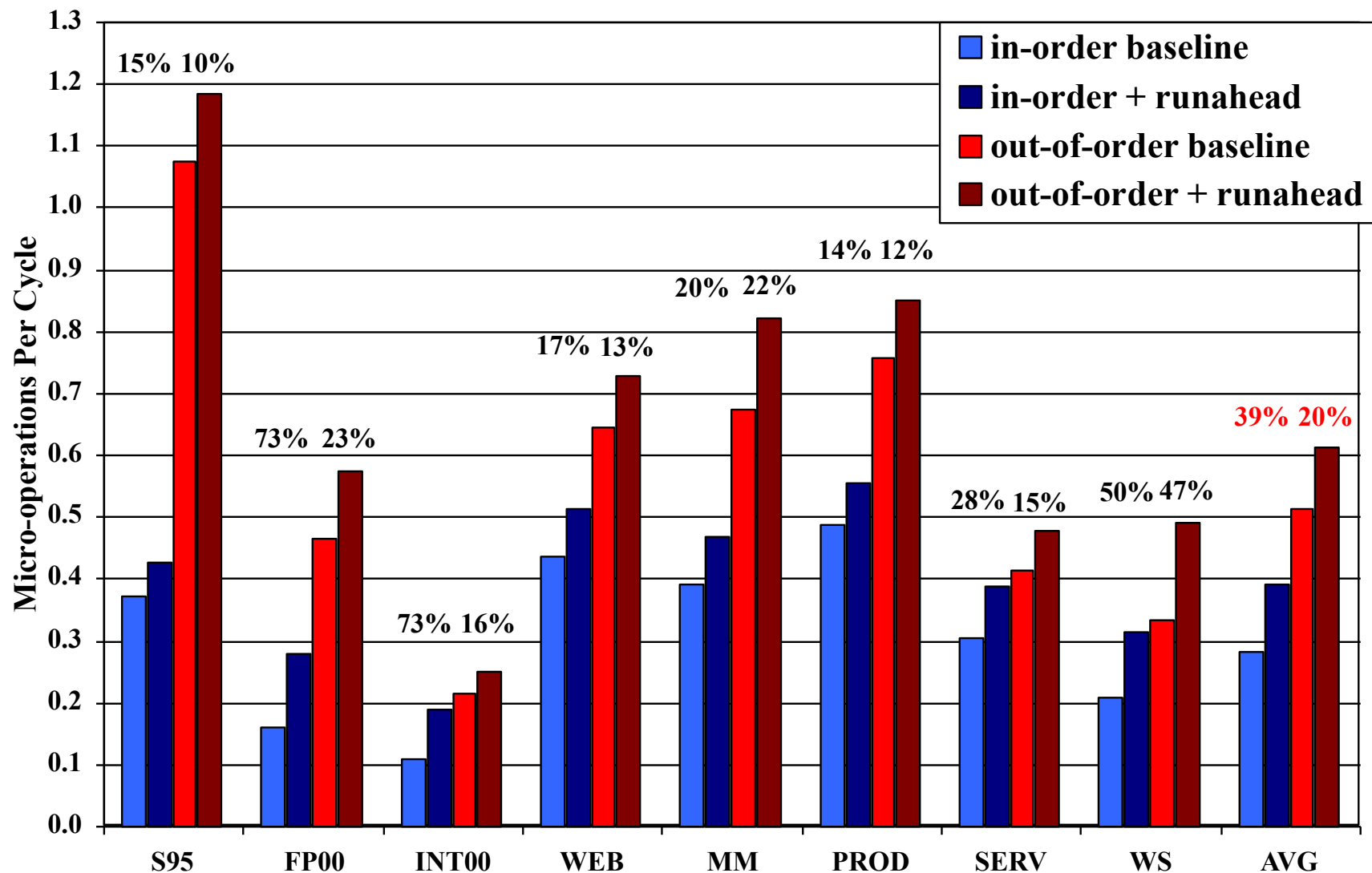
- Limited by branch prediction accuracy
- Cannot prefetch dependent cache misses. Solution?
- **Effectiveness limited by available Memory Level Parallelism**

- Mutlu et al., “**Efficient Runahead Execution: Power-Efficient Memory Latency Tolerance**,” IEEE Micro Jan/Feb 2006.
- Implemented in IBM POWER6, Sun ROCK

Performance of Runahead Execution

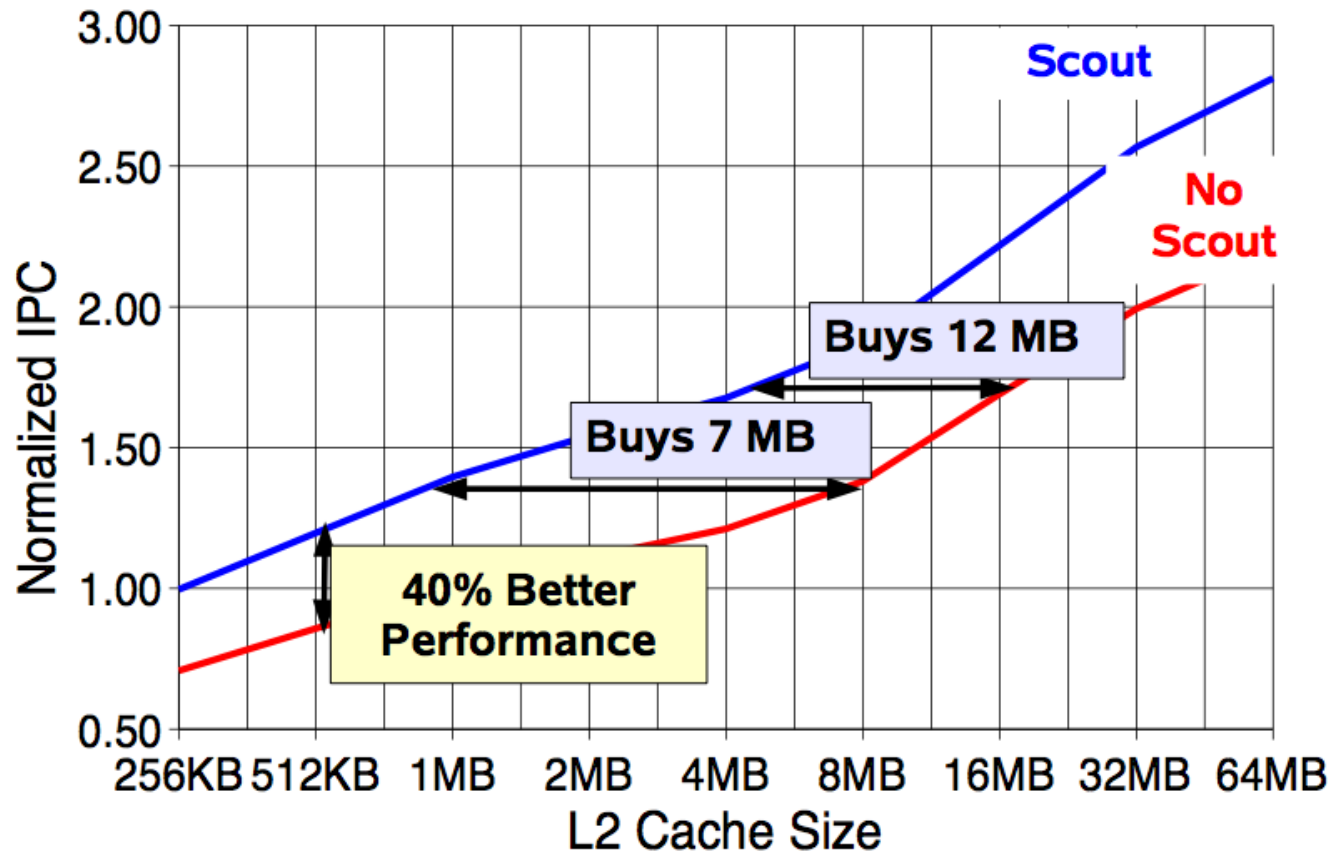


Performance of Runahead Execution (II)



More Powerful Cores in Sun ROCK

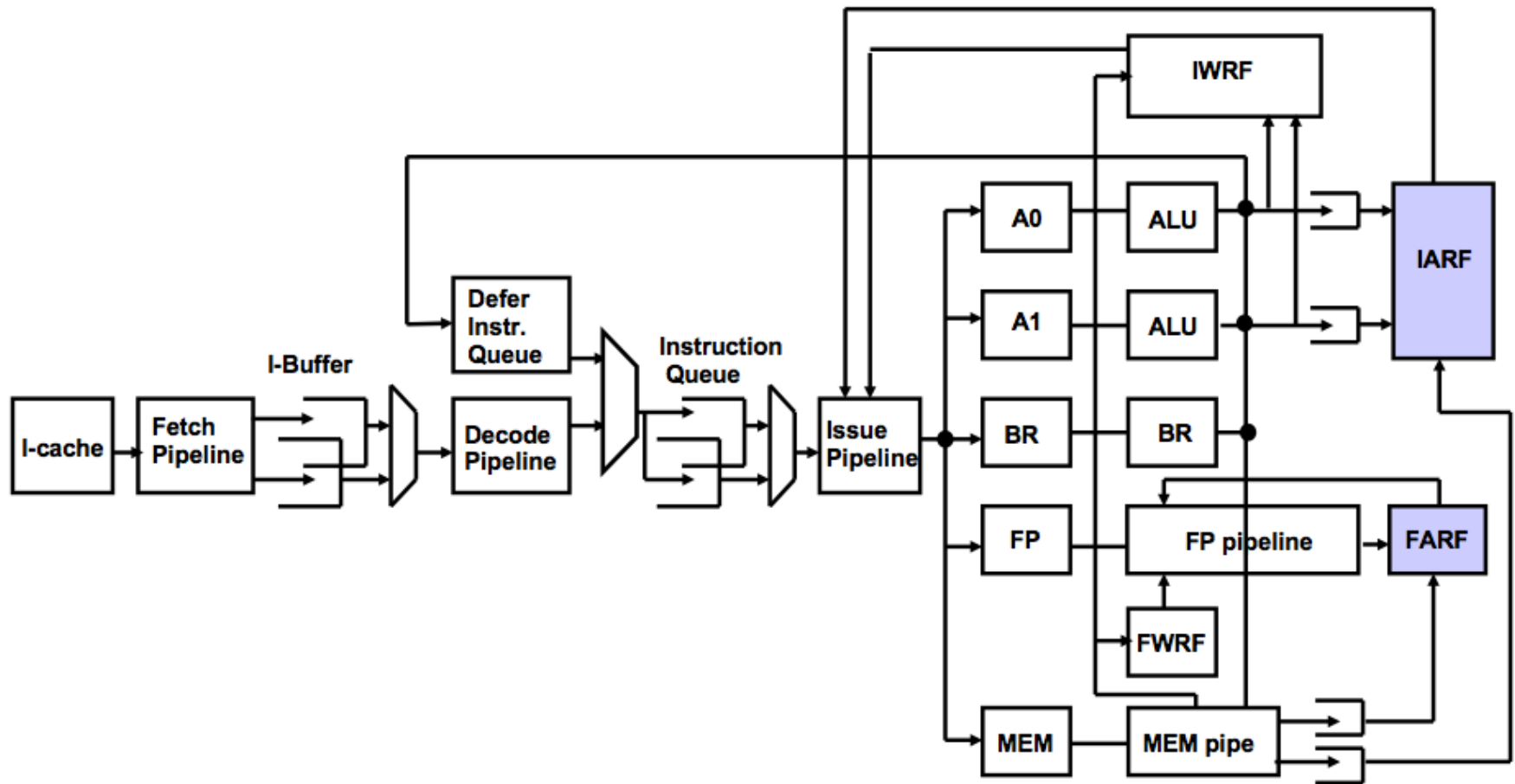
- Chaudhry talk, Aug 2008.



Sun ROCK Cores: Speculative Parallelization

- Load miss in L1 cache starts parallelization using 2 HW threads
- Ahead thread
 - Checkpoints state and executes speculatively
 - Speculatively executes instructions independent of the load miss
 - Defers load miss(es) and dependent instructions to the behind thread
- Behind thread
 - Executes deferred instructions and re-defers them if necessary
- Exploits Memory-Level Parallelism (MLP)
 - Run ahead on load miss and generate additional load misses
- Exploits Instruction-Level Parallelism (ILP)
 - Ahead and behind threads execute independent instructions from different points in program in parallel

ROCK Pipeline



More Powerful Cores in Sun ROCK

■ Advantages

- + Higher single-thread performance (MLP + ILP)
- + Better cache miss tolerance → Can reduce on-chip cache sizes

■ Disadvantages

- Bigger cores → Fewer cores → Lower parallel throughput (in terms of threads).
How about each thread's response time?
- More complex than Niagara cores (but simpler than conventional out-of-order execution) → Longer design time?

More Powerful Cores in Sun ROCK

- Chaudhry et al., “Simultaneous Speculative Threading: A Novel Pipeline Architecture Implemented in Sun's ROCK Processor,” ISCA 2009

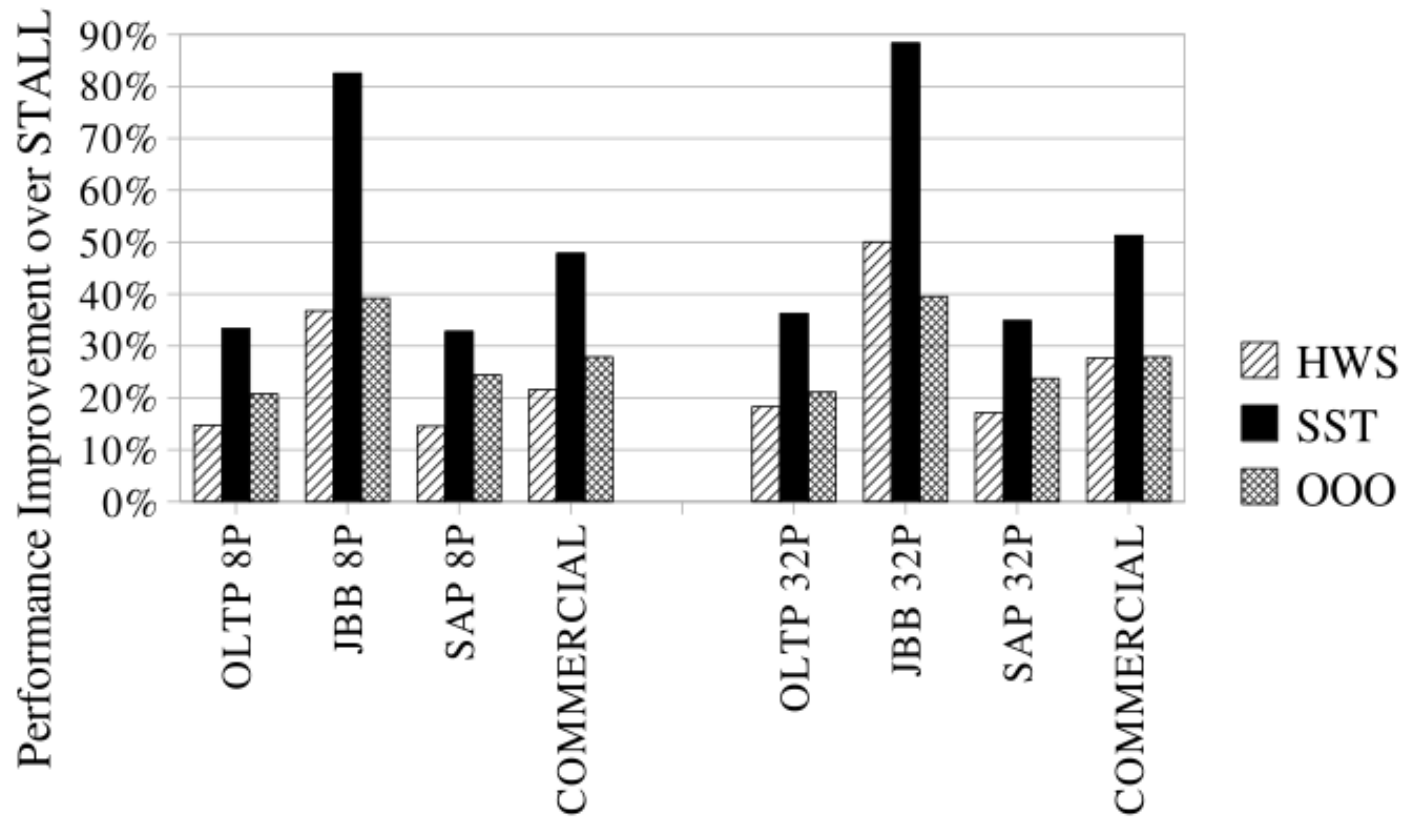
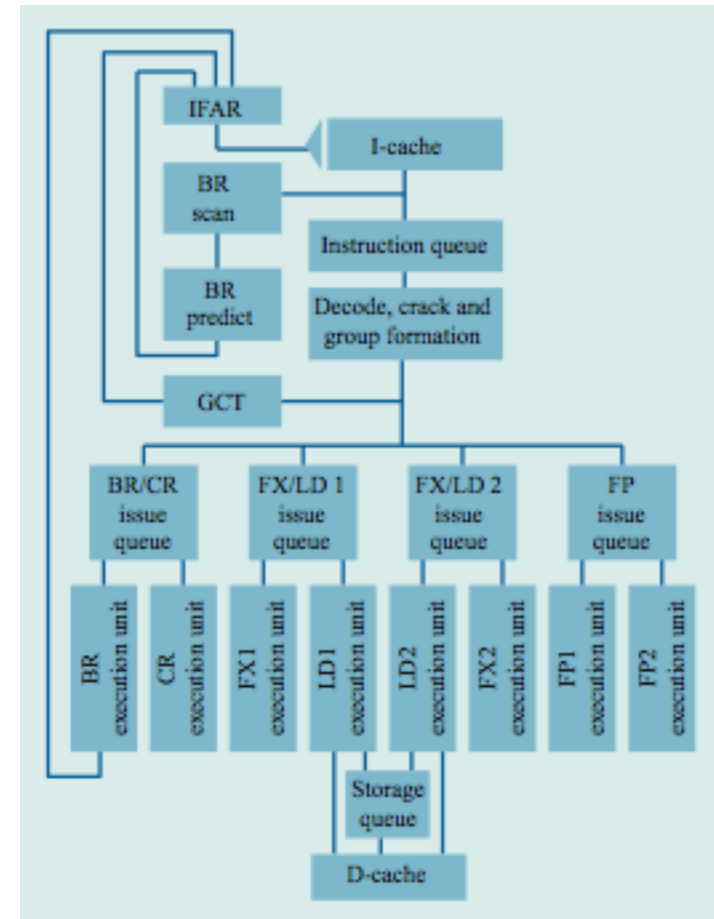
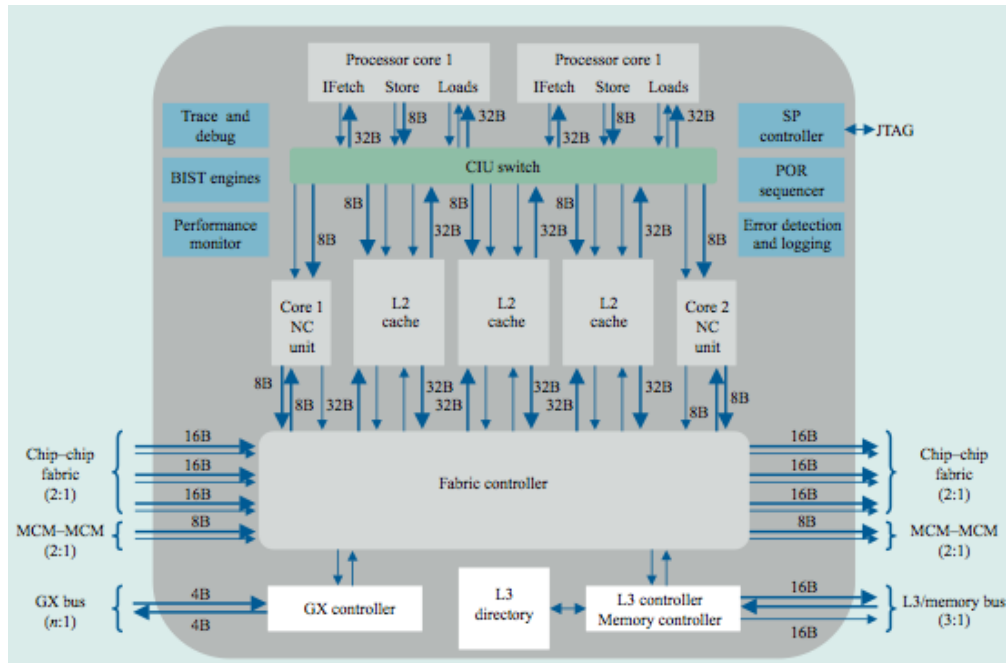


Figure 9: Commercial Performance.

Meet Large: IBM POWER4

- Tendler et al., “POWER4 system microarchitecture,” IBM J R&D, 2002.
- Another symmetric multi-core chip...
- But, fewer and more powerful cores



IBM POWER4

- 2 cores, out-of-order execution
- 100-entry instruction window in each core
- 8-wide instruction fetch, issue, execute
- Large, local+global hybrid branch predictor
- 1.5MB, 8-way L2 cache
- Aggressive stream based prefetching

IBM POWER5

- Kalla et al., “IBM Power5 Chip: A Dual-Core Multithreaded Processor,” IEEE Micro 2004.

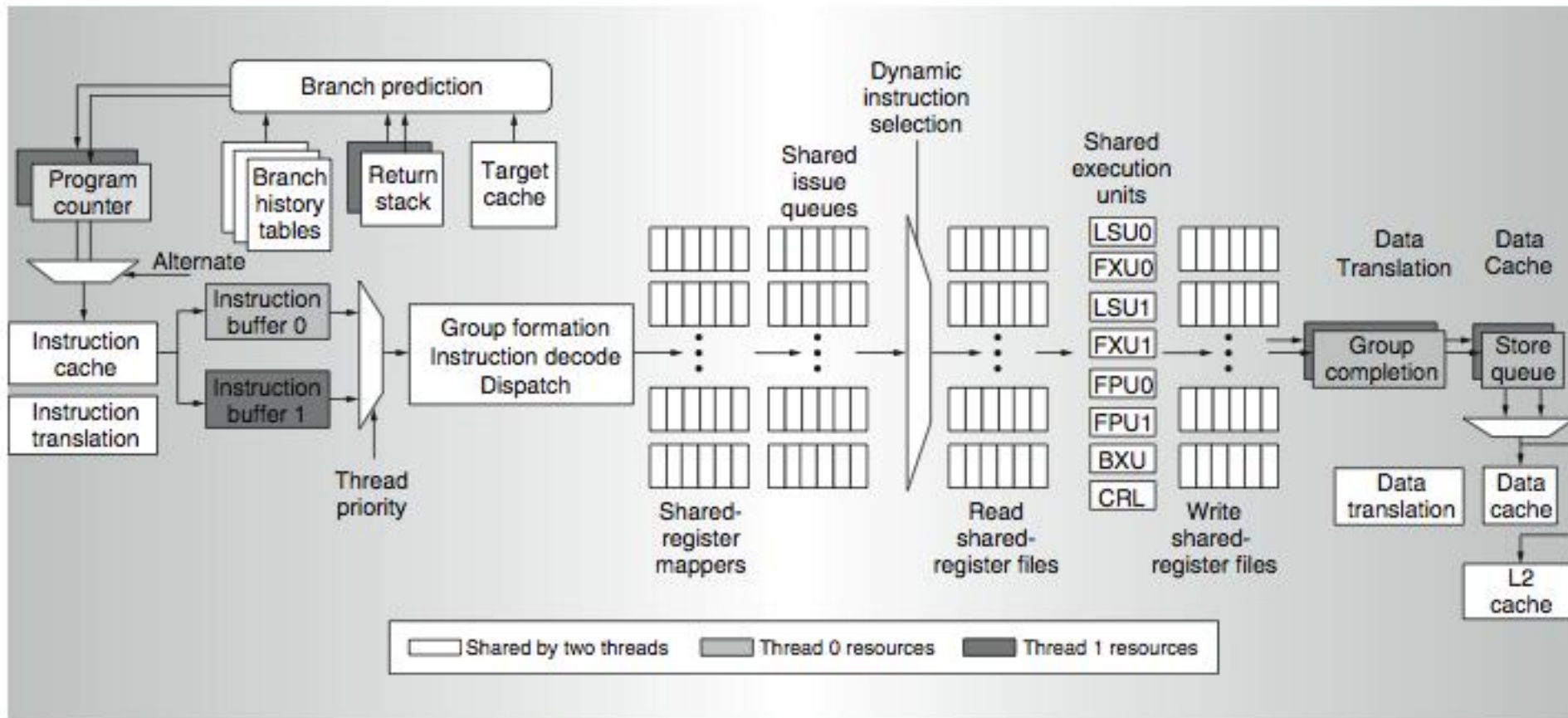
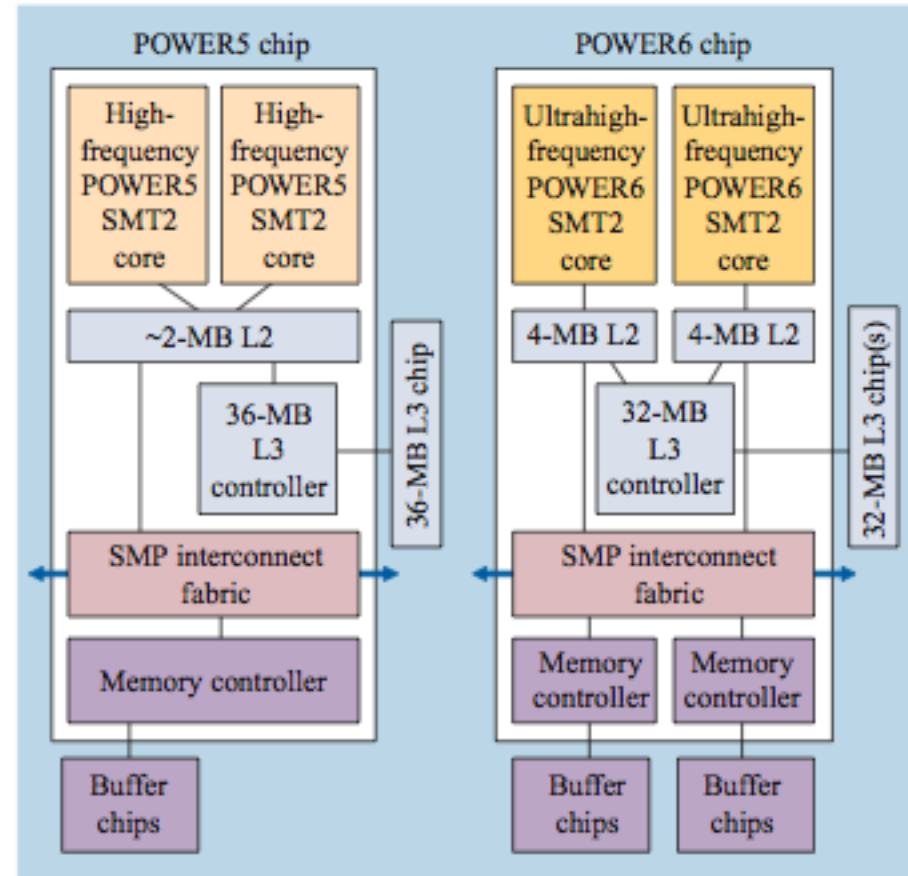


Figure 4. Power5 instruction data flow (BXU = branch execution unit and CRL = condition register logical execution unit).

Meet Large, but Smaller: IBM POWER6

- Le et al., “**IBM POWER6 microarchitecture**,” IBM J R&D, 2007.
- 2 cores, in order, high frequency (4.7 GHz)
- 8 wide fetch
- Simultaneous multithreading in each core
- Runahead execution in each core
 - Similar to Sun ROCK



IBM POWER7

- Kalla et al., “Power7: IBM’s Next-Generation Server Processor,” IEEE Micro 2010.
- 8 out-of-order cores, 4-way SMT in each core
- TurboCore mode
 - Can turn off cores so that other cores can be run at higher frequency

Remember the Demands

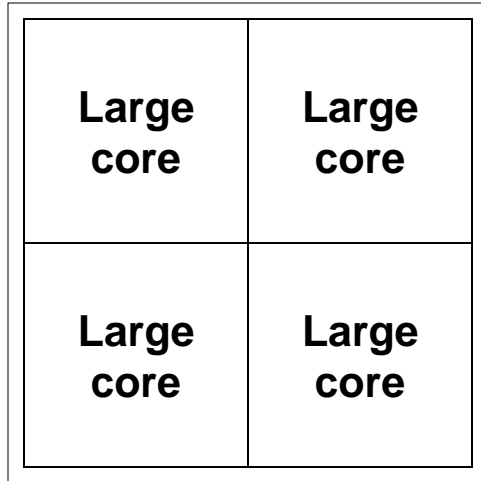
- What we want:
- In a serialized code section → one powerful “large” core
- In a parallel code section → many wimpy “small” cores
- These two conflict with each other:
 - If you have a single powerful core, you cannot have many cores
 - A small core is much more energy and area efficient than a large core
- Can we get the best of both worlds?

Performance vs. Parallelism

Assumptions:

- 1. Small cores takes an area budget of 1 and has performance of 1*
- 2. Large core takes an area budget of 4 and has performance of 2*

Tile-Large Approach



“Tile-Large”

- Tile a few large cores
- IBM Power 5, AMD Barcelona, Intel Core2Quad, Intel Nehalem
- + High performance on single thread, serial code sections (2 units)
- Low throughput on parallel program portions (8 units)

Tile-Small Approach

Small core	Small core	Small core	Small core
Small core	Small core	Small core	Small core
Small core	Small core	Small core	Small core
Small core	Small core	Small core	Small core

“Tile-Small”

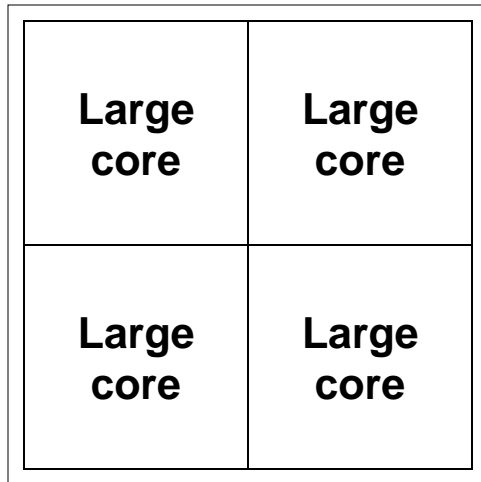
- Tile many small cores
- Sun Niagara, Intel Larrabee, Tiler TILE (tile ultra-small)
 - + High throughput on the parallel part (16 units)
 - Low performance on the serial part, single thread (1 unit)

Can we get the best of both worlds?

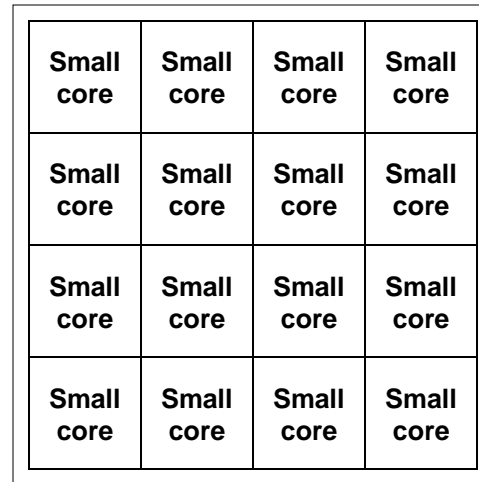
- Tile Large
 - + High performance on single thread, serial code sections (2 units)
 - Low throughput on parallel program portions (8 units)
- Tile Small
 - + High throughput on the parallel part (16 units)
 - Low performance on the serial part, single thread (1 unit),
reduced single-thread performance compared to existing single thread processors
- Idea: Have both large and small on the same chip →
Performance asymmetry

Asymmetric Multi-Core

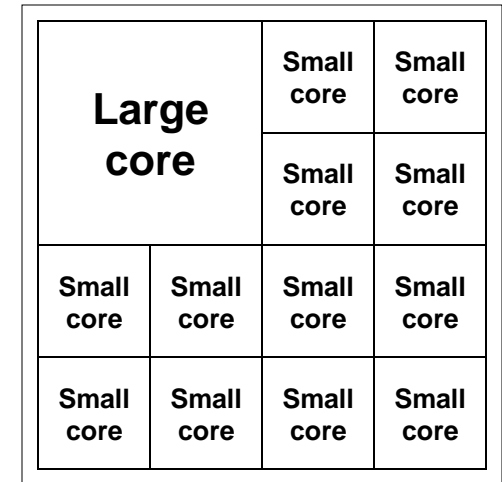
Asymmetric Chip Multiprocessor (ACMP)



“Tile-Large”



“Tile-Small”

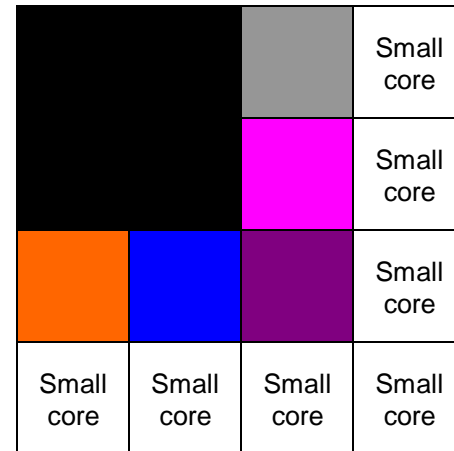
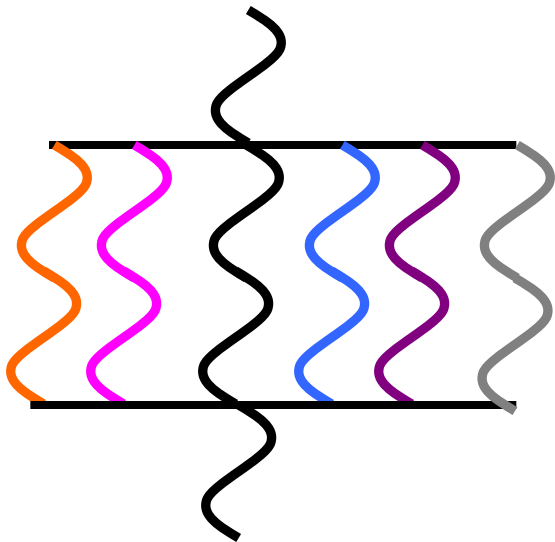


ACMP

- Provide one large core and many small cores
- + Accelerate serial part using the large core (2 units)
- + Execute parallel part on small cores and large core for high throughput (12+2 units)

Accelerating Serial Bottlenecks

Single thread \rightarrow Large core



ACMP Approach

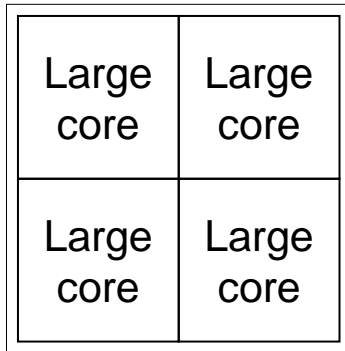
Performance vs. Parallelism

Assumptions:

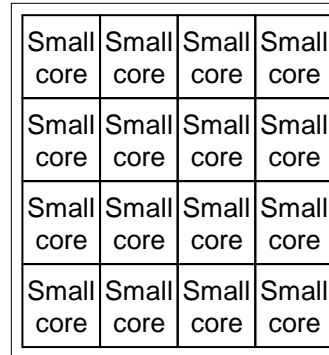
- 1. Small cores takes an area budget of 1 and has performance of 1*
- 2. Large core takes an area budget of 4 and has performance of 2*

ACMP Performance vs. Parallelism

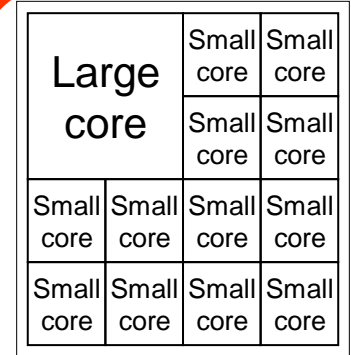
Area-budget = 16 small cores



“Tile-Large”



“Tile-Small”

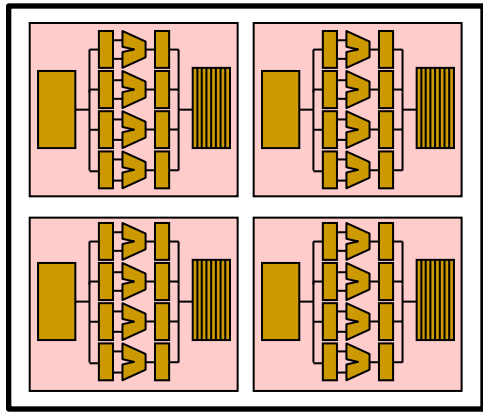


ACMP

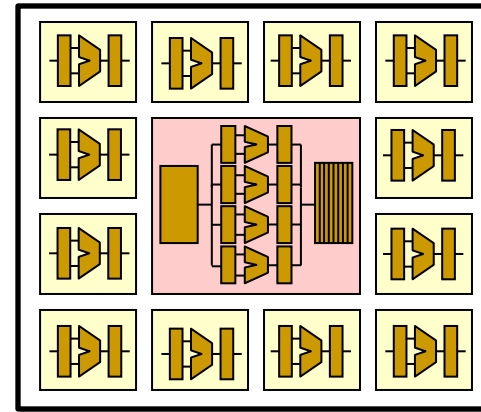
Large Cores	4	0	1
Small Cores	0	16	12
Serial Performance	2	1	2
Parallel Throughput	$2 \times 4 = 8$	$1 \times 16 = 16$	$1 \times 2 + 1 \times 12 = 14$

Some Analysis

- Hill and Marty, “Amdahl’s Law in the Multi-Core Era,” IEEE Computer 2008.
- **Each Chip** Bounded to **N** BCEs (Base Core Equivalents)
- **One R-BCE Core** leaves $N - R$ BCEs
- **Use $N - R$ BCEs for $N - R$ Base Cores**
- Therefore, **$1 + N - R$ Cores per Chip**
- For an $N = 16$ BCE Chip:



Symmetric: Four 4-BCE cores



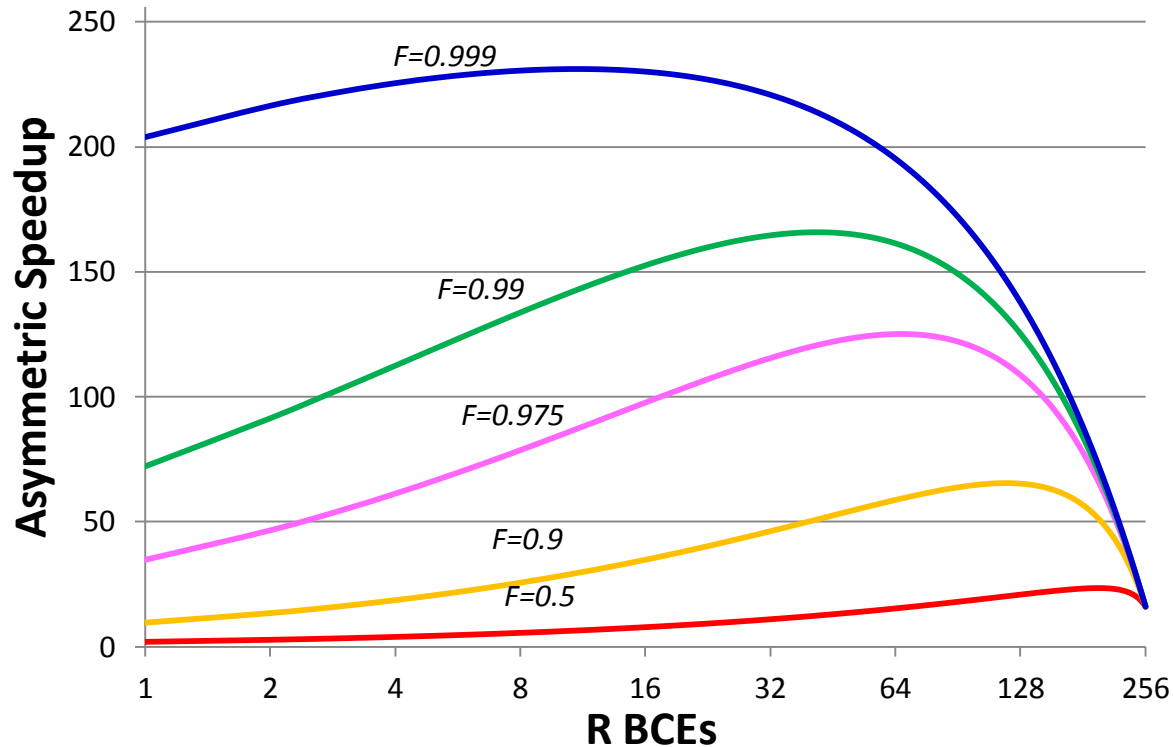
**Asymmetric: One 4-BCE core
& Twelve 1-BCE base cores**

Amdahl's Law Modified

- Serial Fraction $1-F$ same, so time = $(1 - F) / \text{Perf}(R)$
- Parallel Fraction F
 - One core at rate $\text{Perf}(R)$
 - $N-R$ cores at rate 1
 - Parallel time = $F / (\text{Perf}(R) + N - R)$
- Therefore, w.r.t. one base core:

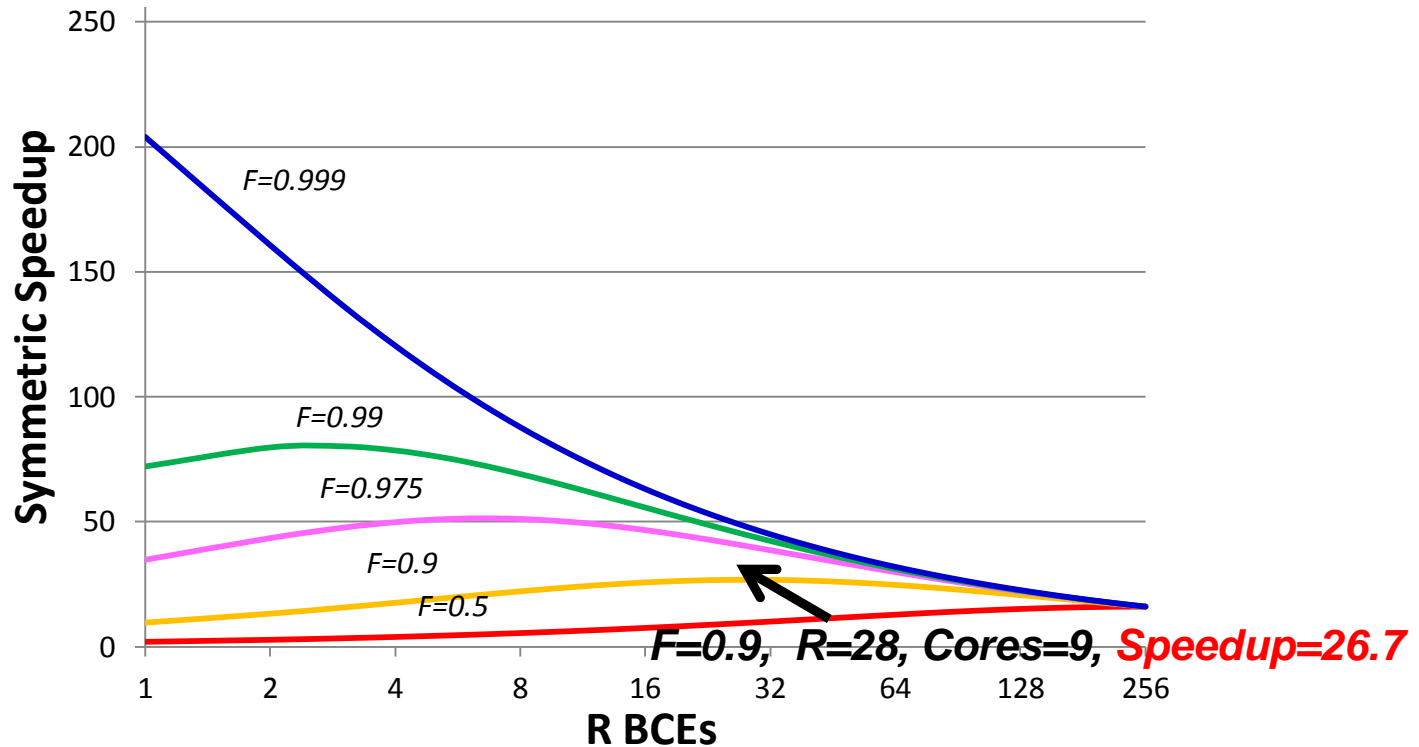
$$\text{Asymmetric Speedup} = \frac{1}{\frac{1 - F}{\text{Perf}(R)} + \frac{F}{\text{Perf}(R) + N - R}}$$

Asymmetric Multicore Chip, $N = 256$ BCEs

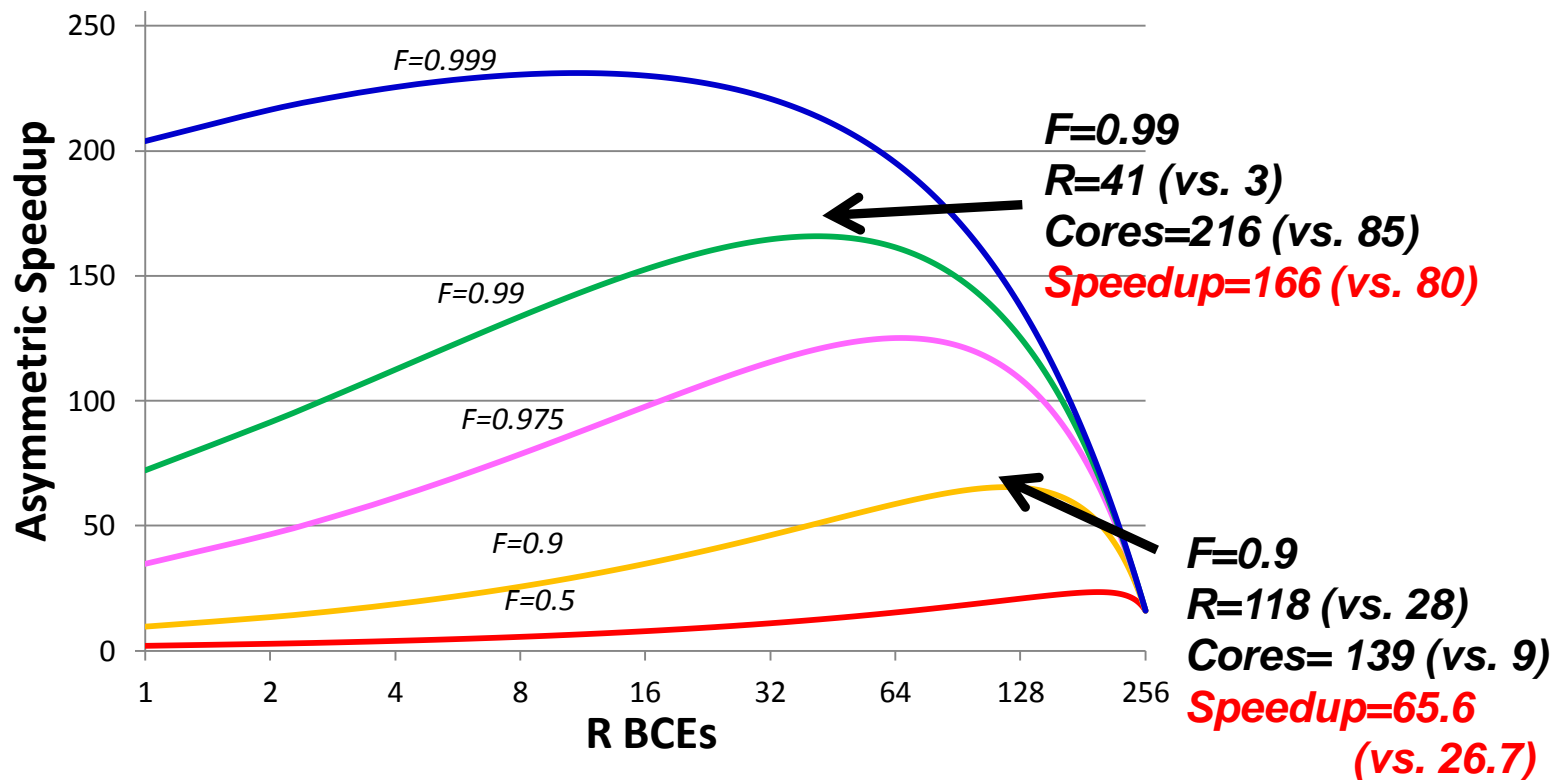


- Number of Cores = 1 (Enhanced) + 256 - R (Base)

Symmetric Multicore Chip, $N = 256$ BCEs



Asymmetric Multicore Chip, $N = 256$ BCEs



- Asymmetric multi-core provides better speedup than symmetric multi-core when N is large

Asymmetric vs. Symmetric Cores

■ Advantages of Asymmetric

- + Can provide better performance when thread parallelism is limited
- + Can be more energy efficient
 - + Schedule computation to the core type that can best execute it

■ Disadvantages

- Need to design more than one type of core. Always?
- Scheduling becomes more complicated
 - What computation should be scheduled on the large core?
 - Who should decide? HW vs. SW?
- Managing locality and load balancing can become difficult if threads move between cores (transparently to software)
- Cores have different demands from shared resources

Caveats of Parallelism, Revisited

■ Amdahl's Law

- f: Parallelizable fraction of a program
- N: Number of processors

$$\text{Speedup} = \frac{1}{1 - f + \frac{f}{N}}$$

- Amdahl, “Validity of the single processor approach to achieving large scale computing capabilities,” AFIPS 1967.
- **Maximum speedup limited by serial portion: Serial bottleneck**
- **Parallel portion is usually not perfectly parallel**
 - **Synchronization** overhead (e.g., updates to shared data)
 - **Load imbalance** overhead (imperfect parallelization)
 - **Resource sharing** overhead (contention among N processors)

Accelerating Parallel Bottlenecks

- Serialized or imbalanced execution in the parallel portion can also benefit from a large core
- Examples:
 - Critical sections that are contended
 - Parallel stages that take longer than others to execute
- Idea: Dynamically identify these code portions that cause serialization and execute them on a large core
- Suleman et al., “Accelerating Critical Section Execution with Asymmetric Multi-Core Architectures,” ASPLOS 2009, IEEE Micro Top Picks 2010.
- Joao et al., “Bottleneck Identification and Scheduling,” ASPLOS 2012.

How to Achieve Asymmetry

■ Static

- ❑ Type and power of cores fixed at design time
- ❑ Two approaches to design “faster cores”:
 - High frequency
 - Build a more complex, powerful core with entirely different arch
- ❑ Is static asymmetry natural? (chip-wide variations in frequency)

■ Dynamic

- ❑ Type and power of cores change dynamically
- ❑ Two approaches to dynamically create “faster cores”:
 - Boost frequency dynamically (limited power budget)
 - Combine small cores to enable a more complex, powerful core
- ❑ Is there a third, fourth, fifth approach?

Asymmetry via Boosting of Frequency

■ Static

- Due to process variations, cores might have different frequency
- Simply hardwire/design cores to have different frequencies

■ Dynamic

- Annavaram et al., “[Mitigating Amdahl’s Law Through EPI Throttling](#),” ISCA 2005.
- Dynamic voltage and frequency scaling

EPI Throttling

- Goal: Minimize execution time of parallel programs while keeping power within a fixed budget
- For best scalar and throughput performance, vary energy expended per instruction (EPI) based on available parallelism
 - $P = \text{EPI} \times \text{IPS}$
 - $P = \text{fixed power budget}$
 - EPI = energy per instruction
 - IPS = aggregate instructions retired per second
- Idea: For a fixed power budget
 - Run sequential phases on high-EPI processor
 - Run parallel phases on multiple low-EPI processors

EPI Throttling via DVFS

- DVFS: Dynamic voltage frequency scaling
- In phases of low thread parallelism
 - Run a few cores at high supply voltage and high frequency
- In phases of high thread parallelism
 - Run many cores at low supply voltage and low frequency

Possible EPI Throttling Techniques

- Grochowski et al., “Best of both Latency and Throughput,” ICCD 2004.

Method	EPI Range	Time to Alter EPI	Throttle Action
Voltage/frequency scaling	1:2 to 1:4	100us (ramp Vcc)	Lower voltage and frequency
Asymmetric cores	1:4 to 1:6	10us (migrate 256KB L2 cache)	Migrate threads from large cores to small cores
Variable-size core	1:1 to 1:2	1us (fill 32KB L1 cache)	Reduce capacity of processor resources
Speculation control	1:1 to 1:1.4	10ns (pipeline latency)	Reduce amount of speculation

Boosting Frequency of a Small Core vs. Large Core

- Frequency boosting implemented on Intel Nehalem, IBM POWER7
- Advantages of Boosting Frequency
 - + Very simple to implement; no need to design a new core
 - + Parallel throughput does not degrade when TLP is high
 - + Preserves locality of boosted thread
- Disadvantages
 - Does not improve performance if thread is memory bound
 - Does not reduce Cycles per Instruction (remember the performance equation?)
 - Changing frequency/voltage can take longer than switching to a large core

Uses of Asymmetry

- So far:
 - Improvement in serial performance (sequential bottleneck)
- What else can we do with asymmetry?
 - Energy reduction?
 - Energy/performance tradeoff?
 - Improvement in parallel portion?

Use of Asymmetry for Energy Efficiency

- Kumar et al., “Single-ISA Heterogeneous Multi-Core Architectures: The Potential for Processor Power Reduction,” MICRO 2003.
- Idea:
 - Implement multiple types of cores on chip
 - Monitor characteristics of the running thread
 - e.g., sample energy/perf on each core periodically
 - Dynamically pick the core that provides the best energy/performance tradeoff for a given phase
 - “Best core” → Depends on optimization metric
- Example: ARM’s big.LITTLE architecture

Use of Asymmetry for Energy Efficiency

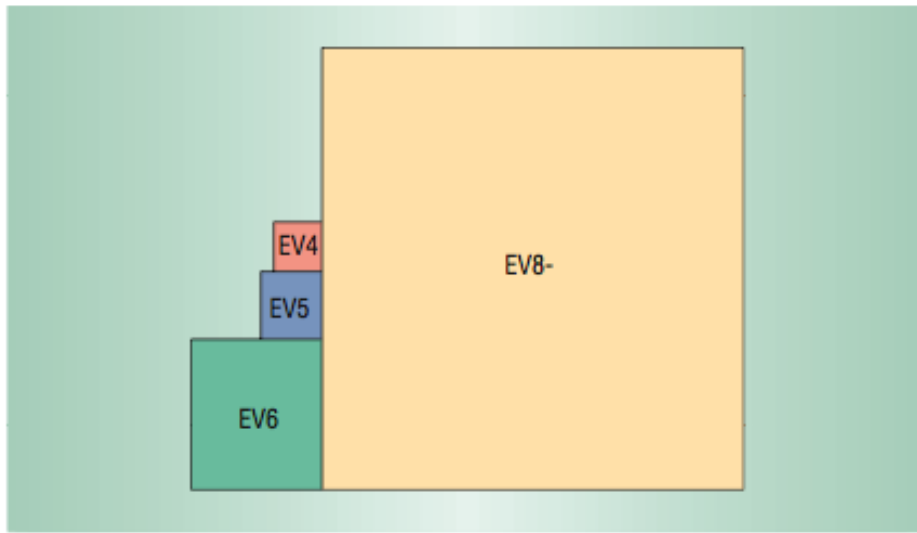
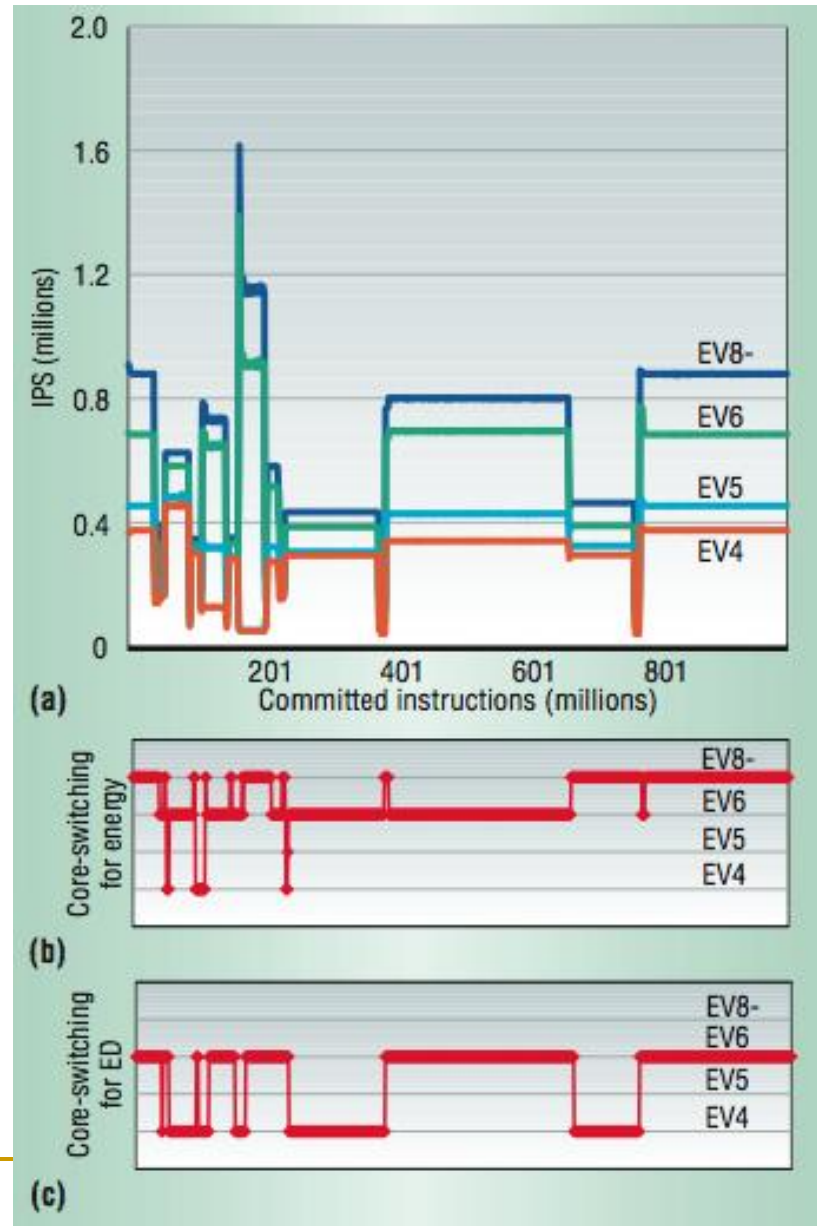


Figure 1. Relative sizes of the Alpha cores scaled to 0.10 μm . EV8 is 80 times bigger but provides only two to three times more single-threaded performance.

Table 1. Power and relative performance of Alpha cores scaled to 0.10 μm . Performance is expressed normalized to EV4 performance.

Core	Peak power (Watts)	Average power (Watts)	Performance (norm. IPC)
EV4	4.97	3.73	1.00
EV5	9.83	6.88	1.30
EV6	17.8	10.68	1.87
EV8	92.88	46.44	2.14



Use of Asymmetry for Energy Efficiency

■ Advantages

- + More flexibility in energy-performance tradeoff
- + Can execute computation to the core that is best suited for it (in terms of energy)

■ Disadvantages/issues

- Incorrect predictions/sampling → wrong core → reduced performance or increased energy
- Overhead of core switching
- Disadvantages of asymmetric CMP (e.g., design multiple cores)
- Need phase monitoring and matching algorithms
 - What characteristics should be monitored?
 - Once characteristics known, how do you pick the core?

Functional vs. Performance Asymmetry

- Functional asymmetry: Place on chip multiple cores with different ISAs/interfaces
- Examples
 - CPU+GPU architectures (Intel Sandybridge, AMD APU, Nvidia Tegra)
 - SoC's with different accelerators (e.g., Qualcomm)
- Example: Nvidia Tegra
 - 72-core GPU
 - 4-core ARM processor

Summary: Multi-Core Evolution

- Symmetric Multi-core
 - Evolution of Sun's and IBM's Multicore systems and design choices
 - Niagara, Niagara 2, ROCK
 - IBM POWERx

- Asymmetric Multi-core
 - Motivation
 - Static vs. Dynamic Asymmetry
 - EPI Throttling
 - Use of Asymmetry for Energy Efficiency
 - Functional vs. Performance Asymmetry

Computer Architecture: Multi-Core Evolution and Design

Prof. Onur Mutlu
Carnegie Mellon University

Backup Slides

Referenced Readings (I)

- Grochowski et al., “Best of both Latency and Throughput,” ICCD 2004.
- Barroso et al., “Piranha: A Scalable Architecture Based on Single-Chip Multiprocessing,” ISCA 2000.
- Barroso et al., “Memory System Characterization of Commercial Workloads,” ISCA 1998.
- Ranganathan et al., “Performance of Database Workloads on Shared-Memory Systems with Out-of-Order Processors,” ASPLOS 1998.
- Kongetira et al., “Niagara: A 32-Way Multithreaded SPARC Processor,” IEEE Micro 2005.
- Chaudhry et al., “Rock: A High-Performance Sparc CMT Processor,” IEEE Micro, 2009.
- Chaudhry et al., “Simultaneous Speculative Threading: A Novel Pipeline Architecture Implemented in Sun's ROCK Processor,” ISCA 2009
- Mutlu et al., “Efficient Runahead Execution: Power-Efficient Memory Latency Tolerance,” IEEE Micro Jan/Feb 2006.
- Mutlu et al., “Runahead Execution,” HPCA 2003.

Referenced Readings (II)

- Tandler et al., “POWER4 system microarchitecture,” IBM J R&D, 2002.
- Kalla et al., “IBM Power5 Chip: A Dual-Core Multithreaded Processor,” IEEE Micro 2004.
- Le et al., “IBM POWER6 microarchitecture,” IBM J R&D, 2007.
- Kalla et al., “Power7: IBM’s Next-Generation Server Processor,” IEEE Micro 2010.
- Hill and Marty, “Amdahl’s Law in the Multi-Core Era,” IEEE Computer 2008.
- Suleman et al., “Accelerating Critical Section Execution with Asymmetric Multi-Core Architectures,” ASPLOS 2009, IEEE Micro Top Picks 2010.
- Joao et al., “Bottleneck Identification and Scheduling,” ASPLOS 2012.
- Annavaram et al., “Mitigating Amdahl’s Law Through EPI Throttling,” ISCA 2005.
- Kumar et al., “Single-ISA Heterogeneous Multi-Core Architectures: The Potential for Processor Power Reduction,” MICRO 2003.
- Amdahl, “Validity of the single processor approach to achieving large scale computing capabilities,” AFIPS 1967.

Related Videos

- Multi-Core Systems and Heterogeneity
 - <http://www.youtube.com/watch?v=LIDxT0hPI2U&list=PLVngZ7BemHHV6N0ejHhwOfLwTr8Q-UKXj&index=1>
 - <http://www.youtube.com/watch?v=Q0zyLVnzkrM&list=PLVngZ7BemHHV6N0ejHhwOfLwTr8Q-UKXj&index=2>

More on EPI Throttling

EPI Throttling (Annavaram et al., ISCA' 05)

■ Static AMP

- ❑ Duty cycles set once prior to program run
- ❑ Parallel phases run on 3P/1.25GHz
- ❑ Sequential phases run on 1P/2GHz
- ❑ Affinity guarantees sequential on 1P and parallel on 3
- ❑ Benchmarks that rapidly transition between sequential and parallel phases

■ Dynamic AMP

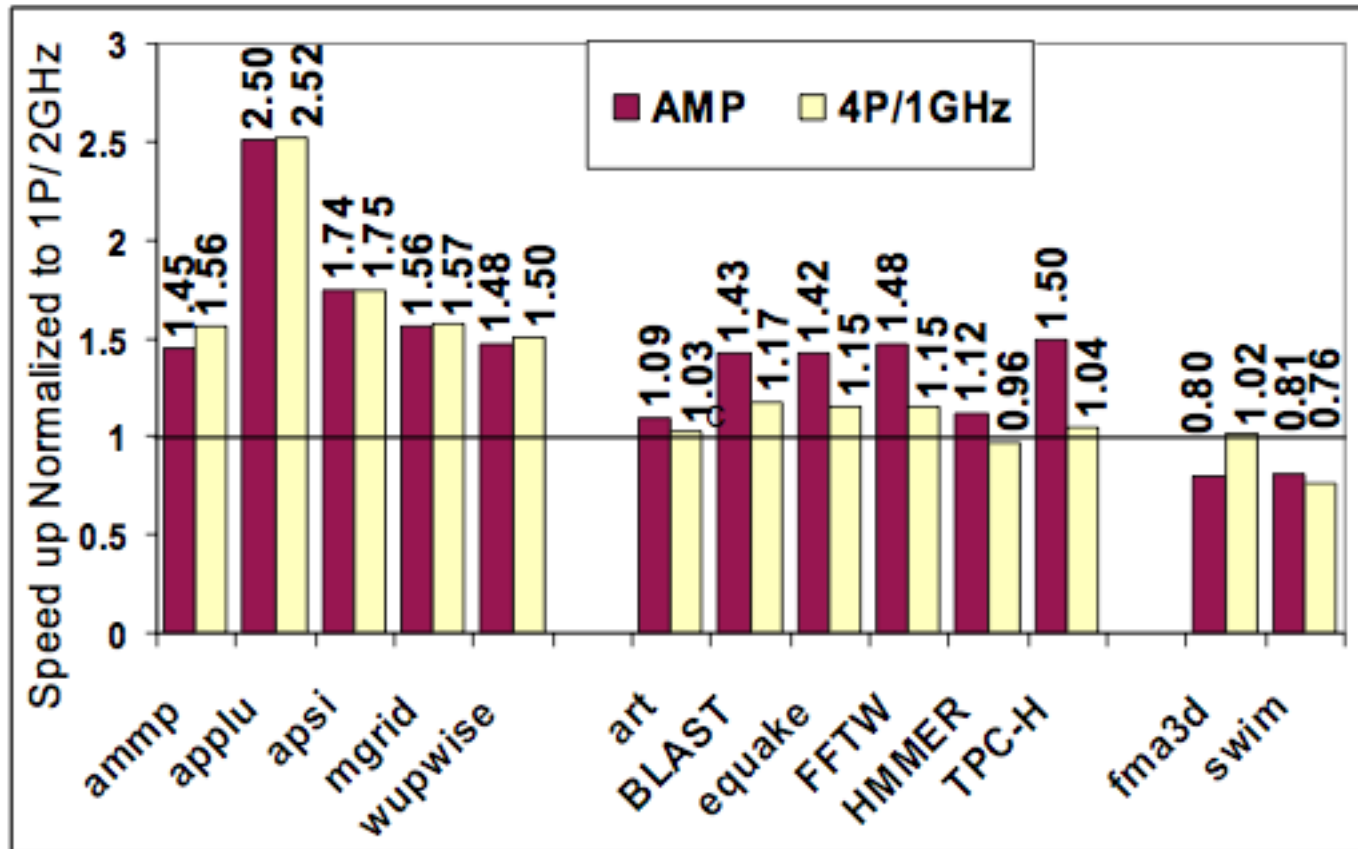
- ❑ Duty cycle changes during program run
- ❑ Parallel phases run on all or a subset of four processors
- ❑ Sequential phases of execution on 1P/2GHz
- ❑ Benchmarks with long sequential and parallel phases

EPI Throttling (Annavaram et al., ISCA' 05)

- Evaluation on Base SMP: 4 Base SMP: 4-way 2GHz Xeon, 2MB L3, 4GB Memory
- Hand-modified programs
 - OMP threads set to 3 for static AMP
 - Calls to set affinity in each thread for static AMP
 - Calls to change duty cycle and to set affinity in dynamic AMP

AMP Configuration	Programs
Static AMP: 1P/2GHz or 3P/1.25GHz	wupwise, swim, mgrid, equake, fma3d, art, ammp, BLAST, HMMER
Dynamic AMP: 1P/2GHz to 4P/1GHz	applu, apsi, FFTW, TPC-H

EPI Throttling (Annavaaram et al., ISCA' 05)



- Frequency boosting AMP improves performance compared to 4-way SMP for many applications

EPI Throttling

- Why does Frequency Boosting (FB) AMP not always improve performance?
- Loss of throughput in static AMP (only 3 processors in parallel portion)
 - Is this really the best way of using FB-AMP?
- Rapid transitions between serial and parallel phases
 - Data/thread migration and throttling overhead
- Boosting frequency does not help memory-bound phases

Use of ACMP to Improve Parallel Portion Performance

- Mutual Exclusion:
 - Threads are not allowed to update shared data concurrently
- Accesses to shared data are encapsulated inside ***critical sections***
- Only one thread can execute a critical section at a given time
- **Idea:** Ship critical sections to a large core
- Suleman et al., “**Accelerating Critical Section Execution with Asymmetric Multi-Core Architectures,**” ASPLOS 2009, IEEE Micro Top Picks 2010.