# Multicore Programming

## Parallel Hardware and Performance

8 Nov 2010 (Part 1)

Peter Sewell      Jaroslav Ševčík      Tim Harris

# Merge sort

16MB input (32-bit integers)

Recurse(left)

Recurse(right)

~98% execution time

Merge to scratch array

Copy back to input array

~2% execution time

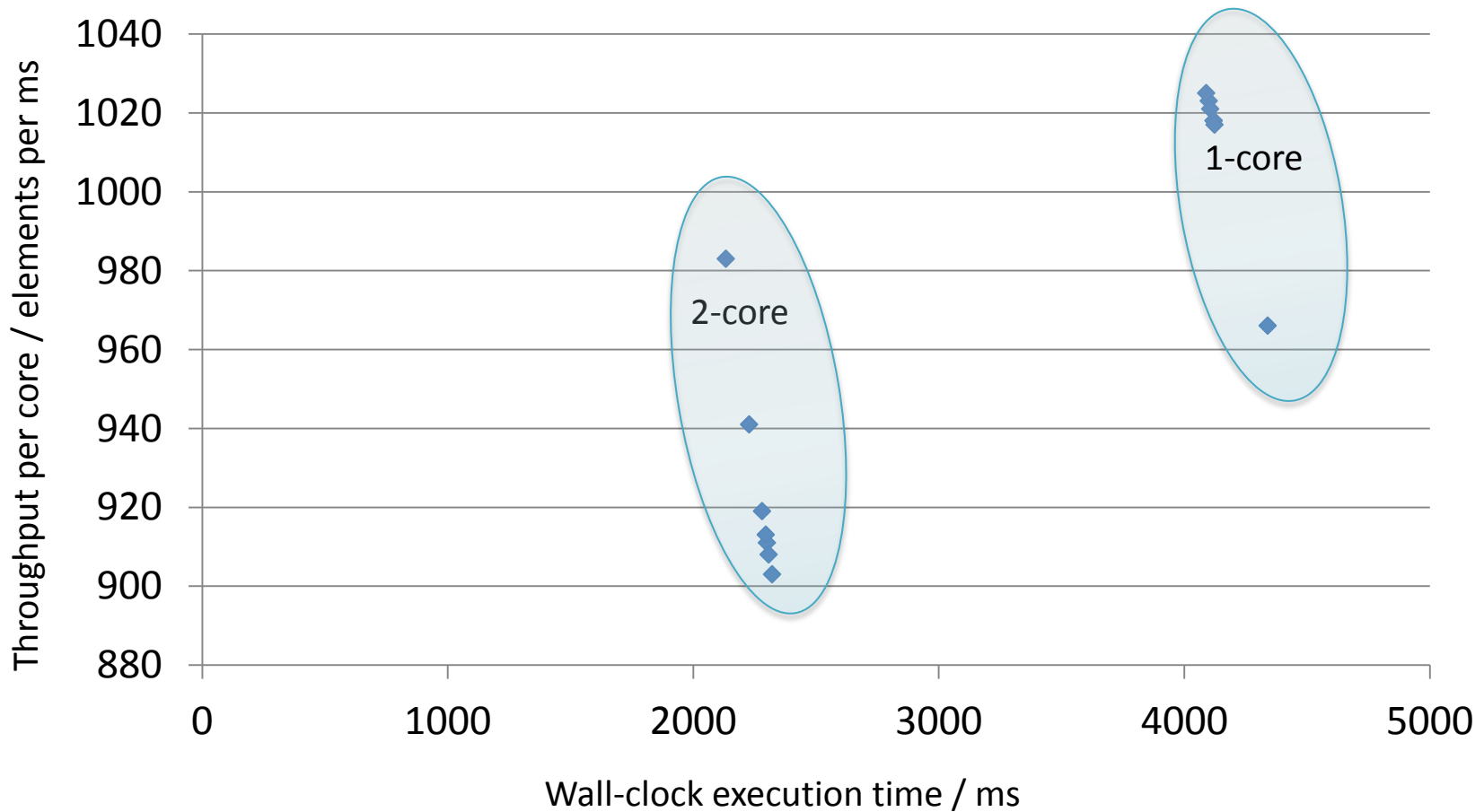# Merge sort, dual core

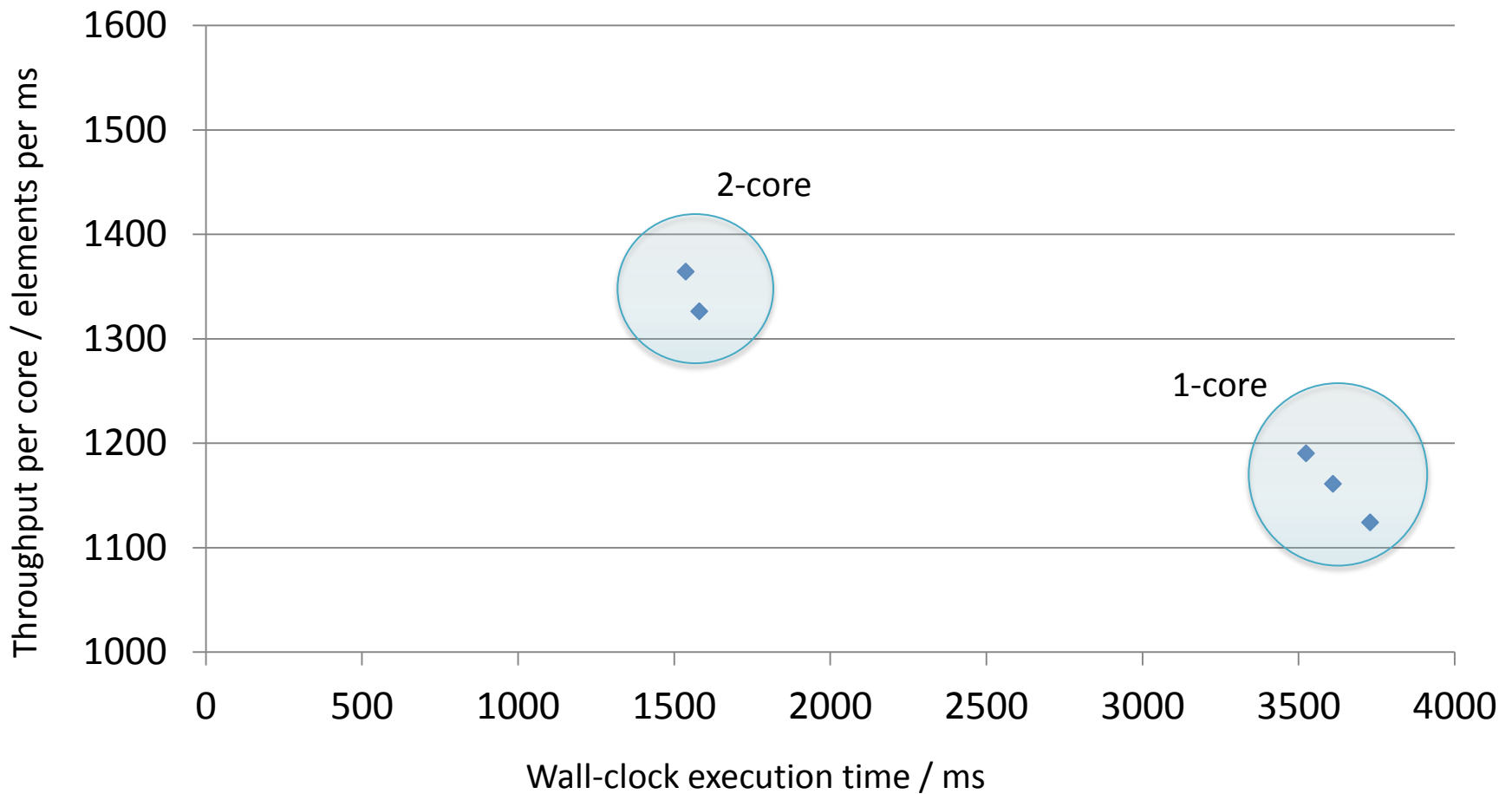16MB input (32-bit integers)

Recurse(left)    Recurse(right)

Merge to scratch array

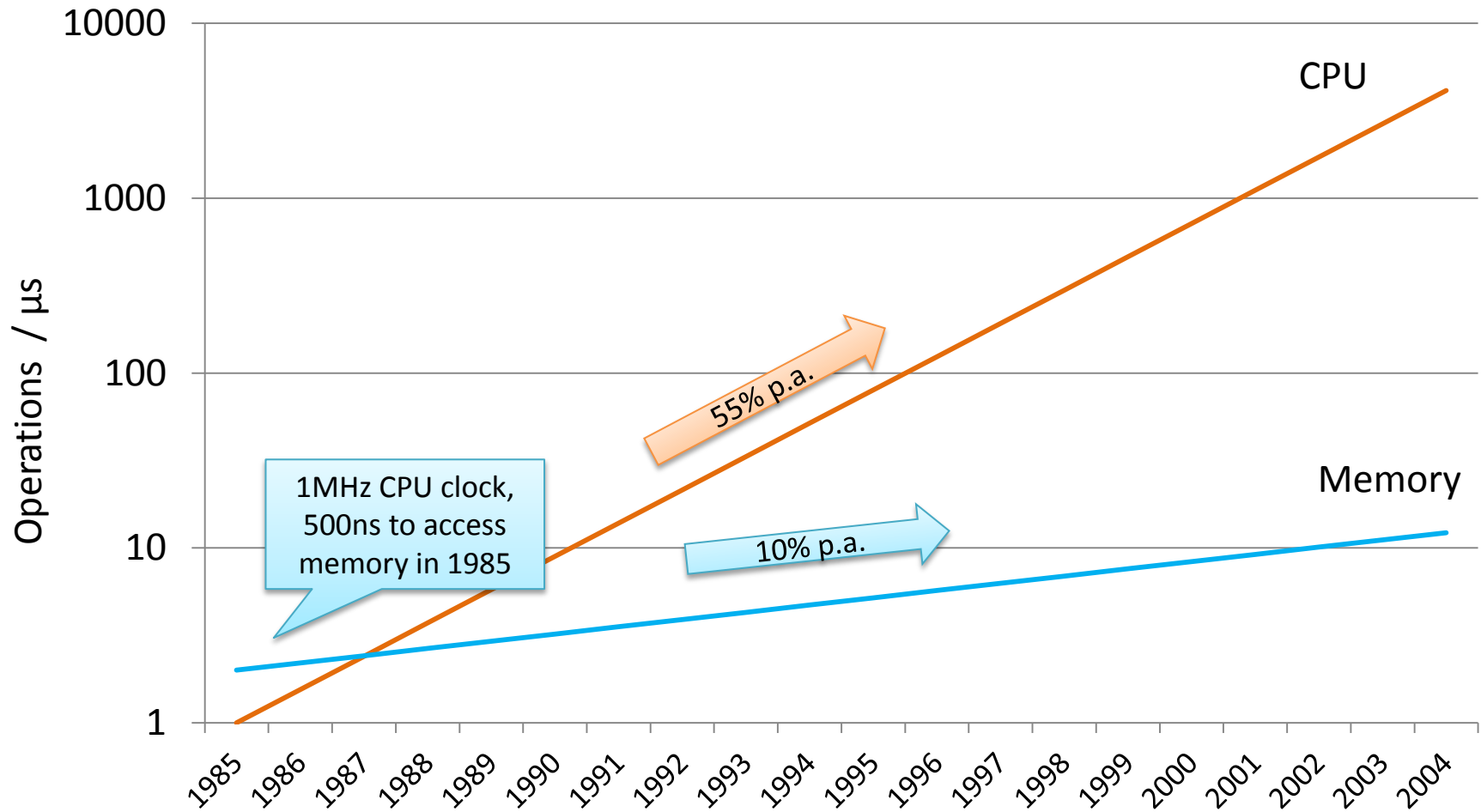Copy back to input array

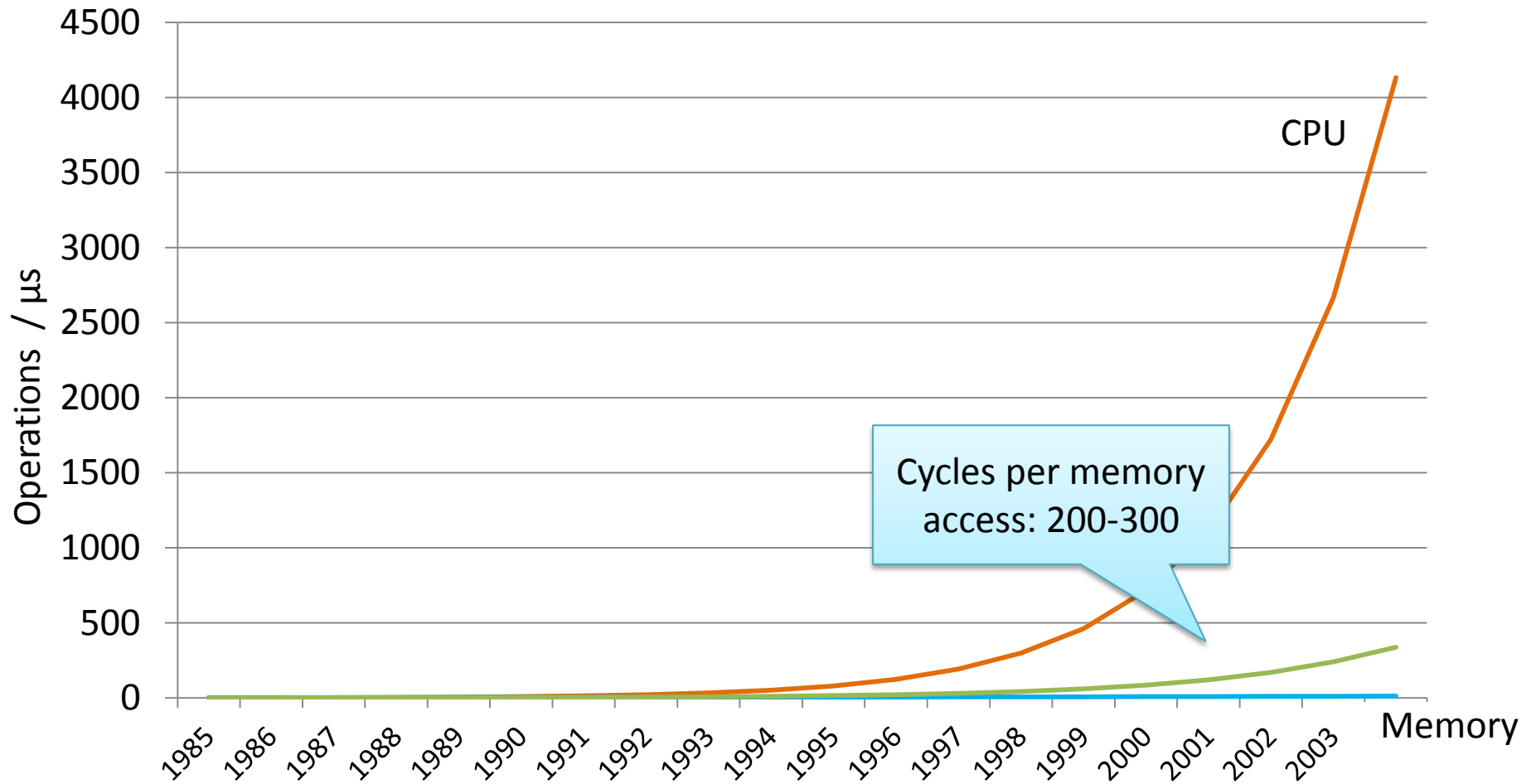# T7300 dual-core laptop

# AMD Phenom 3-core

# The power wall

- Moore's law: area of transistor is about 50% smaller each generation
  - A given chip area accommodates twice as many transistors

- $P_{dyn} = \alpha\, f\, V^2$
  - Shrinking process technology (constant field scaling) allows reducing V to partially counteract increasing f
  - V cannot be reduced arbitrarily
    - Halving V more than halves max f (for a given transistor)
    - Physical limits

- $P_{leak} = V(I_{sub} + I_{ox})$
  - Reduce $I_{sub}$ (sub-threshold leakage): turn off component, increase threshold volatage (reduces max f)
  - Reduce $I_{ox}$ (gate-oxide leakage): increase oxide thickness (but it needs to decrease with process scale)

# The memory wall



Y-axis: Operations / μs — 1, 10, 100, 1000, 10000

X-axis: 1985, 1986, 1987, 1988, 1989, 1990, 1991, 1992, 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004

CPU

Memory

55% p.a.

10% p.a.

1MHz CPU clock, 500ns to access memory in 1985

# The ILP wall

- ILP = "Instruction level parallelism"
- Implicit parallelism between instructions in a single thread
- Identified by the hardware
  - Speculate past memory accesses
  - Speculate past control transfer
- Diminishing returns

# Power wall + ILP wall + memory wall = brick wall

- Power wall means we can't just clock processors faster any longer

- Memory wall means that many workload's perf is dominated by memory access times

- ILP wall means we can't find extra work to keep functional units busy while waiting for memory accesses

# Amdahl's law

- "Sorting takes 70% of the execution time of a sequential program. You replace the sorting algorithm with one that scales perfectly on multi-core hardware. On a machine with $n$ cores, how many cores do you need to use to get a 4x speed-up on the overall algorithm?"

# Amdahl's law, f=70%



Speedup (y-axis) vs #cores (x-axis). Desired 4x speedup shown as a horizontal line at 4.0. Speedup achieved (perfect scaling on 70%) curve rising from 1.0 toward ~2.9.

# Amdahl's law, f=70%

$$MaxSpeedup(f, c) = \frac{1}{(1 - f) + {f}/{c}}$$

f = fraction of code speedup applies to
c = number of cores used

# Amdahl's law, f=70%



Desired 4x speedup

Speedup achieved (perfect scaling on 70%)

Limit as $c \to \infty = 1/(1-f) = 3.33$

Speedup

#cores

# Amdahl's law, f=10%

# Amdahl's law, f=98%

# Amdahl's law & multi-core

Suppose that the same h/w budget (space or power) can make us:

| | |
|---|---|
| 1 | 2 |

1

| | |
|---|---|
| 3 | 4 |

# Perf of big & small cores

Assumption: perf = α √resource

Total perf:
16 * 1/4 = 4

Total perf:
1 * 1 = 1

Core perf (relative to 1 big core)

Resources dedicated to core

1/16  1/8  1/4  1/2  1

# Amdahl's law, f=75%

Why parallelism?

Amdahl's law

Why asymmetric performance?

Parallel algorithms

Multi-processing hardware

# Asymmetric chips
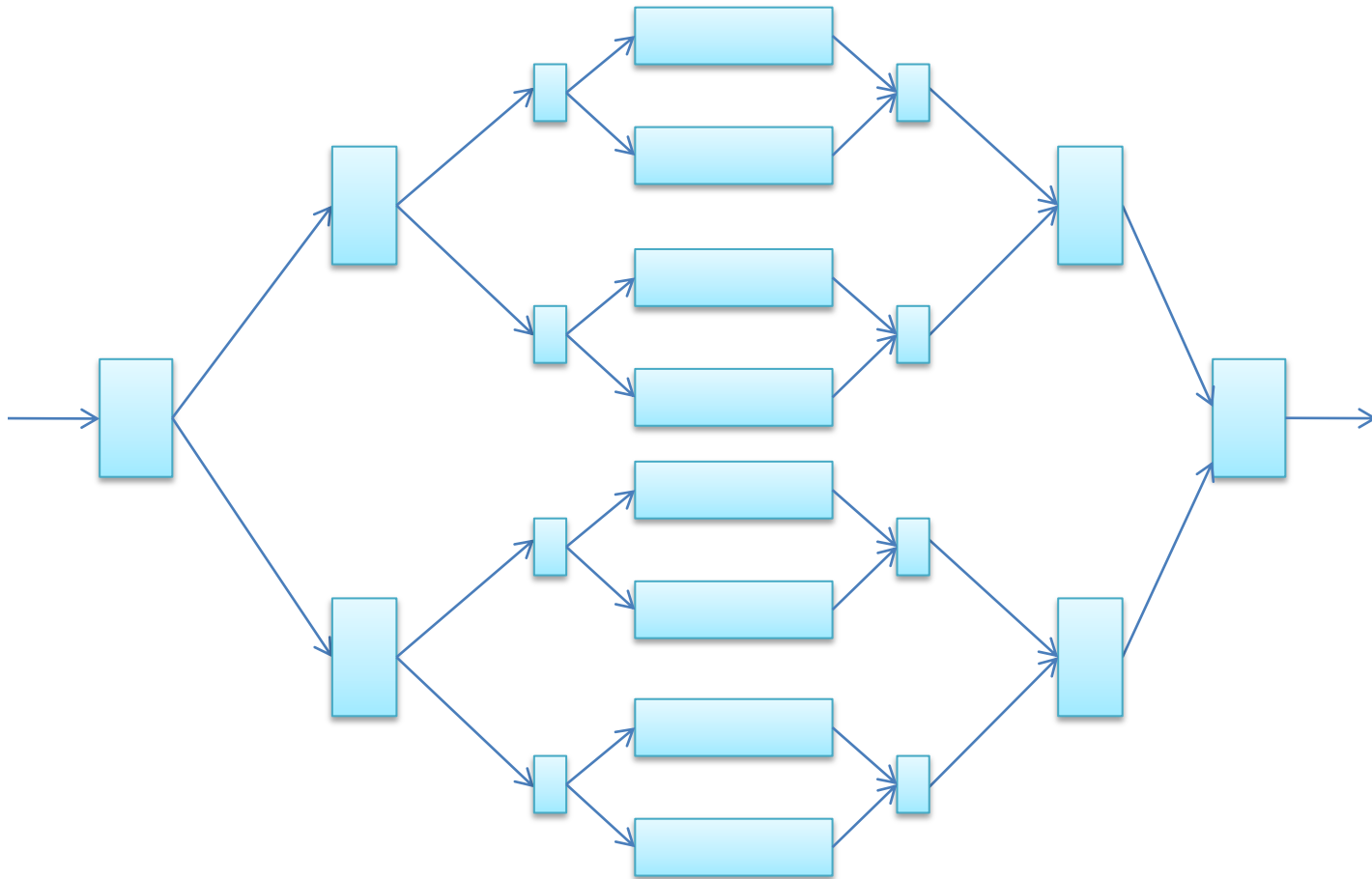
# Amdahl's law, f=75%

# Amdahl's law, f=98%

# Amdahl's law, f=98%

# Amdahl's law, f=98%



Speedup (relative to 1 big core)

1+192

256 small

Leave larger core idle in parallel section

#Cores

Why parallelism?

Amdahl's law

Why asymmetric performance?

**Parallel algorithms**

Multi-processing hardware

# Merge sort (2-core)

# Merge sort (4-core)

# Merge sort (8-core)

# T$_\infty$ (span): critical path length

# $T_1$ (work): time to run sequentially

# $T_{serial}$: optimized sequential code

# A good multi-core parallel algorithm

- $T_1/T_{serial}$ is low
  - What we lose on sequential performance we must make up through parallelism
  - Resource availability may limit the ability to do that

- $T_\infty$ grows slowly with the problem size
  - We tackle bigger problems by using more cores, not by running for longer

# Quick-sort on "good" input



Recurse

$T_1 = O(n \log n)$
$T_\infty = O(n)$

Sequential partition

Sequential append

Recurse

# Quick-sort on "good" input

# Quick-sort on "good" input



Parallel partition → Recurse, Recurse → Parallel append

$$T_1 = O(n \log n)$$
$$T_\infty = O(\log^2 n)$$

# Scheduling from a DAG

# Scheduling from a DAG



Speedup ($T_p / T_1$) vs. # cores

$T_p \geq T_\infty$

$T_p \geq T_1 / P$

# In CILK: $T_p \approx T_1 / p + c\, T_\infty$



Given abundant parallelism ($T_\infty \ll T_1$):
- "Work first principle"
- Keep $T_1$ close to $T_{serial}$
- Put overhead on $T_\infty$ to keep $T_1$ low

$T_p \geq T_\infty$

$T_p \geq T_1 / P$

$c=1$

$c=1.25$

Speedup ($T_p / T_1$)

# cores

Why parallelism?

Amdahl's law

Why asymmetric performance?

Parallel algorithms

Multi-processing hardware

# Multi-threaded h/w



- Multiple threads in a workload with:
  - Poor spatial locality
  - Frequent memory accesses

ALU
ALU
Multi-threaded single core

L1 cache (64KB)

L2 cache (4MB)

Main memory

1
2
3
4
5
…

1
2
3
4
5
…

1
2
3
4
5
…

# Multi-threaded h/w



- Multiple threads with synergistic resource needs
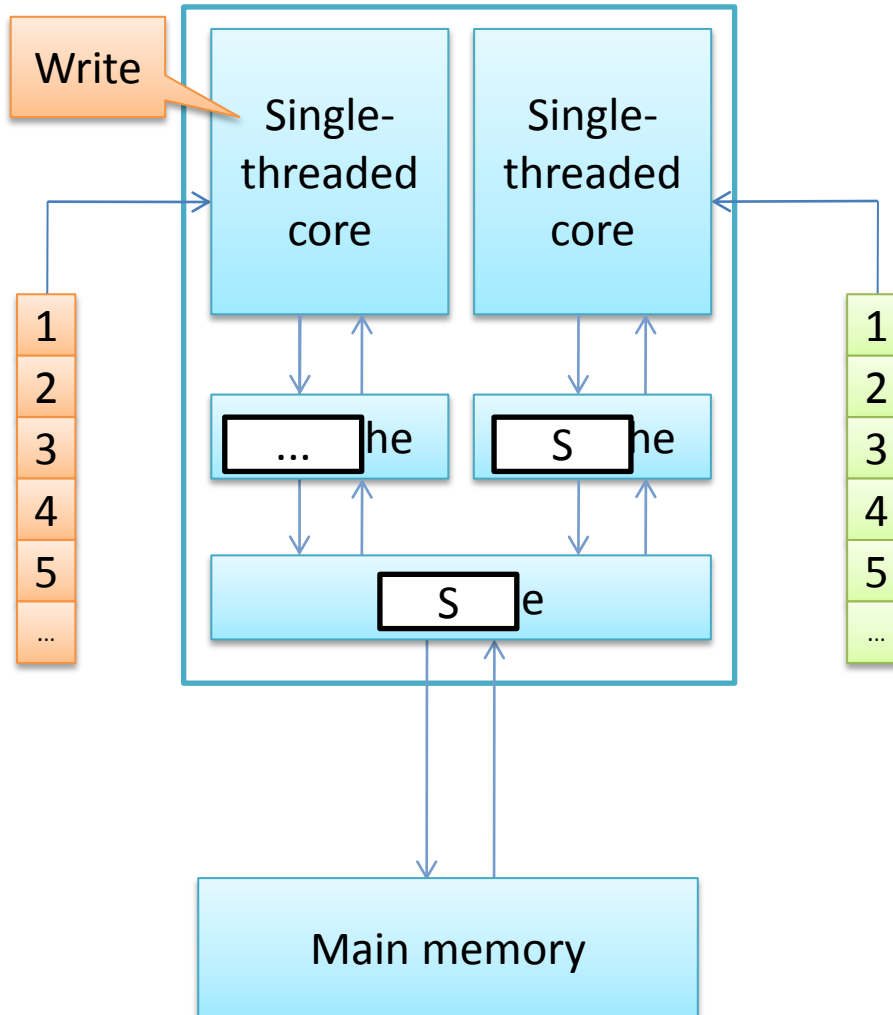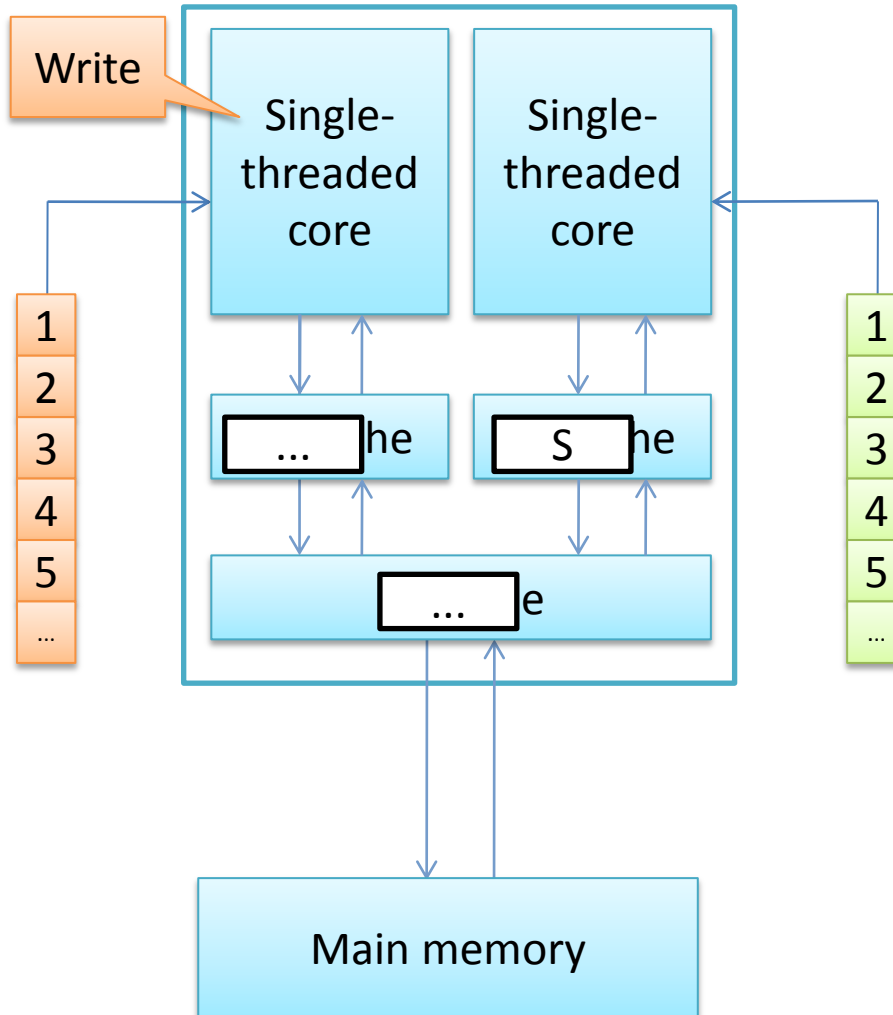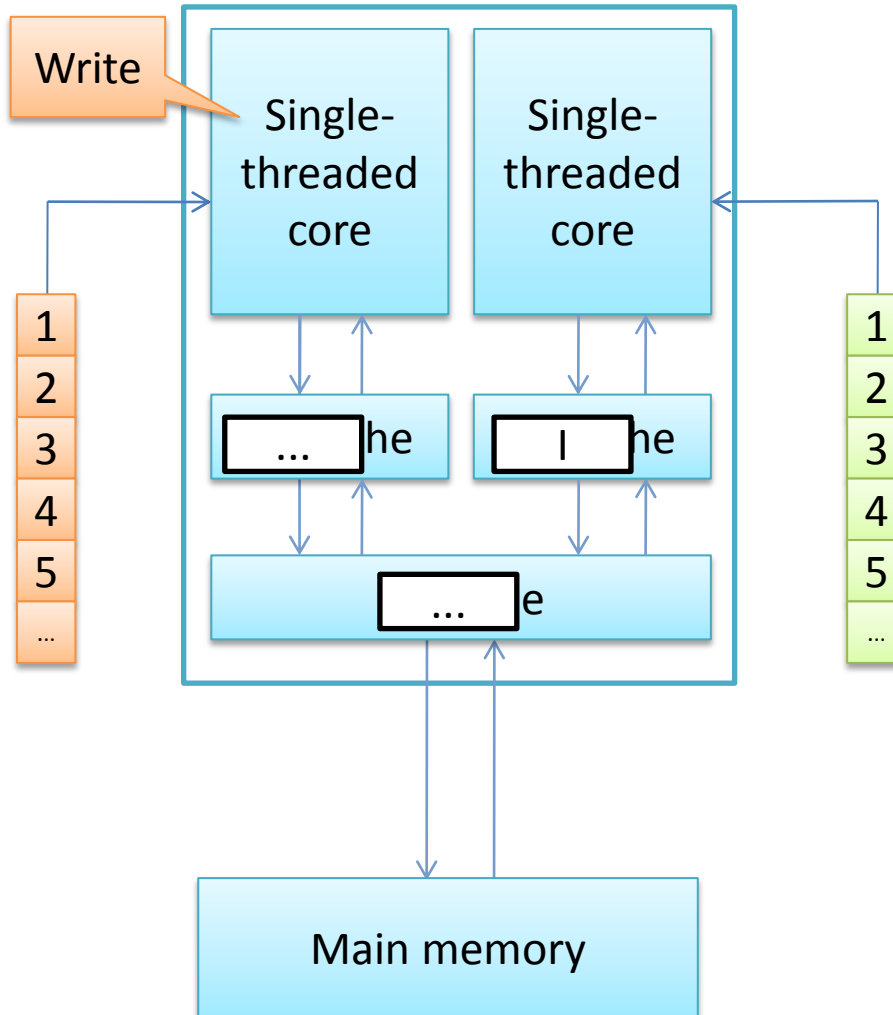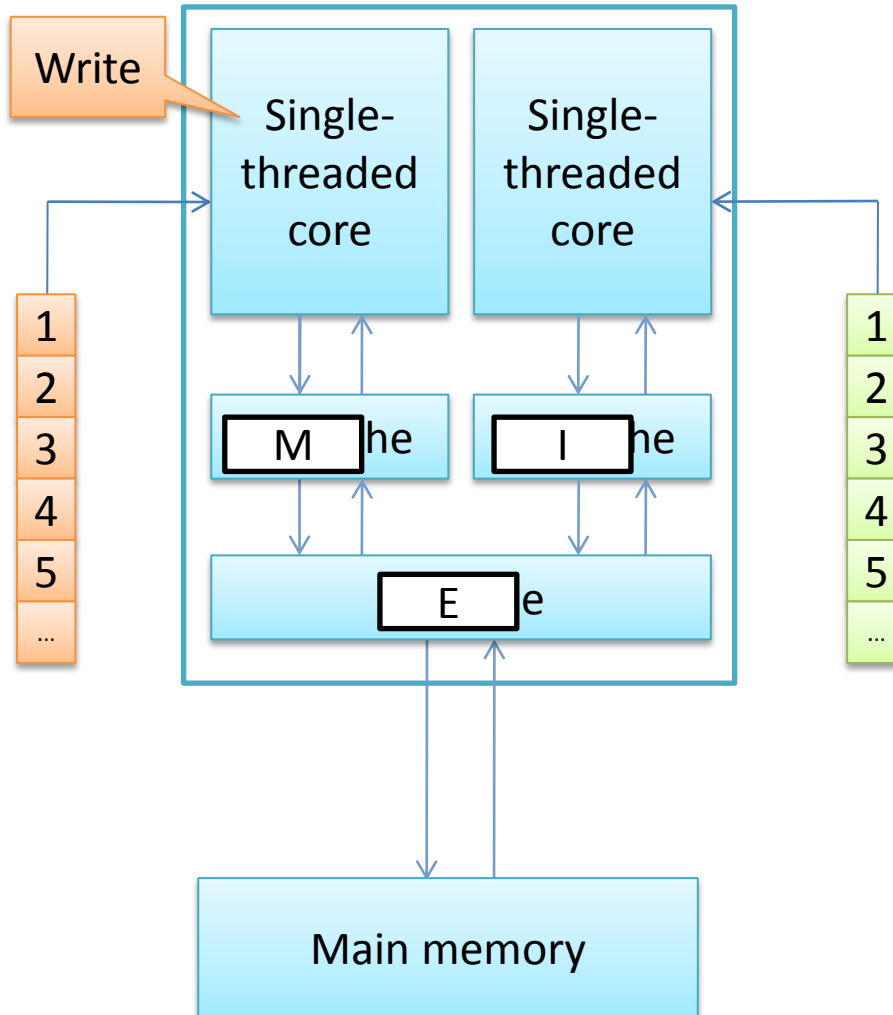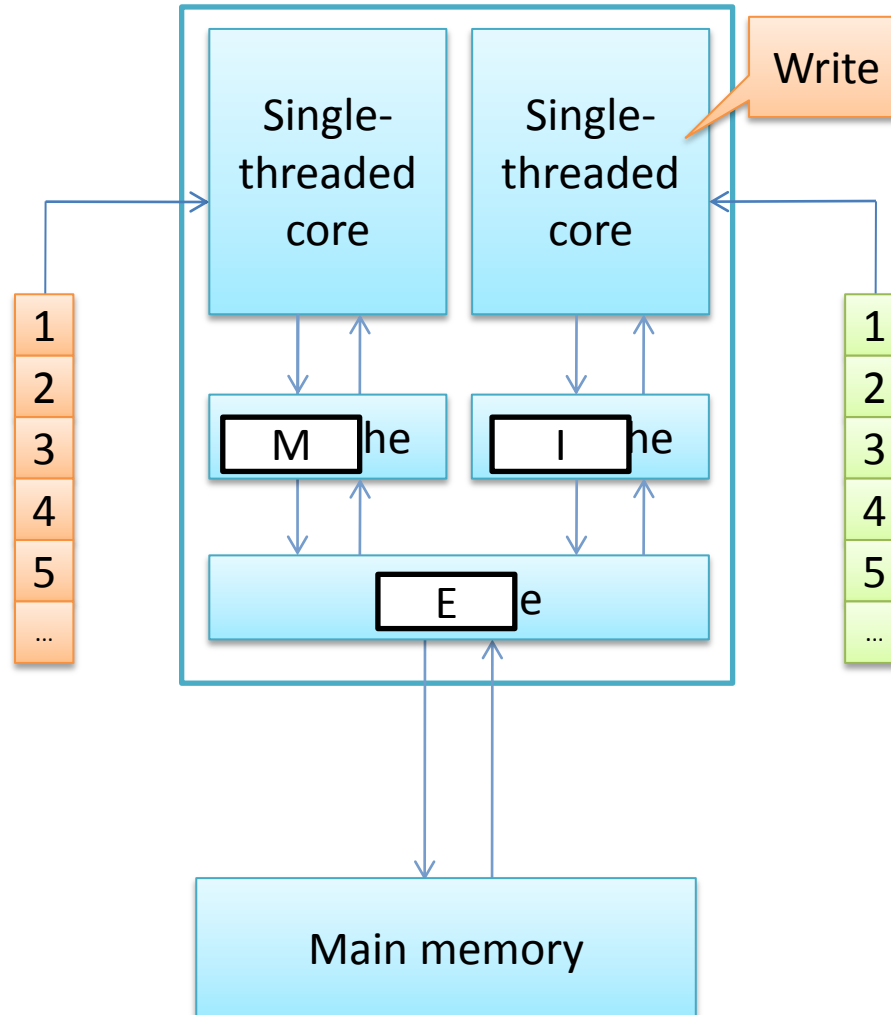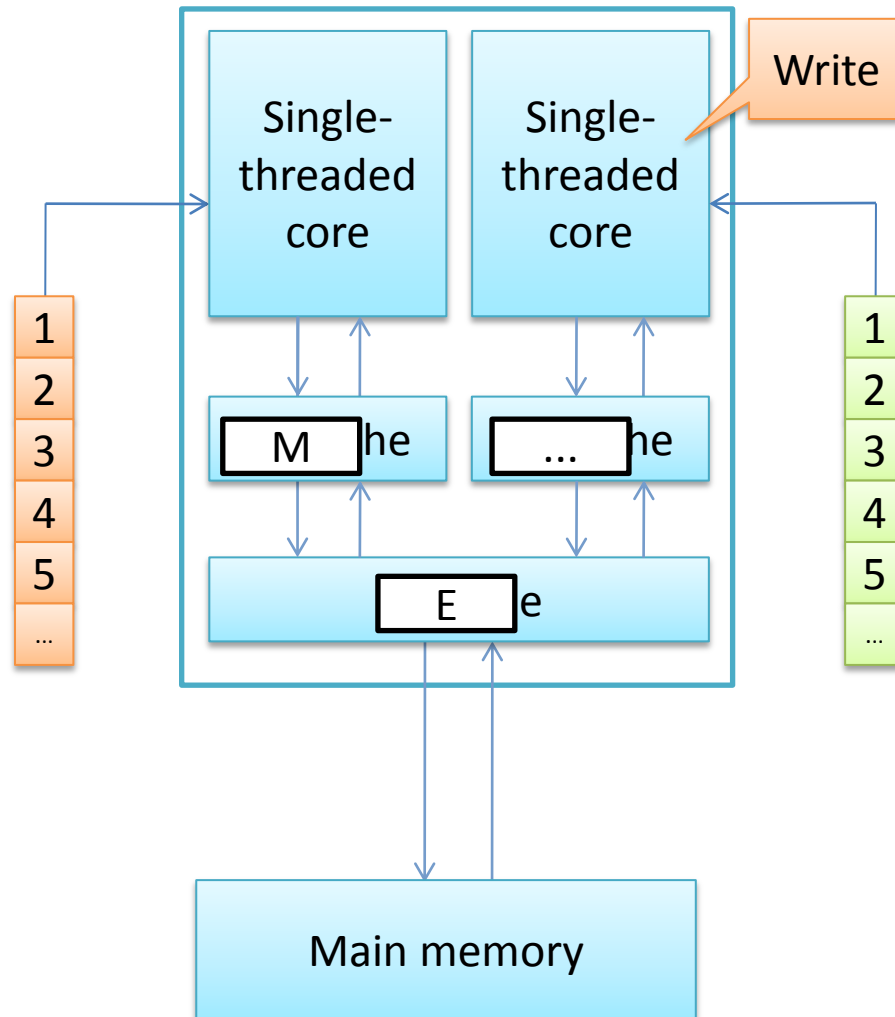
# Multi-core h/w – common L2

# Multi-core h/w – common L2

# Multi-core h/w – common L2

# Multi-core h/w – common L2

# Multi-core h/w – common L2

# Multi-core h/w – common L2
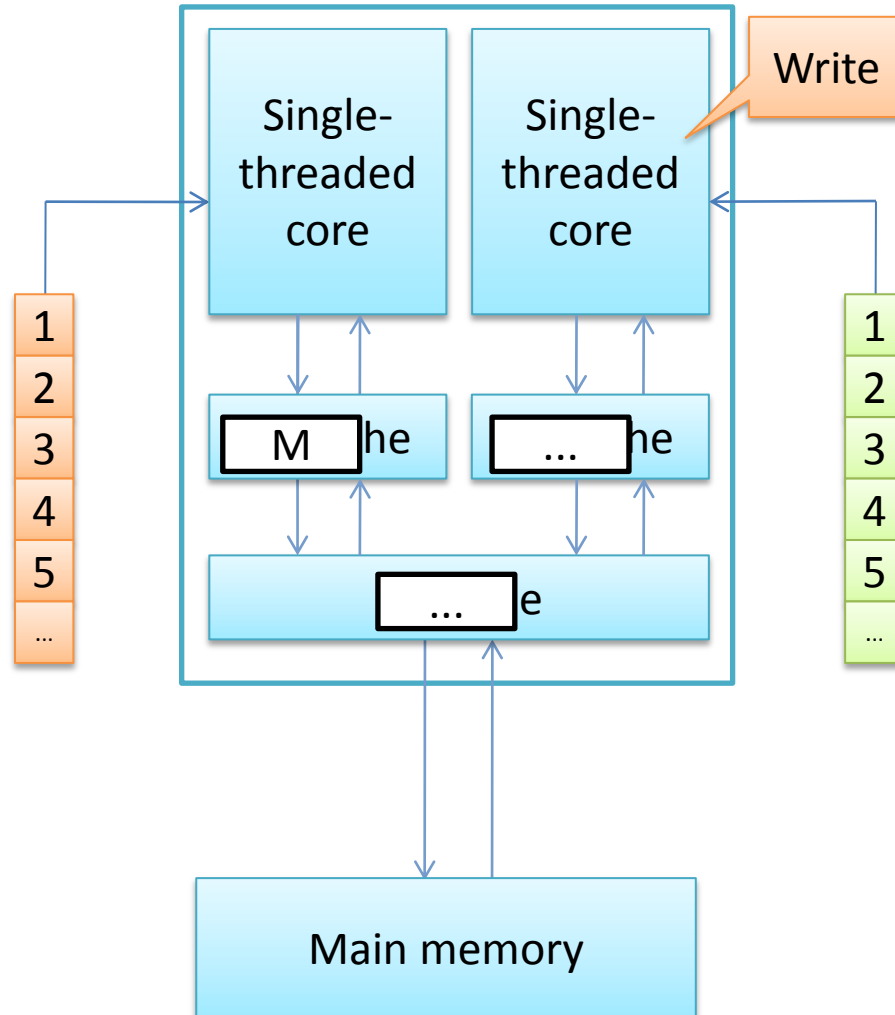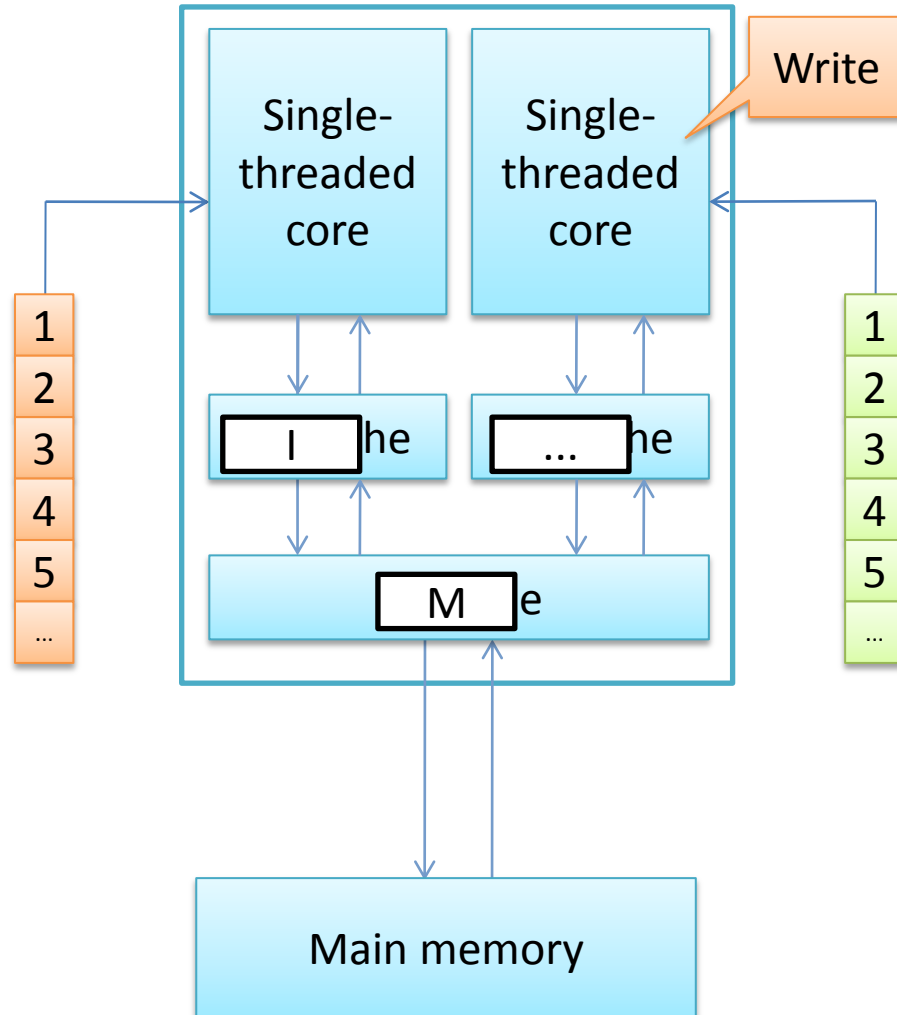
# Multi-core h/w – common L2

# Multi-core h/w – common L2

# Multi-core h/w – common L2

# Multi-core h/w – common L2

# Multi-core h/w – common L2

# Multi-core h/w – common L2

# Multi-core h/w – common L2
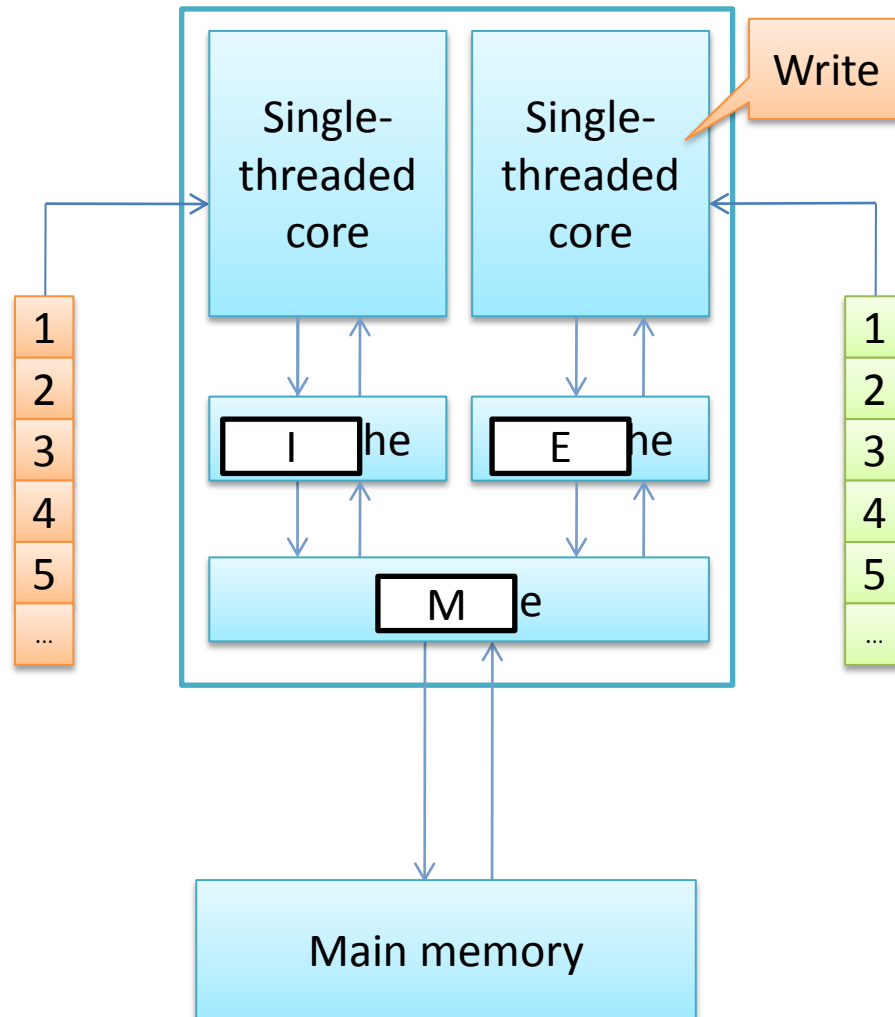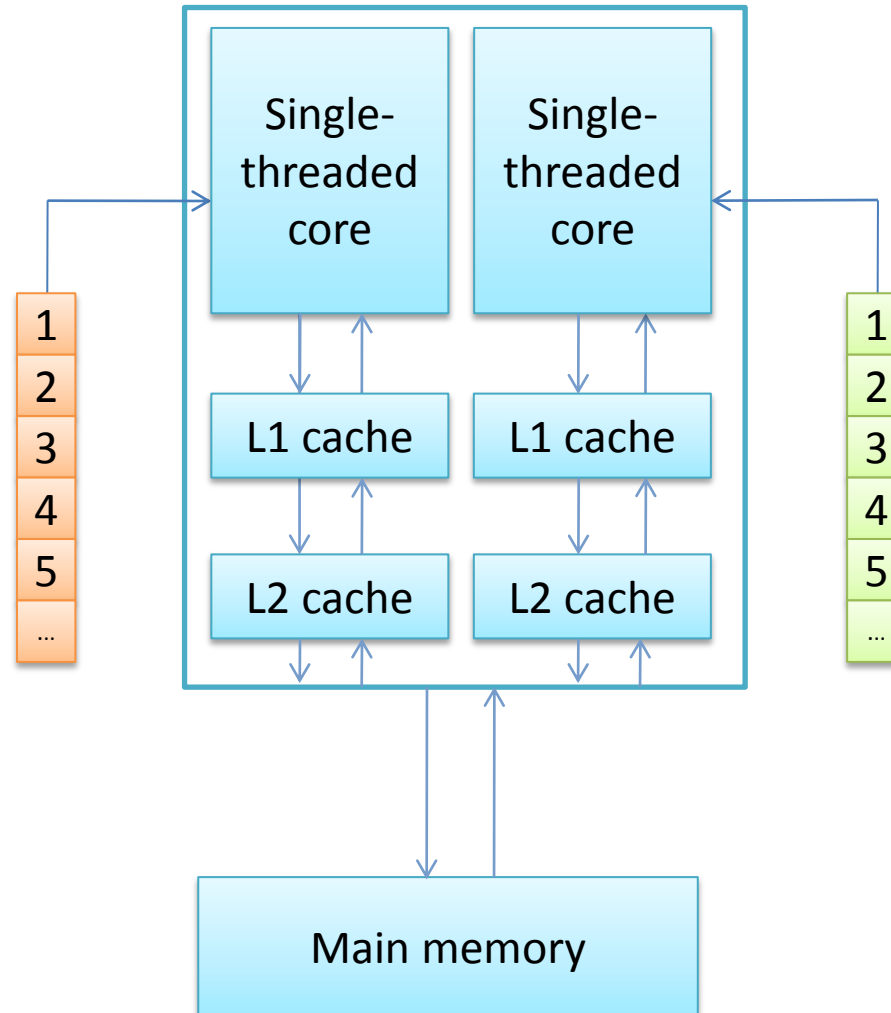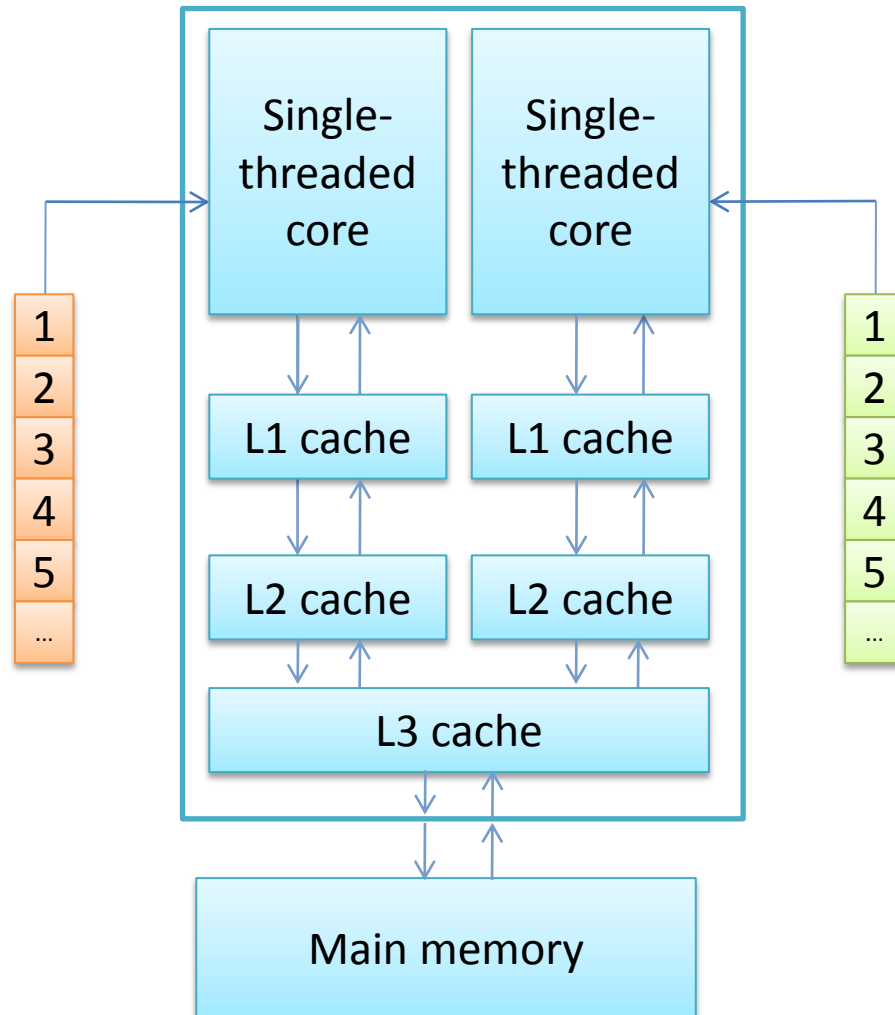
# Multi-core h/w – common L2

# Multi-core h/w – common L2
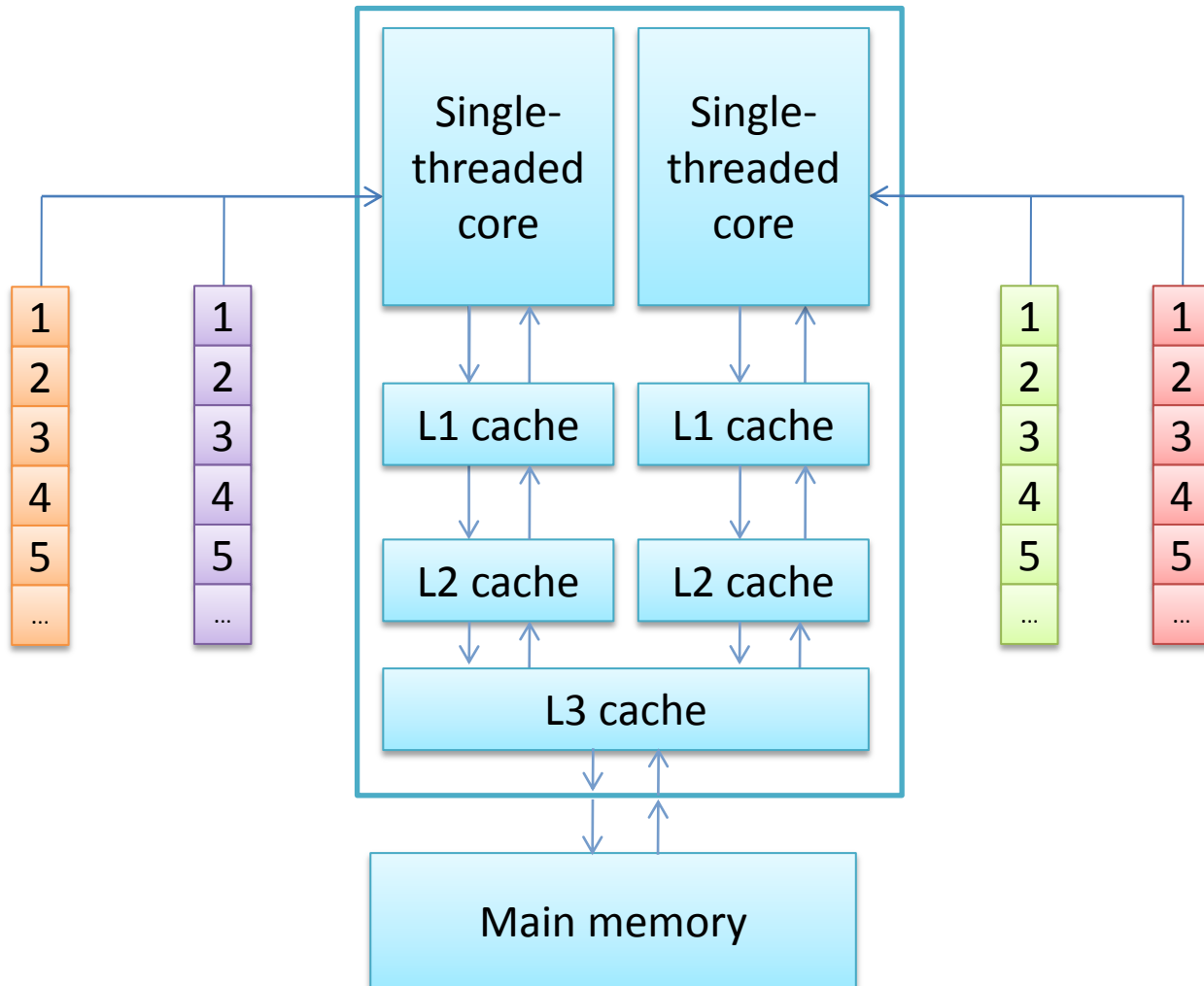
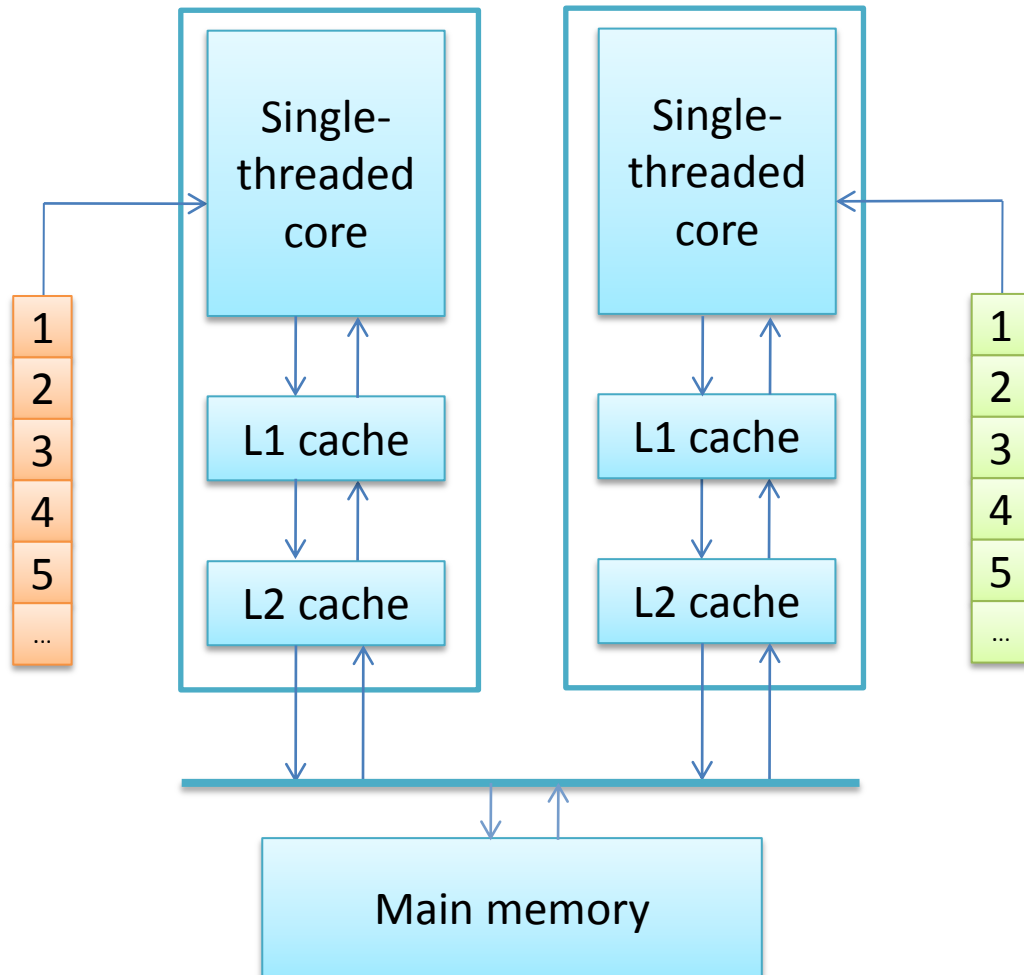# Multi-core h/w – common L2

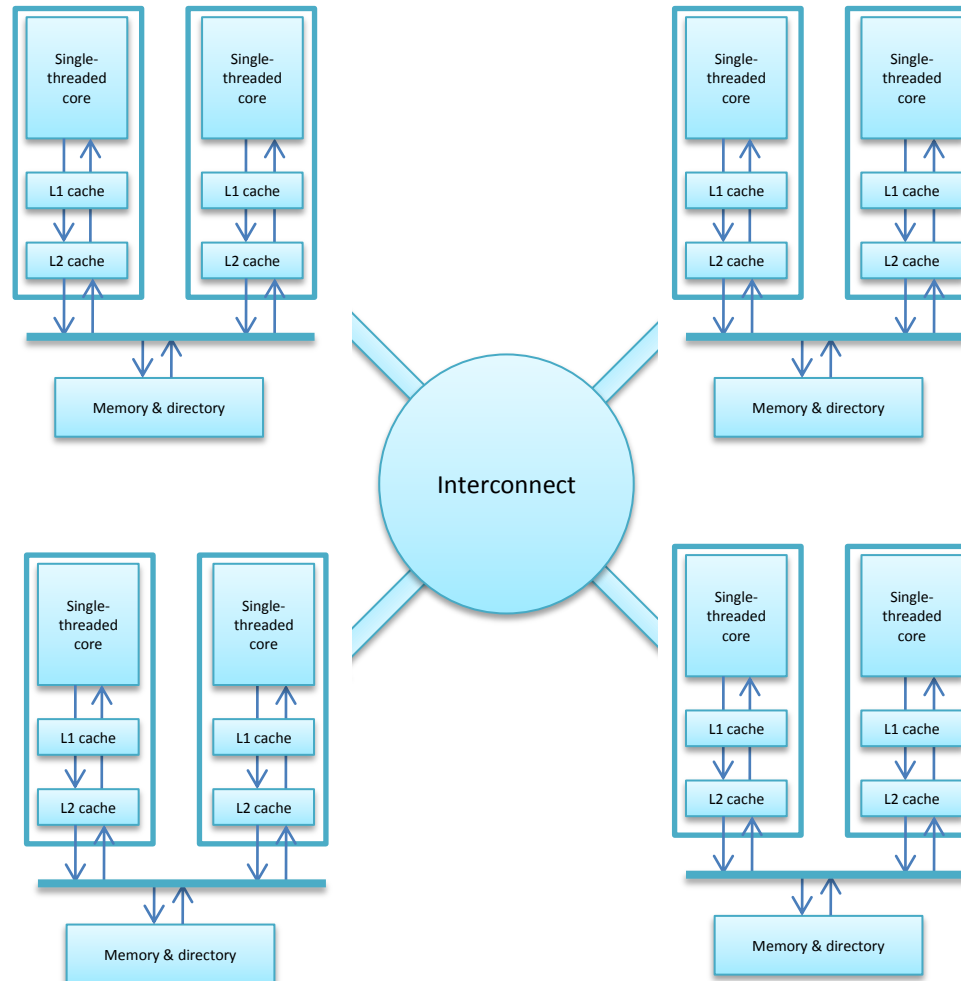# Multi-core h/w – separate L2

# Multi-core h/w – additional L3

Multi-threaded multi-core h/w

# SMP multiprocessor

# NUMA multiprocessor

# Three kinds of parallel hardware

- Multi-threaded cores
  - Increase utilization of a core or memory b/w
  - Peak ops/cycle fixed
- Multiple cores
  - Increase ops/cycle
  - Don't necessarily scale caches and off-chip resources proportionately
- Multi-processor machines
  - Increase ops/cycle
  - Often scale cache & memory capacities and b/w proportionately

# AMD Phenom