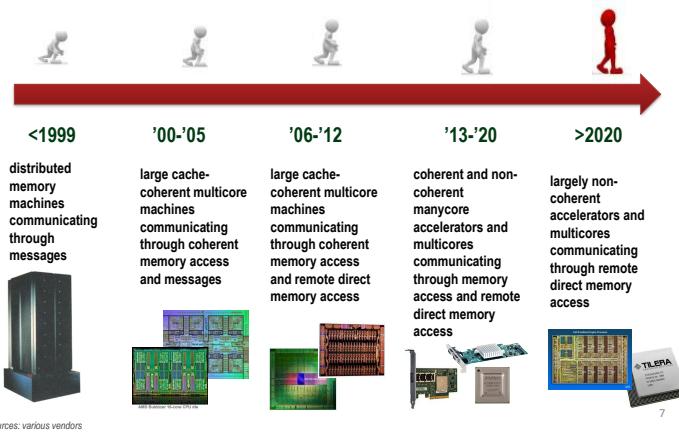




## Architecture Developments



## Computer Architecture vs. Physics

### Physics (technological constraints)

- Cost of data movement
- Capacity of DRAM cells
- Clock frequencies (constrained by end of Dennard scaling)
- Speed of Light
- Melting point of silicon

### Computer Architecture (design of the machine)

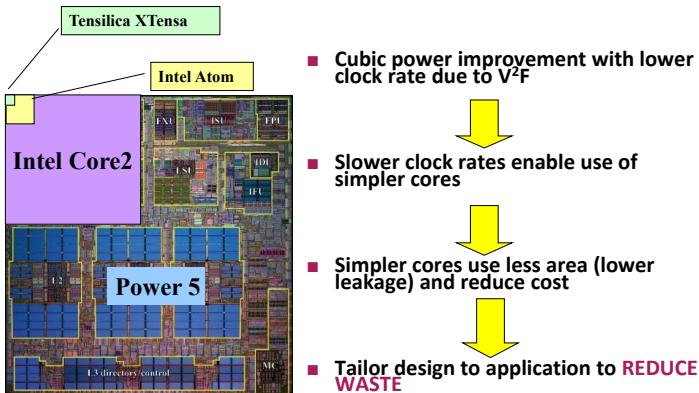
- Power management
- ISA / Multithreading
- SIMD widths

*"Computer architecture, like other architecture, is the art of determining the needs of the user of a structure and then designing to meet those needs as effectively as possible within economic and technological constraints." – Fred Brooks (IBM, 1962)*

*Have converted many former "power" problems into "cost" problems*

8

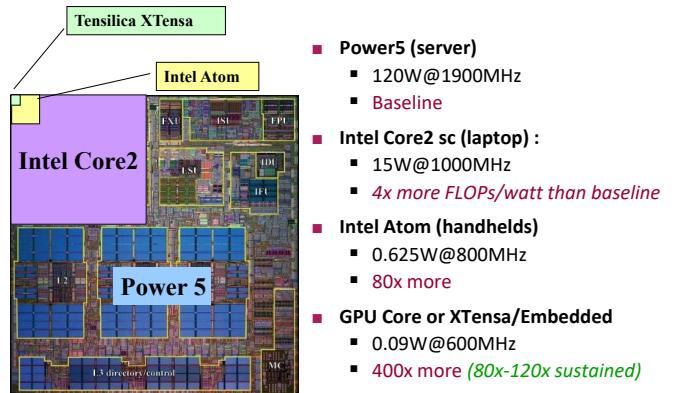
## Low-Power Design Principles (2005)



Credit: John Shalf (LBNL)

9

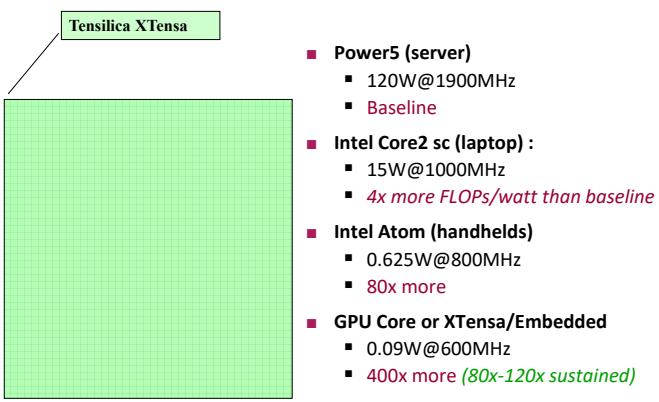
## Low-Power Design Principles (2005)



Credit: John Shalf (LBNL)

10

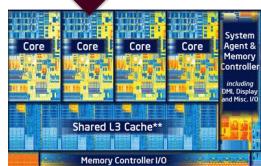
## Low-Power Design Principles (2005)



Even if each simple core is 1/4th as computationally efficient as complex core, you can fit hundreds of them on a single chip and still be 100x more power efficient.

## Heterogeneous Future (LOCs and TOCs)

### Big cores (very few)



Latency Optimized Core (LOC)  
Most energy efficient if you don't have lots of parallelism

Tiny core  
Lots of them!  
0.23mm  
0.2 mm

Throughput Optimized Core (TOC)  
Most energy efficient if you DO have a lot of parallelism!

Credit: John Shalf (LBNL)

11

12

## Data movement – the wires

### ■ Energy Efficiency of copper wire:

- Power = Frequency \* Length / cross-section-area



- Wire resistance increases with length

**Photonics could break through the bandwidth-distance limit**



### ■ Energy efficiency of transistors:

- Power = Frequency \* Current \* Resistance
- Capacitance  $\approx$  Area of Transistor
- Transistor efficiency improves as you shrink it

**Net result is that moving data on wires is starting to cost more energy than computing on said data (interest in Silicon Photonics)**

Credit: John Shaff (LBNL)

## Pin Limits

### ■ Moore's law doesn't apply to adding pins to package

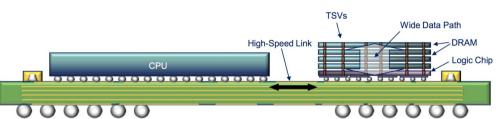
- 30%+ per year nominal Moore's Law
- Pins grow at ~1.5-3% per year at best

### ■ 4000 Pins is aggressive pin package

- Half of those would need to be for power and ground
- Of the remaining 2k pins, run as differential pairs
- Beyond 15Gbps per pin power/complexity costs hurt!
- 10Gbps \* 1k pins is ~1.2TBytes/sec

### ■ 2.5D Integration gets boost in pin density

- But it's a 1 time boost (how much headroom?)
- 4TB/sec? (maybe 8TB/s with single wire signaling?)

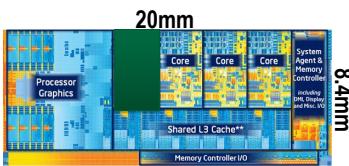


14

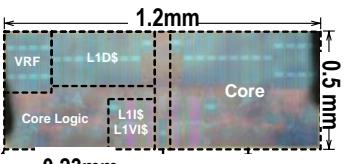
Credit: John Shaff (LBNL)

## Die Photos (3 classes of cores)

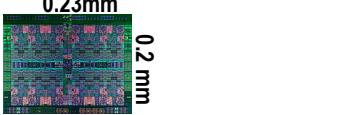
### Big intel



### Small RISC-V



### Tiny tensilica



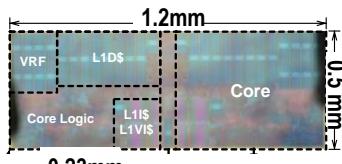
Credit: John Shaff (LBNL)

## Strip down to the core

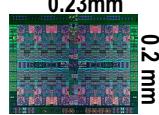
### intel



### RISC-V



### tensilica

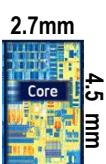


Credit: John Shaff (LBNL)

16

## Actual Size

### intel



### RISC-V



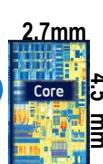
### tensilica



Credit: John Shaff (LBNL)

## Basic Stats

### intel



**Core Energy/Area est.**

Area: 12.25 mm<sup>2</sup>  
Power: 2.5W  
Clock: 2.4 GHz  
E/op: 651 pj

### RISC-V



Area: 0.6 mm<sup>2</sup>  
Power: 0.3W (<0.2W)  
Clock: 1.3 GHz  
E/op: 150 (75) pj

### tensilica



Area: 0.046 mm<sup>2</sup>  
Power: 0.025W  
Clock: 1.0 GHz  
E/op: 22 pj

**Wire Energy**  
Assumptions for 22nm  
100 fJ/bit per mm  
64bit operand  
Energy:  
1mm= $\sim$ 6pj  
20mm= $\sim$ 120pj

Credit: John Shaff (LBNL)

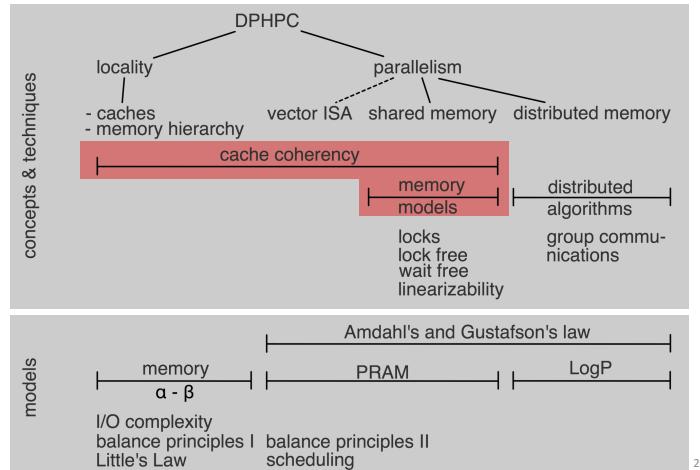
18

## When does data movement dominate?



Credit: John Shaff (LBNL)

## DPHPC Overview



19

20

## Goals of this lecture

- Memory Trends
- Cache Coherence
- Memory Consistency

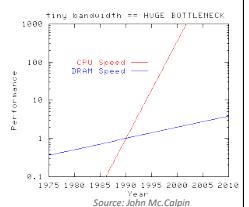
## Memory – CPU gap widens

- Measure processor speed as “throughput”
  - FLOPS/s, IOPS/s, ...
  - Moore’s law - ~60% growth per year



### Today's architectures

- POWER8: 338 dp GFLOP/s – 230 GB/s memory bw
- BW i7-5775C: 883 GFLOPS/s ~50 GB/s memory bw
- Trend: memory performance grows 10% per year



21

## Issues (AMD Interlagos as Example)

- How to measure bandwidth?
  - Data sheet (often peak performance, may include overheads)
    - Frequency times bus width: 51 GiB/s
  - Microbenchmark performance
    - Stride 1 access (32 MiB): 32 GiB/s
    - Random access (8 B out of 32 MiB): 241 MiB/s
    - Why?
  - Application performance
    - As observed (performance counters)
    - Somewhere in between stride 1 and random access
- How to measure Latency?
  - Data sheet (often optimistic, or not provided)
    - <100ns
  - Random pointer chase
    - 110 ns with one core, 258 ns with 32 cores!

## Conjecture: Buffering is a must!

- Two most common examples:
- Write Buffers
  - Delayed write back saves memory bandwidth
  - Data is often overwritten or re-read
- Caching
  - Directory of recently used locations
  - Stored as blocks (cache lines)

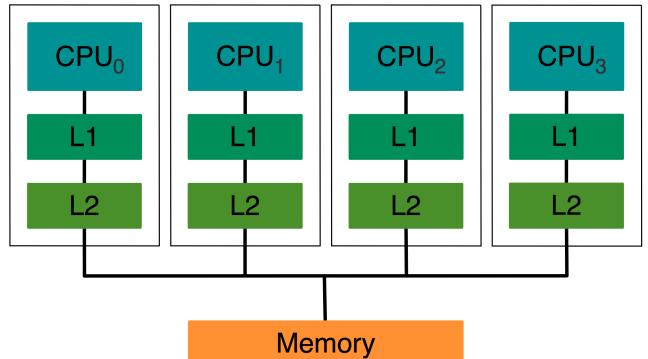
23

24

## Cache Coherence

- Different caches may have a copy of the same memory location!
- Cache coherence
  - Manages existence of multiple copies
- Cache architectures
  - Multi level caches
  - Shared vs. private (partitioned)
  - Inclusive vs. exclusive
  - Write back vs. write through
  - Victim cache to reduce conflict misses
  - ...

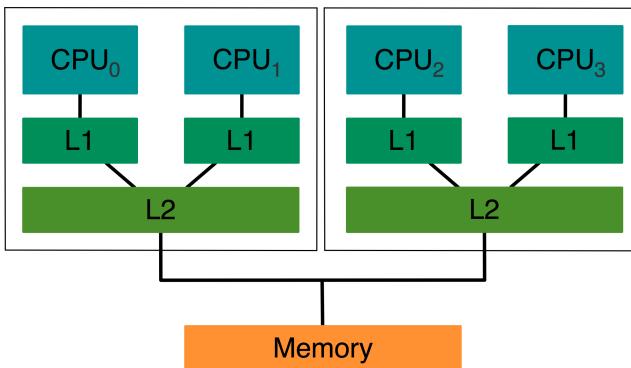
## Exclusive Hierarchical Caches



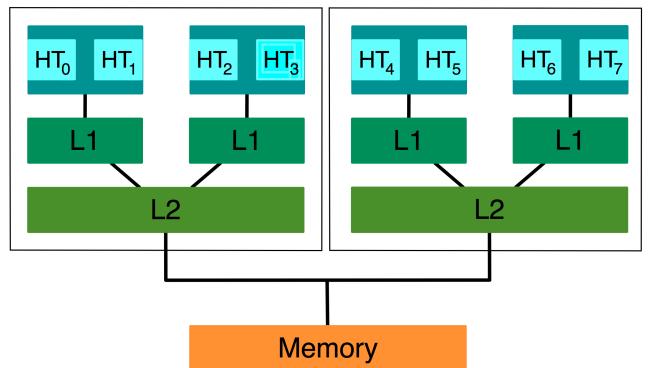
25

26

## Shared Hierarchical Caches



## Shared Hierarchical Caches with MT



27

28

## Caching Strategies (repeat)

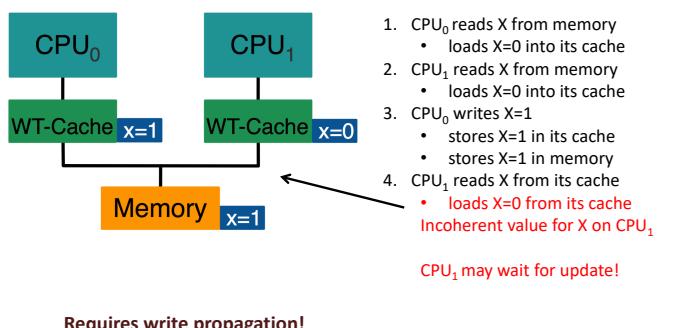
- Remember:
  - Write Back?
  - Write Through?
- Cache coherence requirements
 

A memory system is coherent if it guarantees the following:

  - Write propagation (updates are eventually visible to all readers)
  - Write serialization (writes to the same location must be observed in order)

*Everything else: memory model issues (later)*

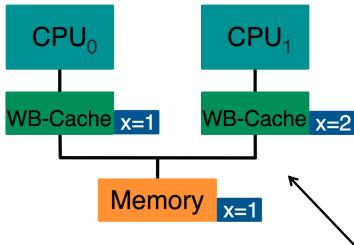
## Write Through Cache



29

30

## Write Back Cache



Requires write serialization!

1. CPU<sub>0</sub> reads X from memory
    - loads X=0 into its cache
  2. CPU<sub>1</sub> reads X from memory
    - loads X=0 into its cache
  3. CPU<sub>0</sub> writes X=1
    - stores X=1 in its cache
  4. CPU<sub>1</sub> writes X=2
    - stores X=2 in its cache
  5. CPU<sub>1</sub> writes back cache line
    - stores X=2 in memory
  6. CPU<sub>0</sub> writes back cache line
    - stores X=1 in memory
- Later store X=2 from CPU<sub>1</sub> lost*

31

## A simple (?) example

- Assume C99:

```
struct twoint {
    int a;
    int b;
}
```

- Two threads:

- Initially: a=b=0
- Thread 0: write 1 to a
- Thread 1: write 1 to b

- Assume non-coherent write back cache

- What may end up in main memory?

32

## Cache Coherence Protocol

- Programmer can hardly deal with unpredictable behavior!
- Cache controller maintains data integrity
  - All writes to different locations are visible

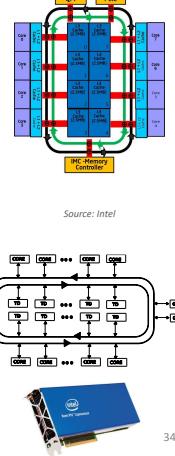
### Fundamental Mechanisms

- Snooping
  - Shared bus or (broadcast) network
- Directory-based
  - Record information necessary to maintain coherence:
    - E.g., owner and state of a line etc.

## Fundamental CC mechanisms

- Snooping

- Shared bus or (broadcast) network
- Cache controller “snoops” all transactions
- Monitors and changes the state of the cache’s data
- Works at small scale, challenging at large-scale
  - E.g., Intel Broadwell



34

33

## Cache Coherence Parameters

- Concerns/Goals
  - Performance
  - Implementation cost (chip space, more important: dynamic energy)
  - Correctness
  - (Memory model side effects)
- Issues
  - Detection (when does a controller need to act)
  - Enforcement (how does a controller guarantee coherence)
  - Precision of block sharing (per block, per sub-block?)
  - Block size (cache line size?)

## An Engineering Approach: Empirical start

- Problem 1: stale reads

- Cache 1 holds value that was already modified in cache 2
- Solution:
  - Disallow this state*
  - Invalidate all remote copies before allowing a write to complete*

- Problem 2: lost update

- Incorrect write back of modified line writes main memory in different order from the order of the write operations or overwrites neighboring data
- Solution:
  - Disallow more than one modified copy*

35

36

## Invalidation vs. update (I)

### ■ Invalidation-based:

- On each write of a shared line, it has to invalidate copies in remote caches
- Simple implementation for bus-based systems:
  - Each cache snoops*
  - Invalidate lines written by other CPUs*
  - Signal sharing for cache lines in local cache to other caches*

### ■ Update-based:

- Local write updates copies in remote caches
  - Can update all CPUs at once*
  - Multiple writes cause multiple updates (more traffic)*

## Invalidation vs. update (II)

### ■ Invalidation-based:

- Only write misses hit the bus (works with write-back caches)
- Subsequent writes to the same cache line are local
  - → Good for multiple writes to the same line (in the same cache)

### ■ Update-based:

- All sharers continue to hit cache line after one core writes
  - Implicit assumption: shared lines are accessed often*
- Supports producer-consumer pattern well
- Many (local) writes may waste bandwidth!

### ■ Hybrid forms are possible!

37

38

## MESI Cache Coherence

### ■ Most common hardware implementation of discussed requirements

aka. "Illinois protocol"

### Each line has one of the following states (in a cache):

#### ■ Modified (M)

- Local copy has been modified, no copies in other caches
- Memory is stale

#### ■ Exclusive (E)

- No copies in other caches
- Memory is up to date

#### ■ Shared (S)

- Unmodified copies *may* exist in other caches
- Memory is up to date

#### ■ Invalid (I)

- Line is not in cache

39

40

## Terminology

### ■ Clean line:

- Content of cache line and main memory is identical (also: memory is up to date)
- Can be evicted without write-back

### ■ Dirty line:

- Content of cache line and main memory differ (also: memory is stale)
- Needs to be written back eventually
  - Time depends on protocol details*

### ■ Bus transaction:

- A signal on the bus that can be observed by all caches
- Usually blocking

### ■ Local read/write:

- A load/store operation originating at a core connected to the cache

## Transitions in response to local reads

### ■ State is M

- No bus transaction

### ■ State is E

- No bus transaction

### ■ State is S

- No bus transaction

### ■ State is I

- Generate bus read request (BusRd)
  - May force other cache operations (see later)*
- Other cache(s) signal "sharing" if they hold a copy
- If shared was signaled, go to state S
- Otherwise, go to state E

### ■ After update: return read value

## Transitions in response to local writes

### ■ State is M

- No bus transaction

### ■ State is E

- No bus transaction
- Go to state M

### ■ State is S

- Line already local & clean
- There may be other copies
- Generate bus read request for upgrade to exclusive (BusRdX\*)
- Go to state M

### ■ State is I

- Generate bus read request for exclusive ownership (BusRdX)
- Go to state M

41

42

## Transitions in response to snooped BusRd

- State is M
  - Write cache line back to main memory
  - Signal "shared"
  - Go to state S (or E)
- State is E
  - Signal "shared"
  - Go to state S and signal "shared"
- State is S
  - Signal "shared"
- State is I
  - Ignore

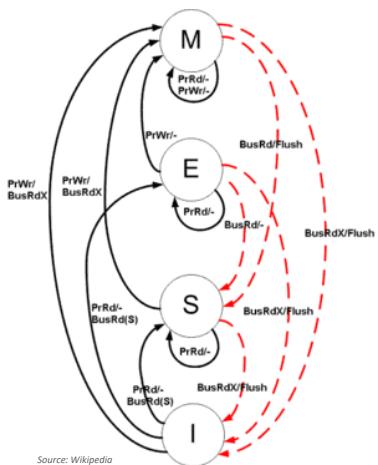
## Transitions in response to snooped BusRdX

- State is M
  - Write cache line back to memory
  - Discard line and go to I
- State is E
  - Discard line and go to I
- State is S
  - Discard line and go to I
- State is I
  - Ignore
- BusRdX\* is handled like BusRdX!

43

44

## MESI State Diagram (FSM)



45

46

## Small Exercise

- Initially: all in I state

Action	P1 state	P2 state	P3 state	Bus action	Data from
P1 reads x	E	I	I	BusRd	Memory
P2 reads x	S	S	I	BusRd	Cache
P1 writes x	M	I	I	BusRdX*	Cache
P1 reads x	M	I	I	-	Cache
P3 writes x	I	I	M	BusRdX	Memory

## Small Exercise

- Initially: all in I state

Action	P1 state	P2 state	P3 state	Bus action	Data from
P1 reads x	E	I	I	BusRd	Memory
P2 reads x	S	S	I	BusRd	Cache
P1 writes x	M	I	I	BusRdX*	Cache
P1 reads x	M	I	I	-	Cache
P3 writes x	I	I	M	BusRdX	Memory

47

48

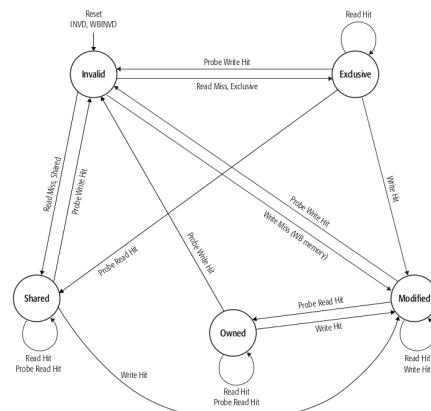
## Optimizations?

- Class question: what could be optimized in the MESI protocol to make a system faster?

## Related Protocols: MOESI (AMD)

- Extended MESI protocol
- Cache-to-cache transfer of modified cache lines
  - Cache in M or O state always transfers cache line to requesting cache
  - No need to contact (slow) main memory
- Avoids write back when another process accesses cache line
  - Good when cache-to-cache performance is higher than cache-to-memory  
*E.g., shared last level cache!*

## MOESI State Diagram



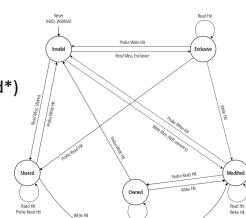
Source: AMD64 Architecture Programmer's Manual

49

50

## Related Protocols: MOESI (AMD)

- **Modified (M): Modified Exclusive**
  - No copies in other caches, local copy dirty
  - Memory is stale, cache supplies copy (reply to BusRd\*)
- **Owner (O): Modified Shared**
  - Exclusive right to make changes
  - Other S copies may exist ("dirty sharing")
  - Memory is stale, cache supplies copy (reply to BusRd\*)
- **Exclusive (E):**
  - Same as MESI (one local copy, up to date memory)
- **Shared (S):**
  - Unmodified copy may exist in other caches
  - Memory is up to date unless an O copy exists in another cache
- **Invalid (I):**
  - Same as MESI



51

52

## Multi-level caches

- **Most systems have multi-level caches**
  - Problem: only "last level cache" is connected to bus or network
  - Yet, snoop requests are relevant for inner-levels of cache (L1)
  - Modifications of L1 data may not be visible at L2 (and thus the bus)
- **L1/L2 modifications**
  - On BusRd check if line is in M state in L1  
*It may be in E or S in L2!*
  - On BusRdX(\*) send invalidations to L1
  - Everything else can be handled in L2
- **If L1 is write through, L2 could "remember" state of L1 cache line**
  - May increase traffic though

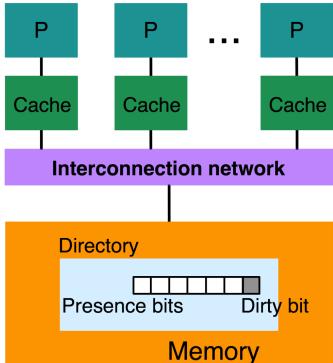
## Directory-based cache coherence

- **Snooping does not scale**
  - Bus transactions must be *globally* visible
  - Implies broadcast
- **Typical solution: tree-based (hierarchical) snooping**
  - Root becomes a bottleneck
- **Directory-based schemes are more scalable**
  - Directory (entry for each CL) keeps track of all owning caches
  - Point-to-point update to involved processors  
*No broadcast*
  - Can use specialized (high-bandwidth) network, e.g., HT, QPI ...

53

54

## Basic Scheme



55

- System with  $N$  processors  $P_i$
- For each memory block (size: cache line) maintain a directory entry
  - $N$  presence bits
  - Set if block in cache of  $P_i$
  - 1 dirty bit
- First proposed by Censier and Feautrier (1978)

## Directory-based CC: Read miss

- $P_i$  intends to read, misses
- If dirty bit (in directory) is off
  - Read from main memory
  - Set presence[i]
  - Supply data to reader
- If dirty bit is on
  - Recall cache line from  $P_j$  (determine by presence[j])
  - Update memory
  - Unset dirty bit, block shared
  - Set presence[i]
  - Supply data to reader

56

## Directory-based CC: Write miss

- $P_i$  intends to write, misses
- If dirty bit (in directory) is off
  - Send invalidations to all processors  $P_j$  with presence[j] turned on
  - Unset presence bit for all processors
  - Set dirty bit
  - Set presence[i], owner  $P_i$
- If dirty bit is on
  - Recall cache line from owner  $P_j$
  - Update memory
  - Unset presence[j]
  - Set presence[i], dirty bit remains set
  - Supply data to writer

57

## Discussion

- Scaling of memory bandwidth
  - No centralized memory
- Directory-based approaches scale with restrictions
  - Require presence bit for each cache
  - Number of bits determined at design time
  - Directory requires memory (size scales linearly)
  - Shared vs. distributed directory
- Software-emulation
  - Distributed shared memory (DSM)
  - Emulate cache coherence in software (e.g., TreadMarks)
  - Often on a per-page basis, utilizes memory virtualization and paging

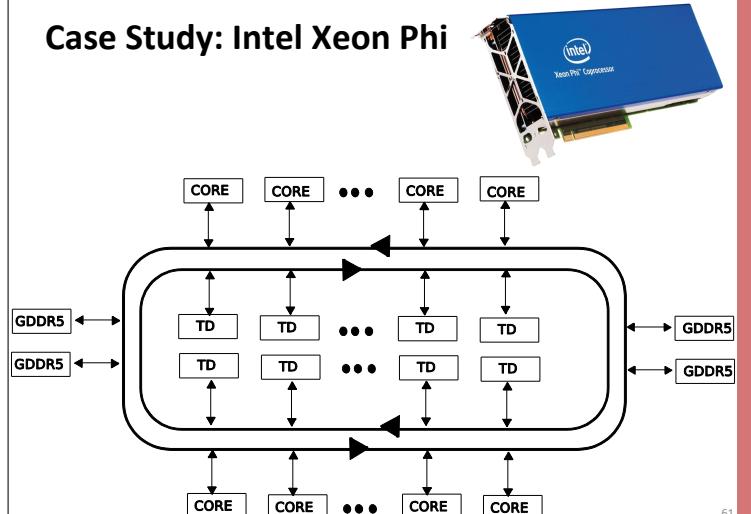
59

## Open Problems (for projects or theses)

- Tune algorithms to cache-coherence schemes
  - What is the optimal parallel algorithm for a given scheme?
  - Parameterize for an architecture
- Measure and classify hardware
  - Read Maranget et al. "A Tutorial Introduction to the ARM and POWER Relaxed Memory Models" and have fun!
  - RDMA consistency is barely understood!
  - GPU memories are not well understood!
    - Huge potential for new insights!
- Can we program (easily) without cache coherence?
  - How to fix the problems with inconsistent values?
  - Compiler support (issues with arrays)?

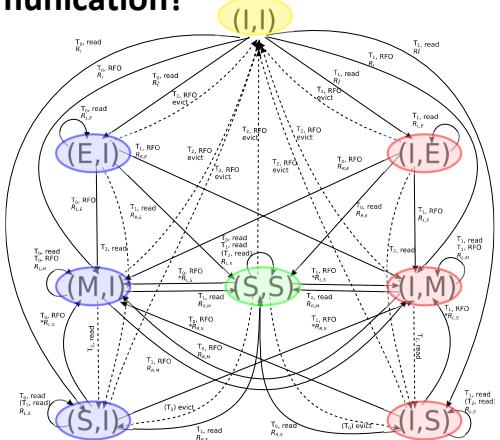
60

## Case Study: Intel Xeon Phi

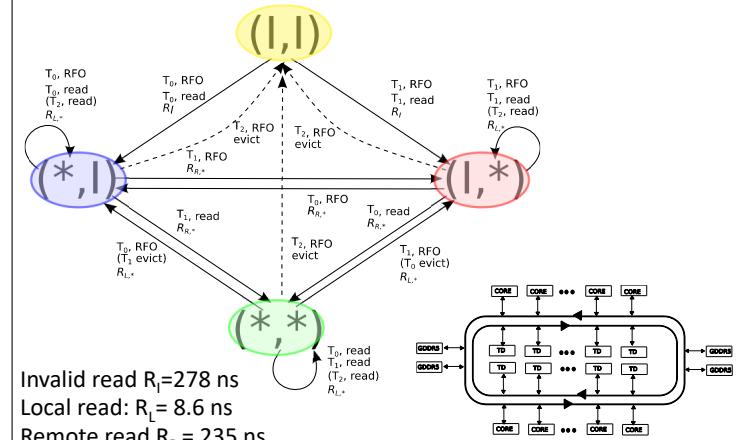


61

## Communication?

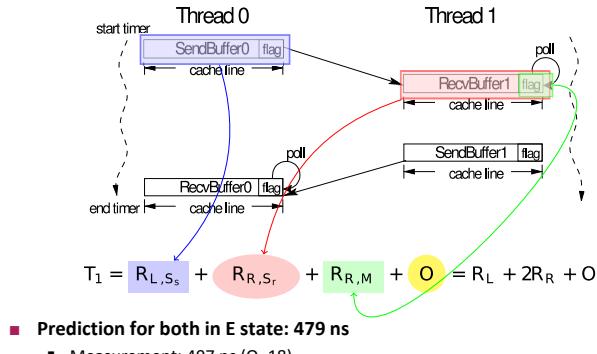


Ramos, Hoeffer: "Modeling Communication in Cache-Coherent SMP Systems - A Case-Study with Xeon Phi", HPDC'13



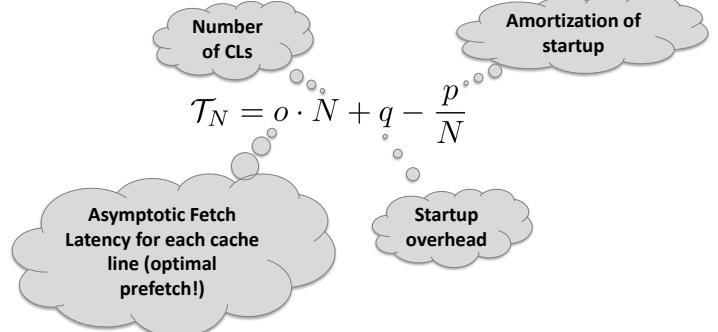
Inspired by Molka et al.: "Memory performance and cache coherency effects on an Intel Nehalem multiprocessor system"

## Single-Line Ping Pong



## Multi-Line Ping Pong

- More complex due to prefetch



## Multi-Line Ping Pong

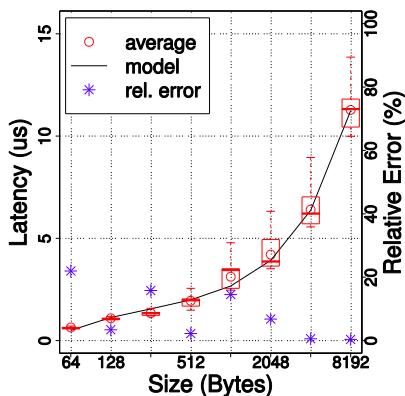
$$T_N = o \cdot N + q - \frac{p}{N}$$

### E state:

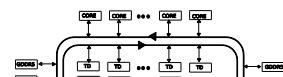
- $o=76$  ns
- $q=1,521$  ns
- $p=1,096$  ns

### I state:

- $o=95$  ns
- $q=2,750$  ns
- $p=2,017$  ns



## DTD Contention ☺



$$T_C(n_{th}) = c \cdot n_{th} + b - \frac{a}{n_{th}}$$

### E state:

- $a=0$ ns
- $b=320$ ns
- $c=56.2$ ns

