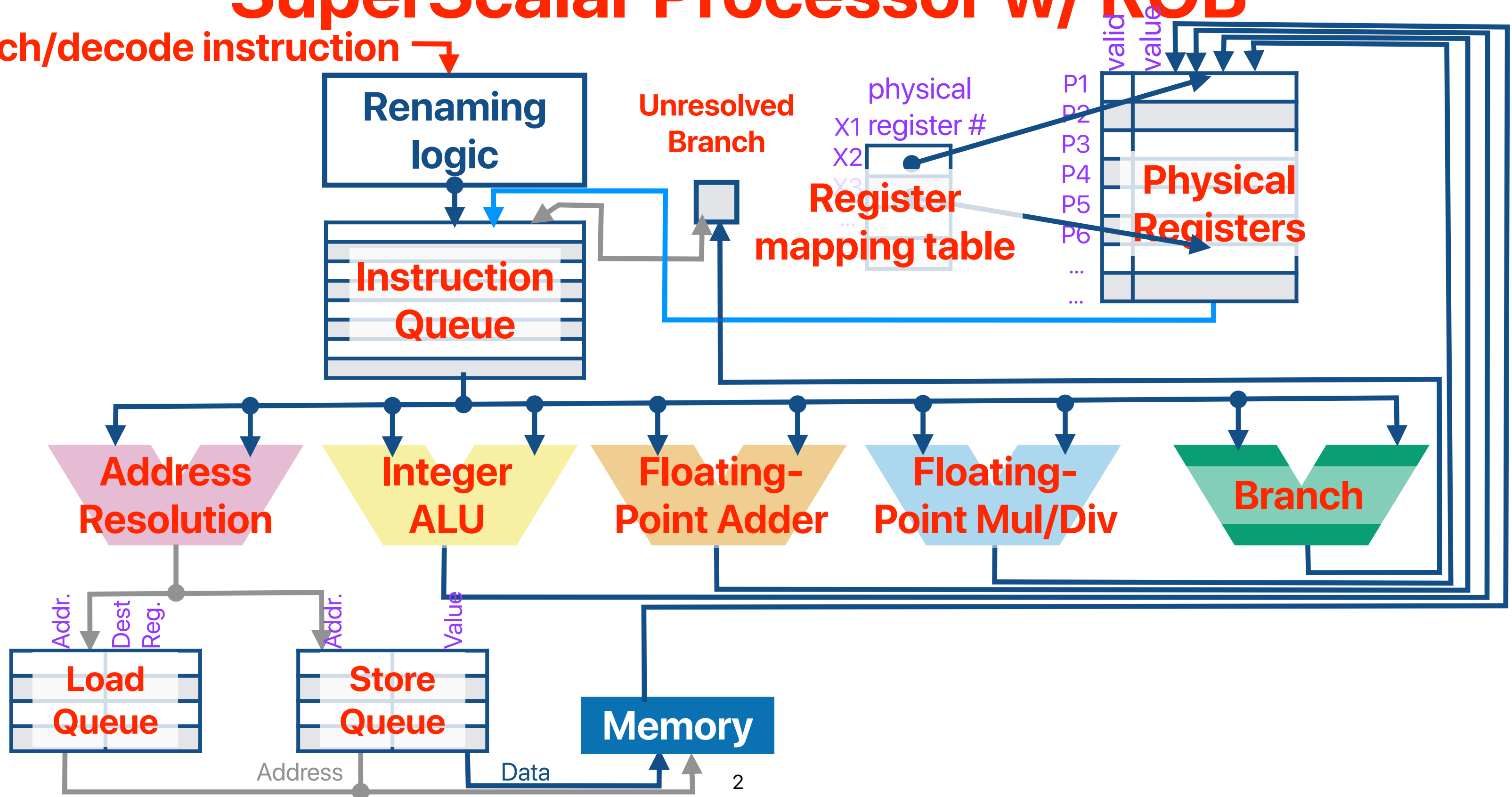


Multithreaded Architectures and Programming on Multithreaded Architectures

Hung-Wei

SuperScalar Processor w/ ROB

Fetch/decode instruction →



Recap: What about "linked list"

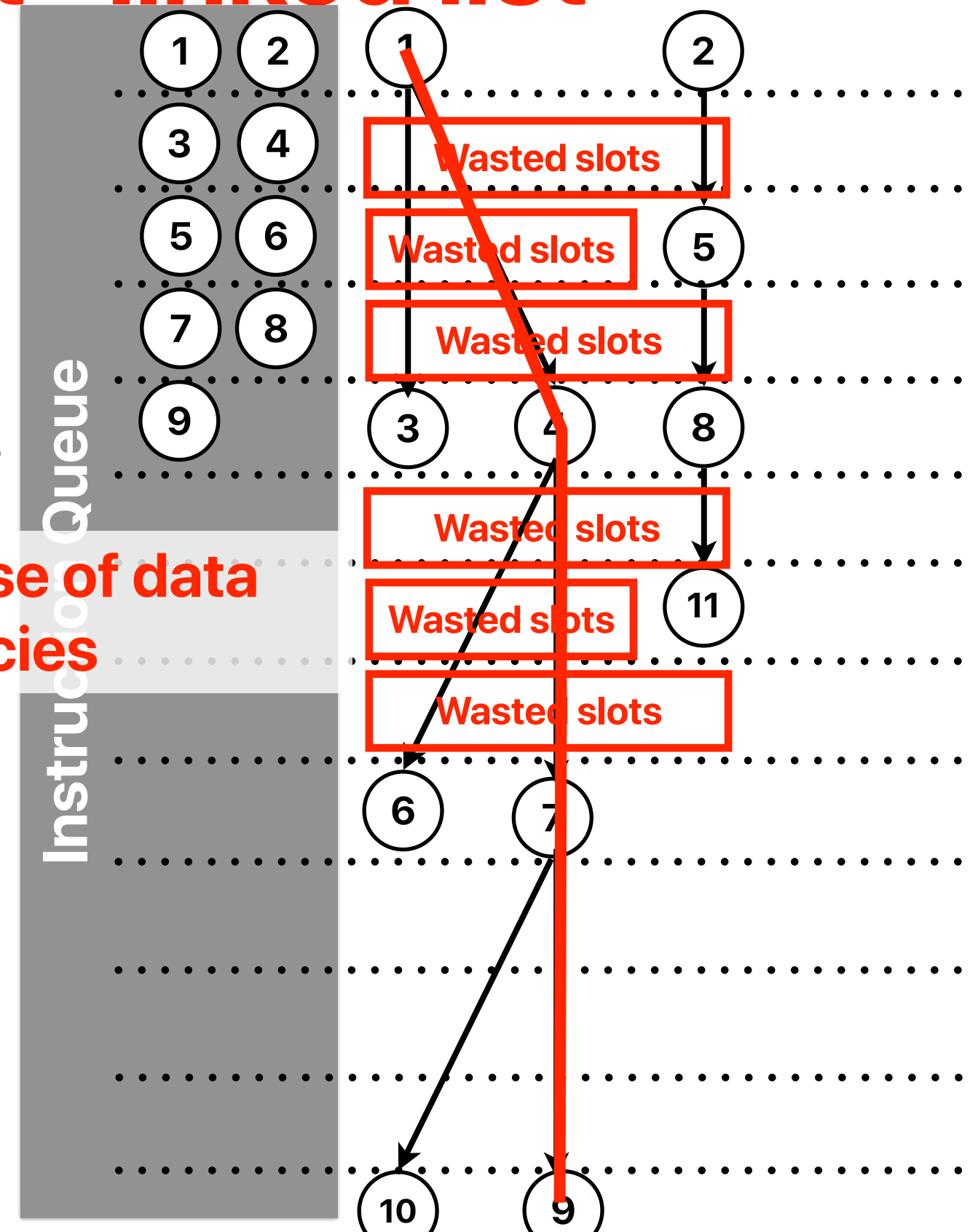
Static instructions

```
LOOP: ld    X10, 8(X10)
      addi  X7, X7, 1
      bne   X10, X0, LOOP
```

Dynamic instructions

```
① ld    X10, 8(X10)
② addi  X7, X7, 1
③ bne   X10, X0, LOOP
④ ld    X10, 8(X10)
⑤ addi  X7, X7, 1
⑥ bne   X10, X0, LOOP
⑦ ld    X10, 8(X10)
⑧ addi  X7, X7, 1
⑨ bne   X10, X0, LOOP
```

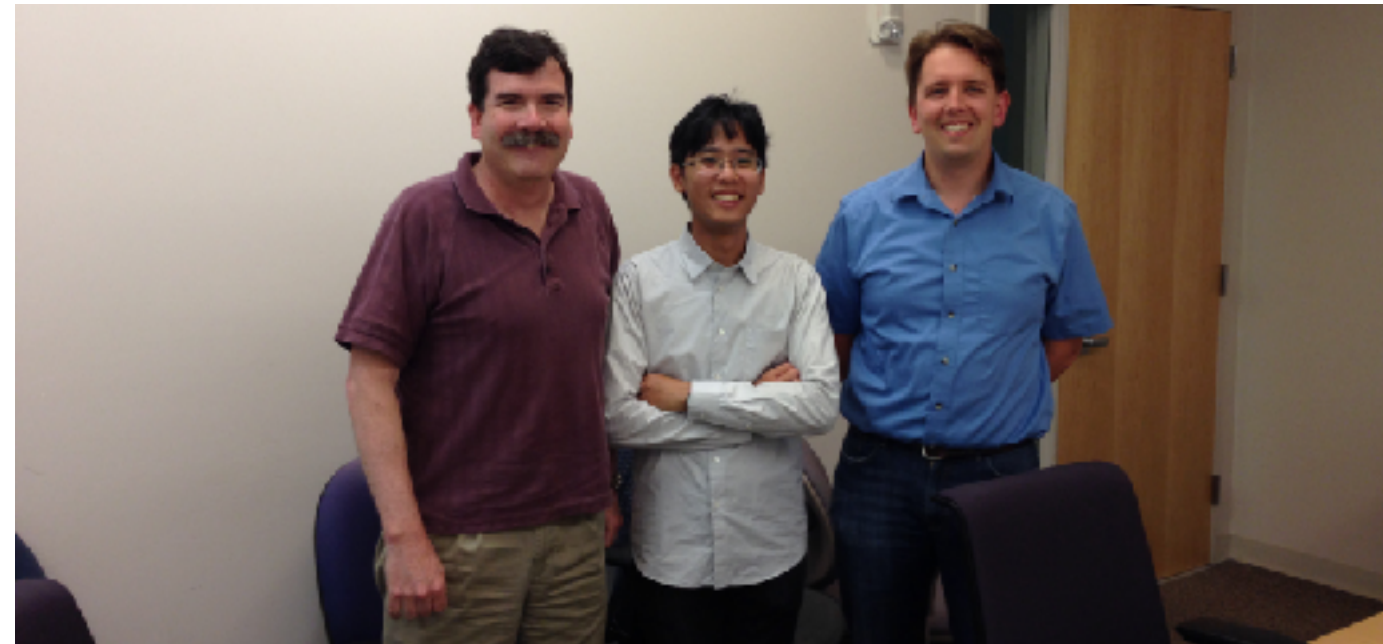
ILP is low because of data dependencies



Demo: ILP within a program

- perf is a tool that captures performance counters of your processors and can generate results like branch mis-prediction rate, cache miss rates and ILP.

Simultaneous multithreading

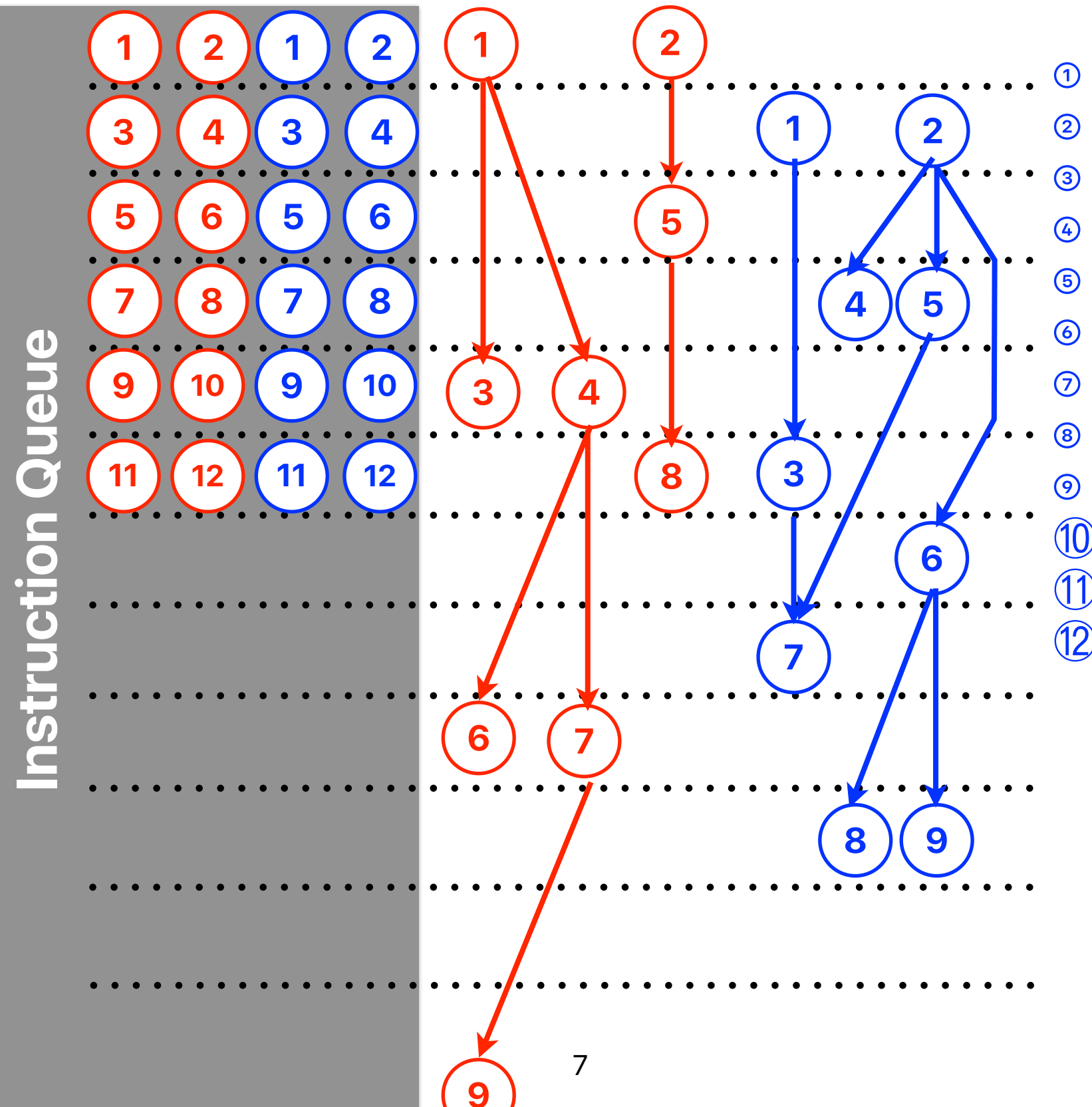


Simultaneous multithreading

- The processor can schedule instructions from different threads/processes/programs
- Fetch instructions from different threads/processes to fill the not utilized part of pipeline
 - Exploit “thread level parallelism” (TLP) to solve the problem of insufficient ILP in a single thread
 - You need to create an illusion of multiple processors for OSs

Simultaneous multithreading

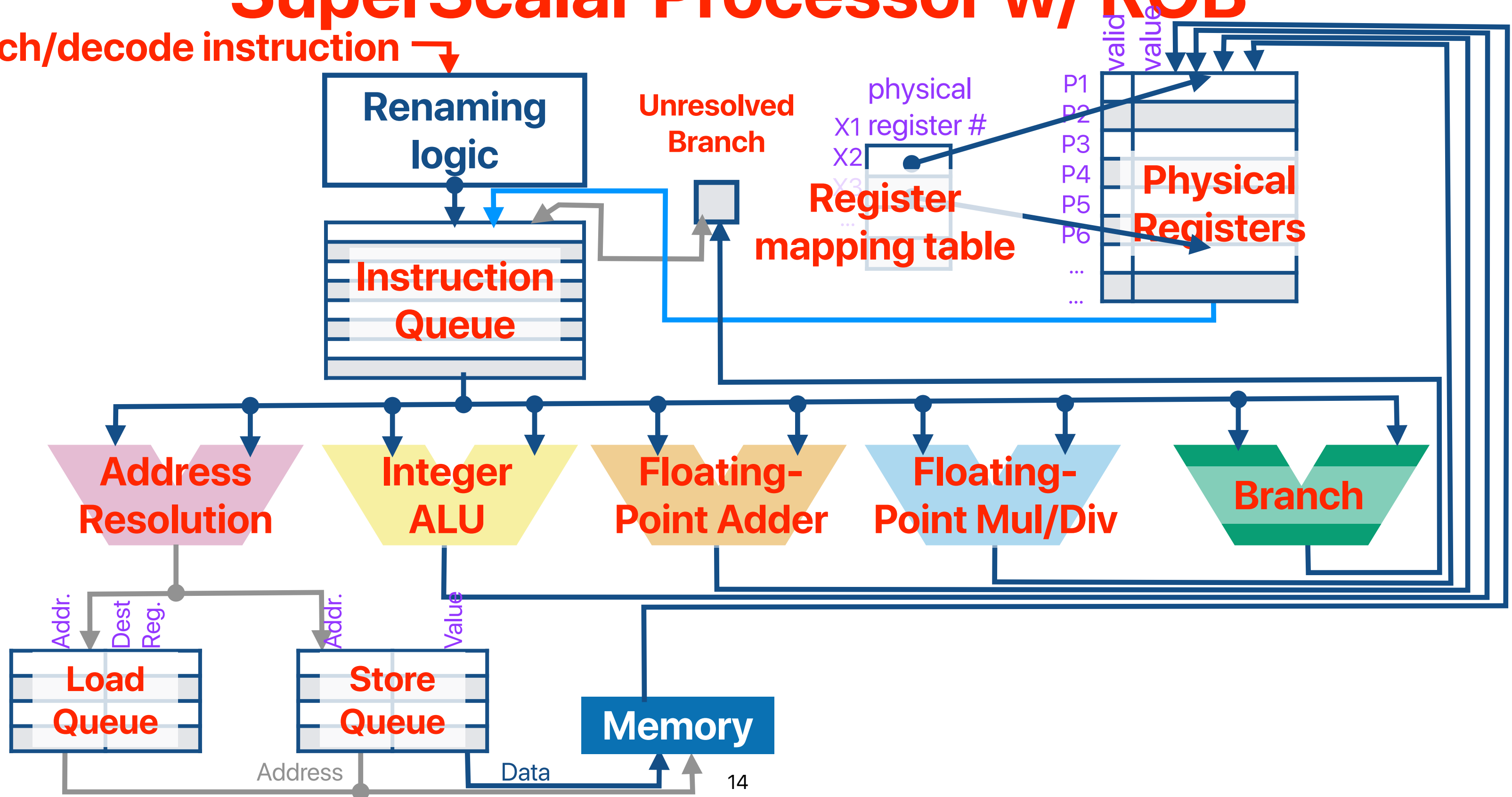
① ld X10, 8(X10)
② addi X7, X7, 1
③ bne X10, X0, LOOP
④ ld X10, 8(X10)
⑤ addi X7, X7, 1
⑥ bne X10, X0, LOOP
⑦ ld X10, 8(X10)
⑧ addi X7, X7, 1
⑨ bne X10, X0, LOOP



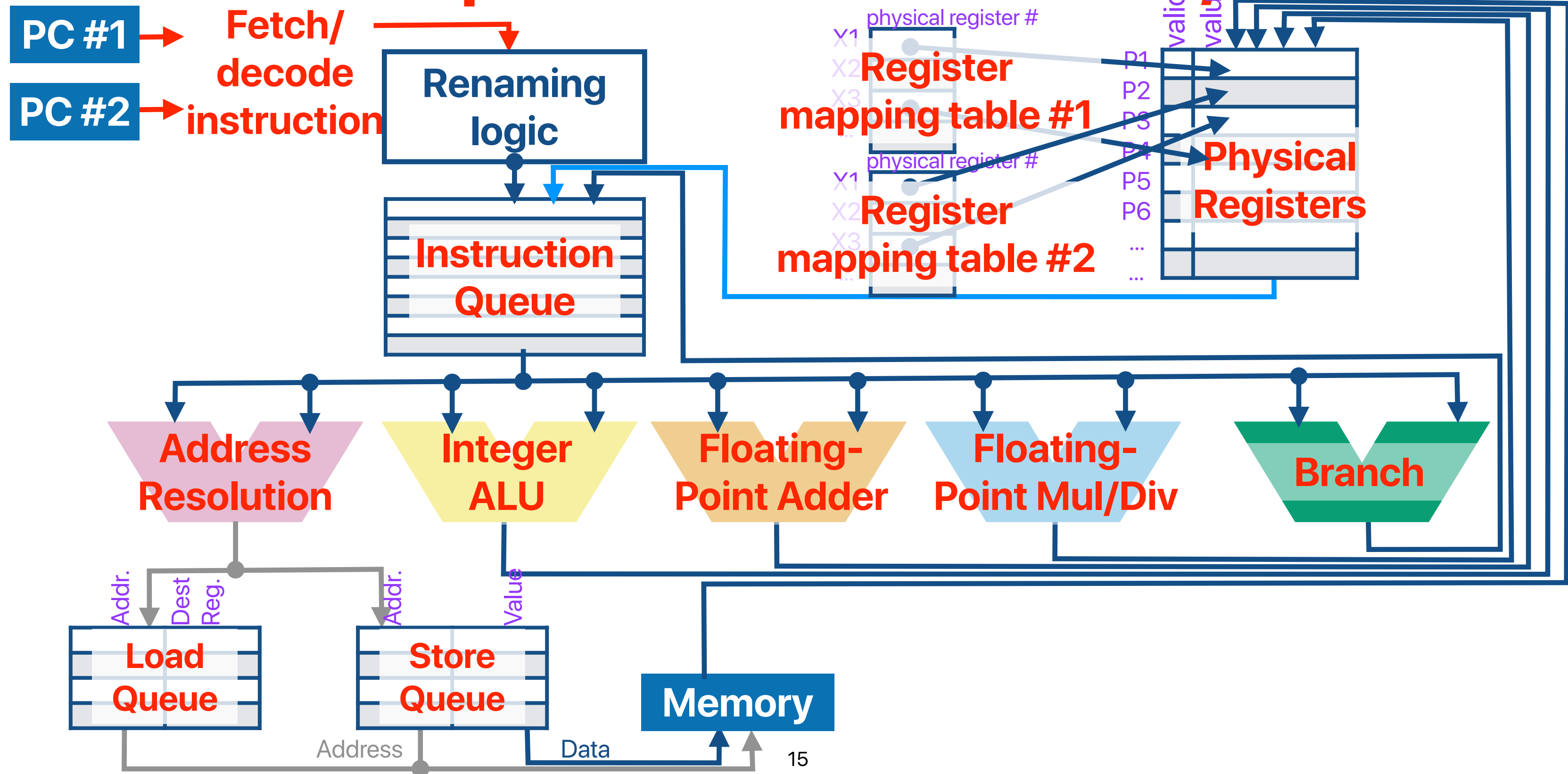
① ld X1, 0(X10)
② addi X10, X10, 8
③ add X20, X20, X1
④ bne X10, X2, LOOP
⑤ ld X1, 0(X10)
⑥ addi X10, X10, 8
⑦ add X20, X20, X1
⑧ bne X10, X2, LOOP
⑨ ld X1, 0(X10)
⑩ addi X10, X10, 8
⑪ add X20, X20, X1
⑫ bne X10, X2, LOOP

SuperScalar Processor w/ ROB

Fetch/decode instruction →



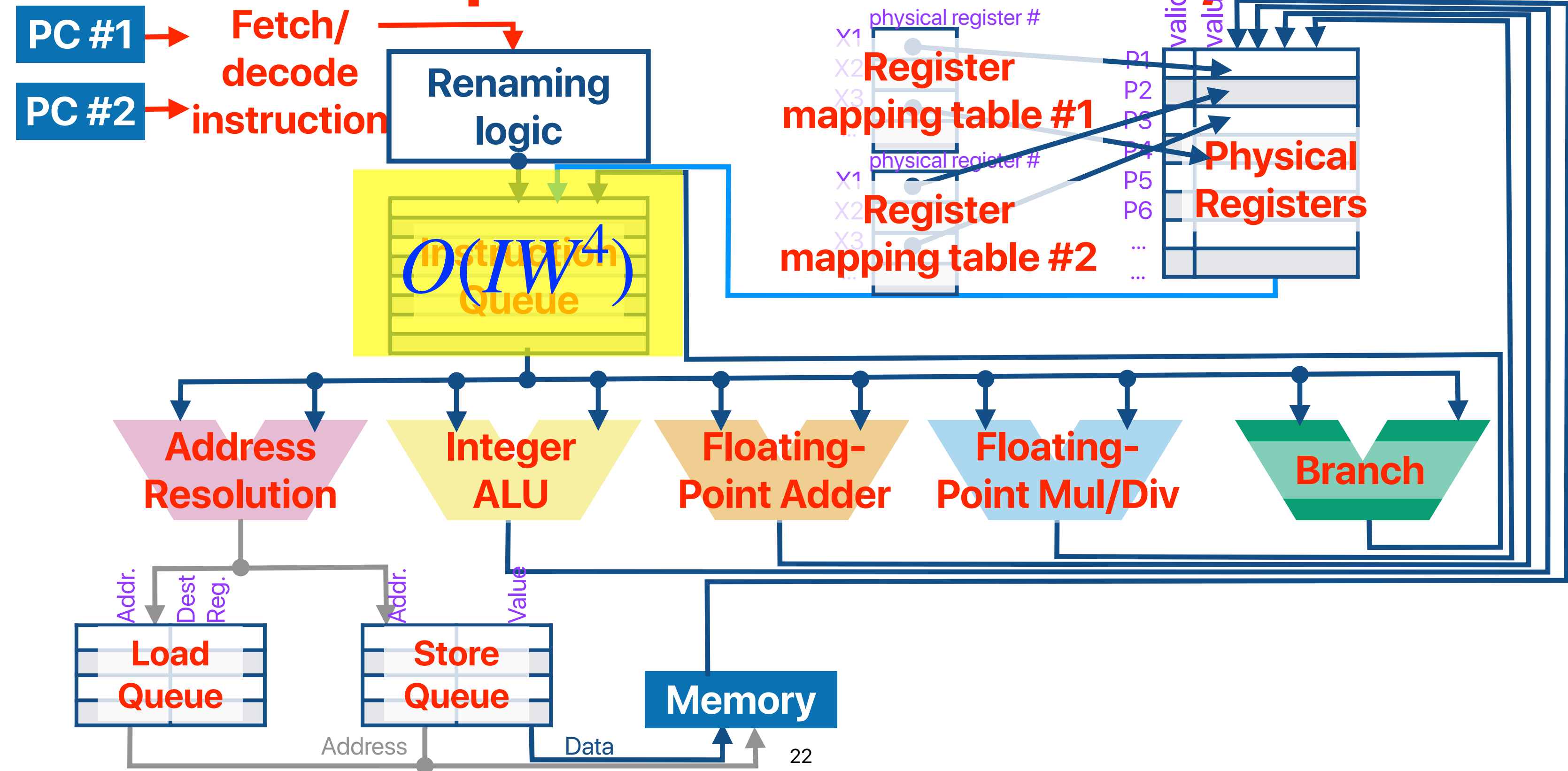
SMT SuperScalar Processor w/ ROB



SMT

- Improve the throughput of execution
 - May increase the latency of a single thread
- Less branch penalty per thread
- Increase hardware utilization
- Simple hardware design: Only need to duplicate PC/Register Files
- Real Case:
 - Intel HyperThreading (supports up to two threads per core)
 - Intel Pentium 4, Intel Atom, Intel Core i7
 - AMD RyZen (Zen microarchitecture)

SMT SuperScalar Processor w/ ROB



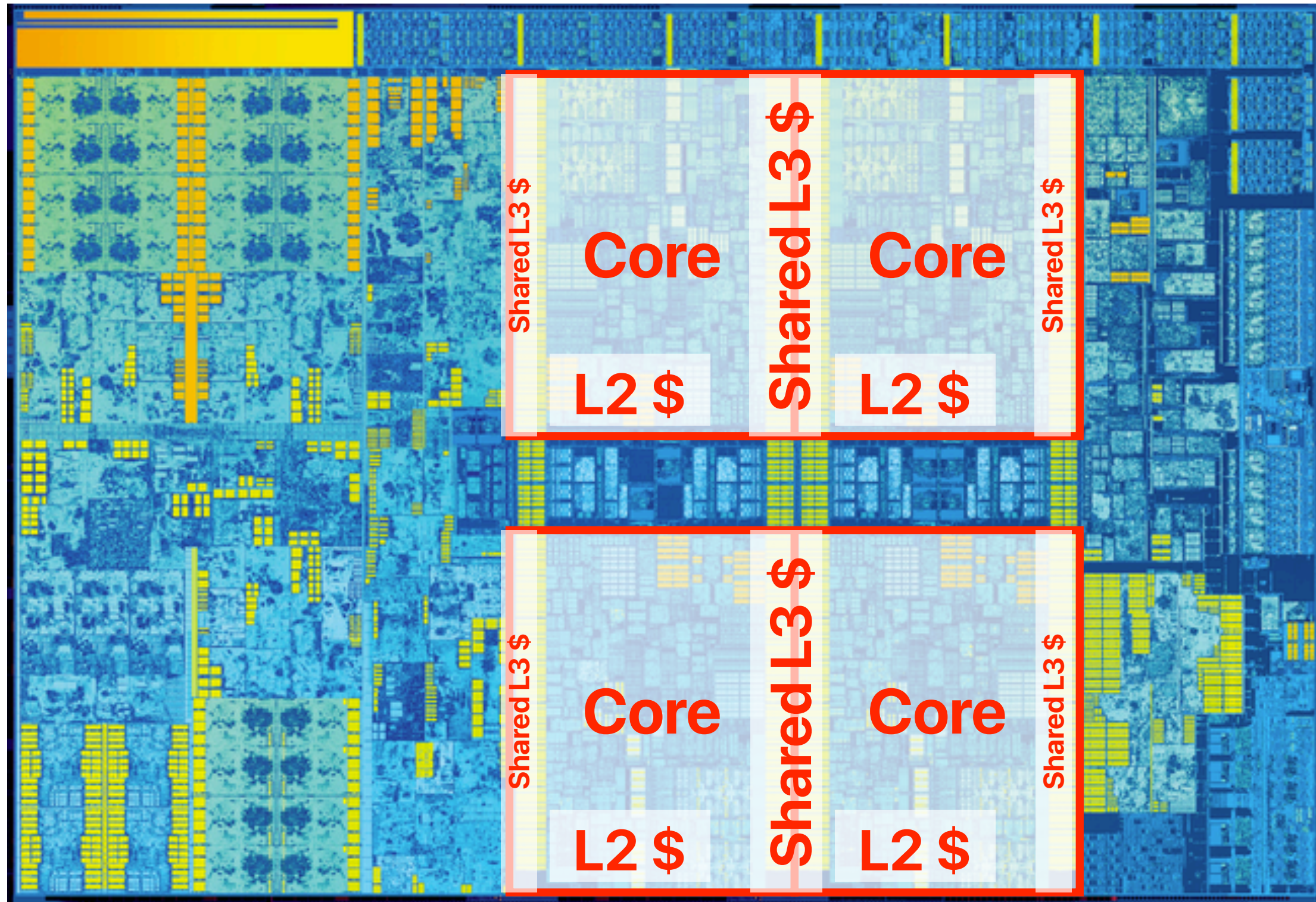
Wider-issue processors won't give you much more

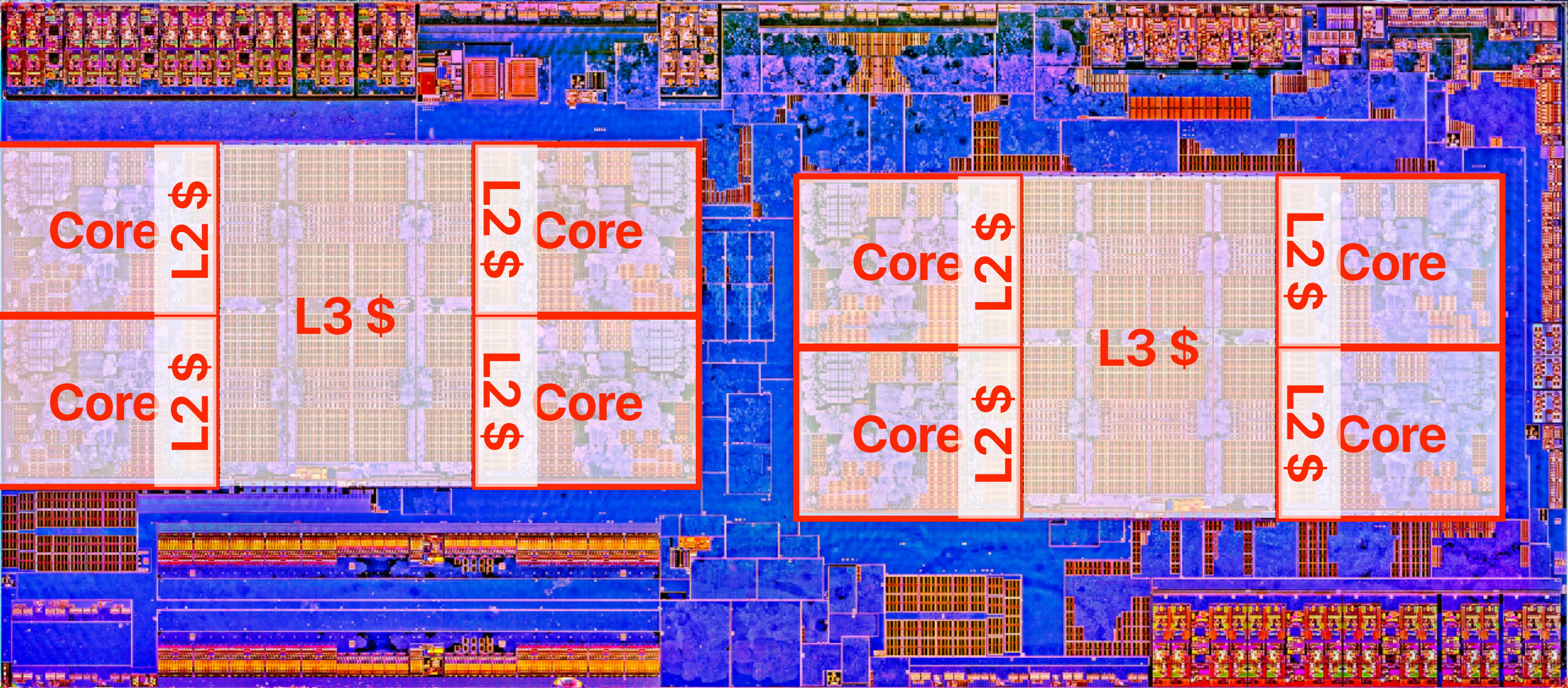
Program	IPC	BP Rate %	I cache %MPCI	D cache %MPCI	L2 cache %MPCI
compress	0.9	85.9	0.0	3.5	1.0
eqntott	1.3	79.8	0.0	0.8	0.7
m88ksim	1.4	91.7	2.2	0.4	0.0
MPsim	0.8	78.7	5.1	2.3	2.3
applu	0.9	79.2	0.0	2.0	1.7
apsi	0.6	95.1	1.0	4.1	2.1
swim	0.9	99.7	0.0	1.2	1.2
tomcatv	0.8	99.6	0.0	7.7	2.2
pmake	1.0	86.2	2.3	2.1	0.4

Program	IPC	BP Rate %	I cache %MPCI	D cache %MPCI	L2 cache %MPCI
compress	1.2	86.4	0.0	3.9	1.1
eqntott	1.8	80.0	0.0	1.1	1.1
m88ksim	2.3	92.6	0.1	0.0	0.0
MPsim	1.2	81.6	3.4	1.7	2.3
applu	1.7	79.7	0.0	2.8	2.8
apsi	1.2	95.6	0.2	3.1	2.6
swim	2.2	99.8	0.0	2.3	2.5
tomcatv	1.3	99.7	0.0	4.2	4.3
pmake	1.4	82.7	0.7	1.0	0.6

Table 5. Performance of a single 2-issue superscalar processor. Table 6. Performance of the 6-issue superscalar processor.

Intel SkyLake





Core \$
L2

Core \$
L2

L3 \$

Core \$
L2

Core \$
L2

Core \$
L2

Core \$
L2

L3 \$

Core \$
L2

Core \$
L2

AMD

RYZEN

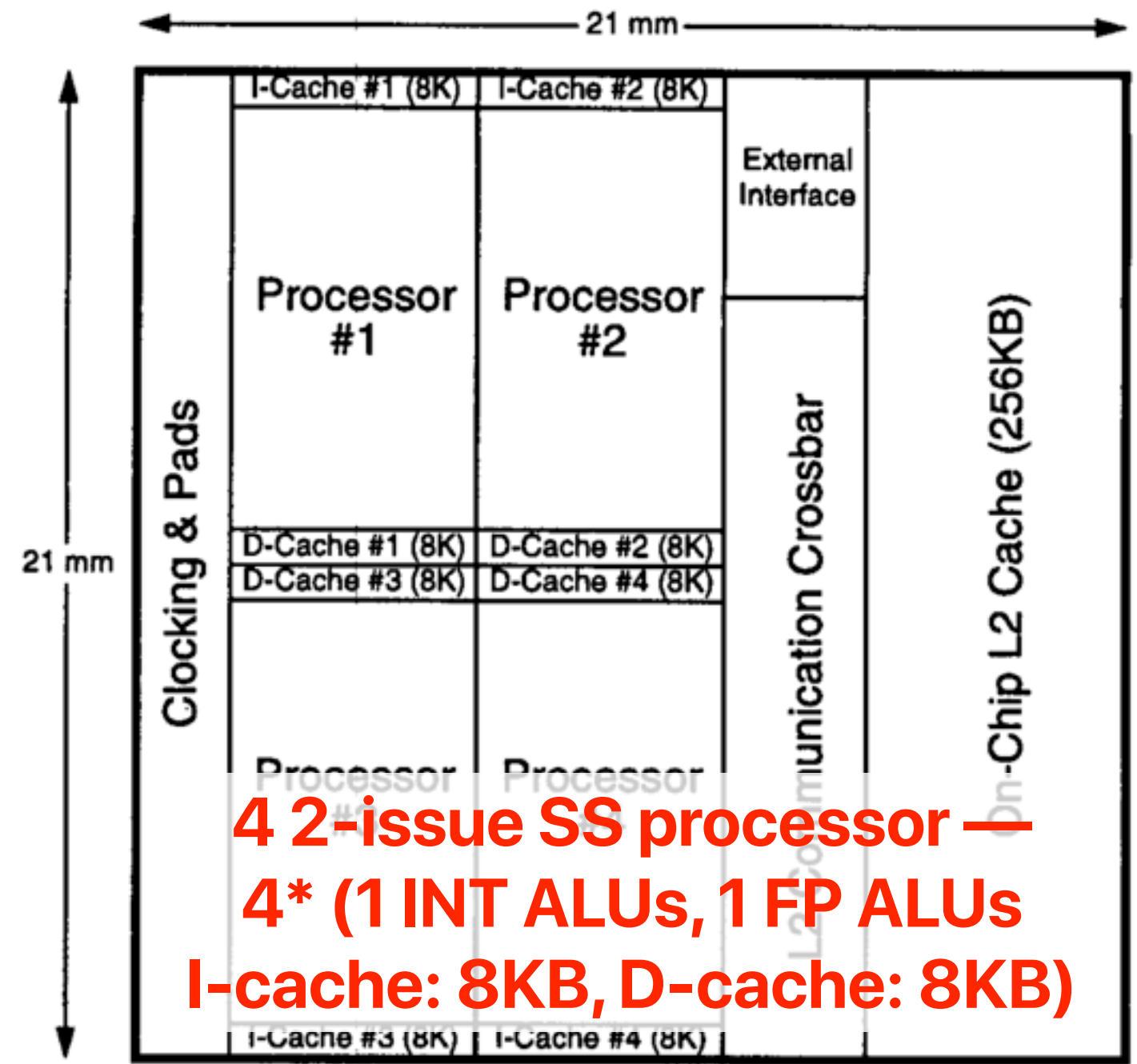
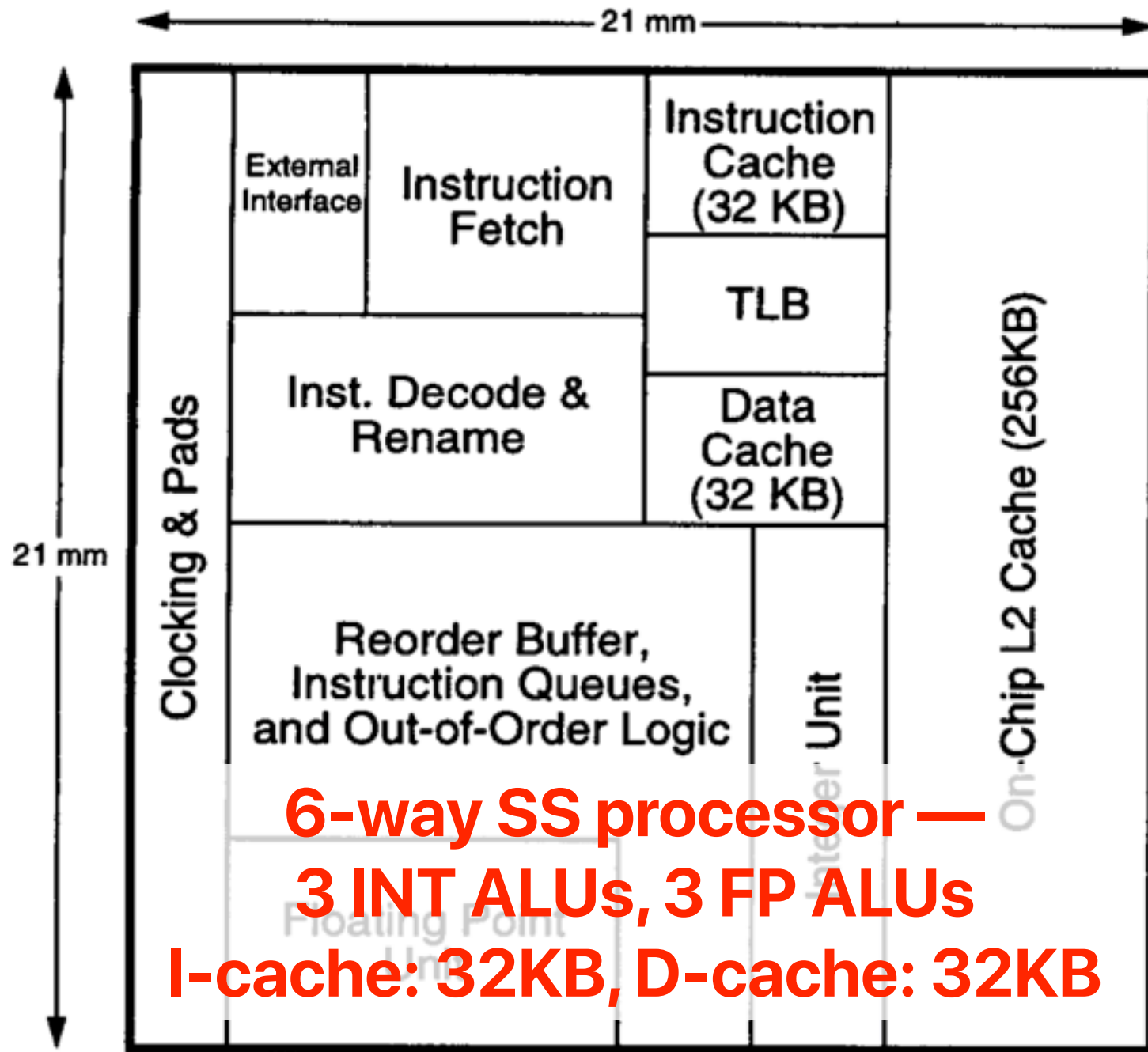
The case for a Single-Chip Multiprocessor

Kunle Olukotun, Basem A. Nayfeh, Lance Hammond, Ken Wilson, and Kunyung

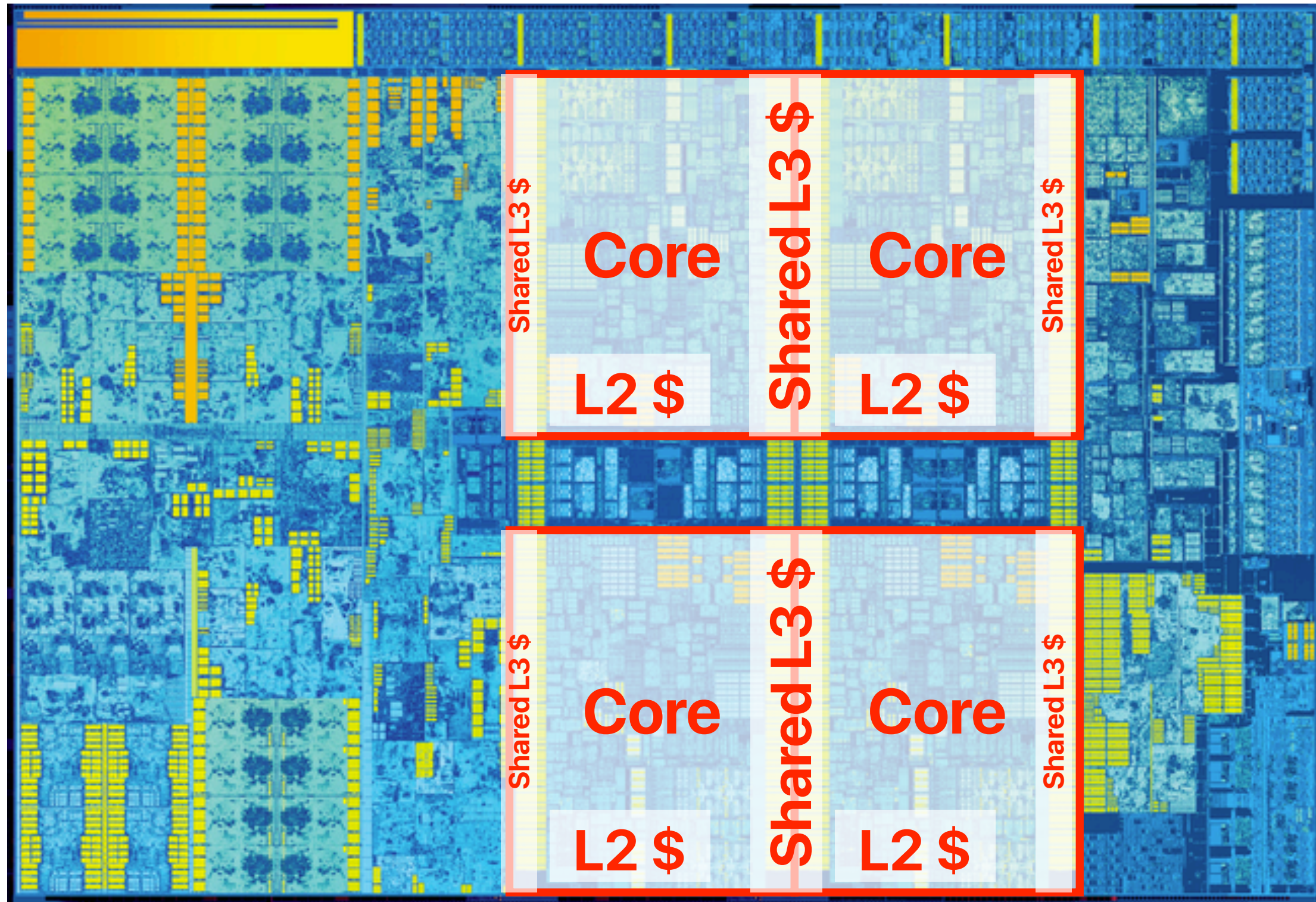
Chang

Stanford University

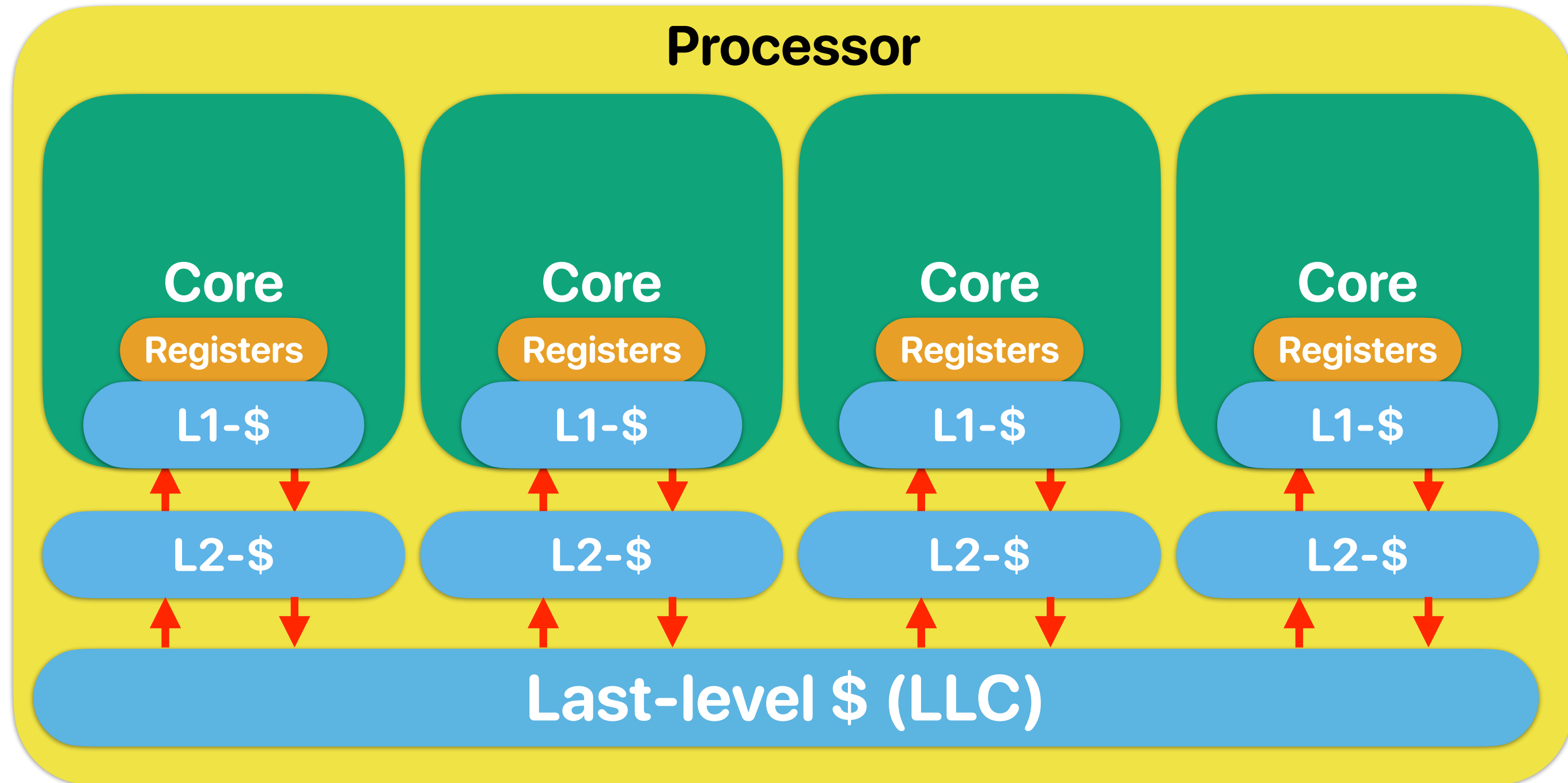
Wide-issue SS processor v.s. multiple narrower-issue SS processors



Intel SkyLake



Concept of CMP

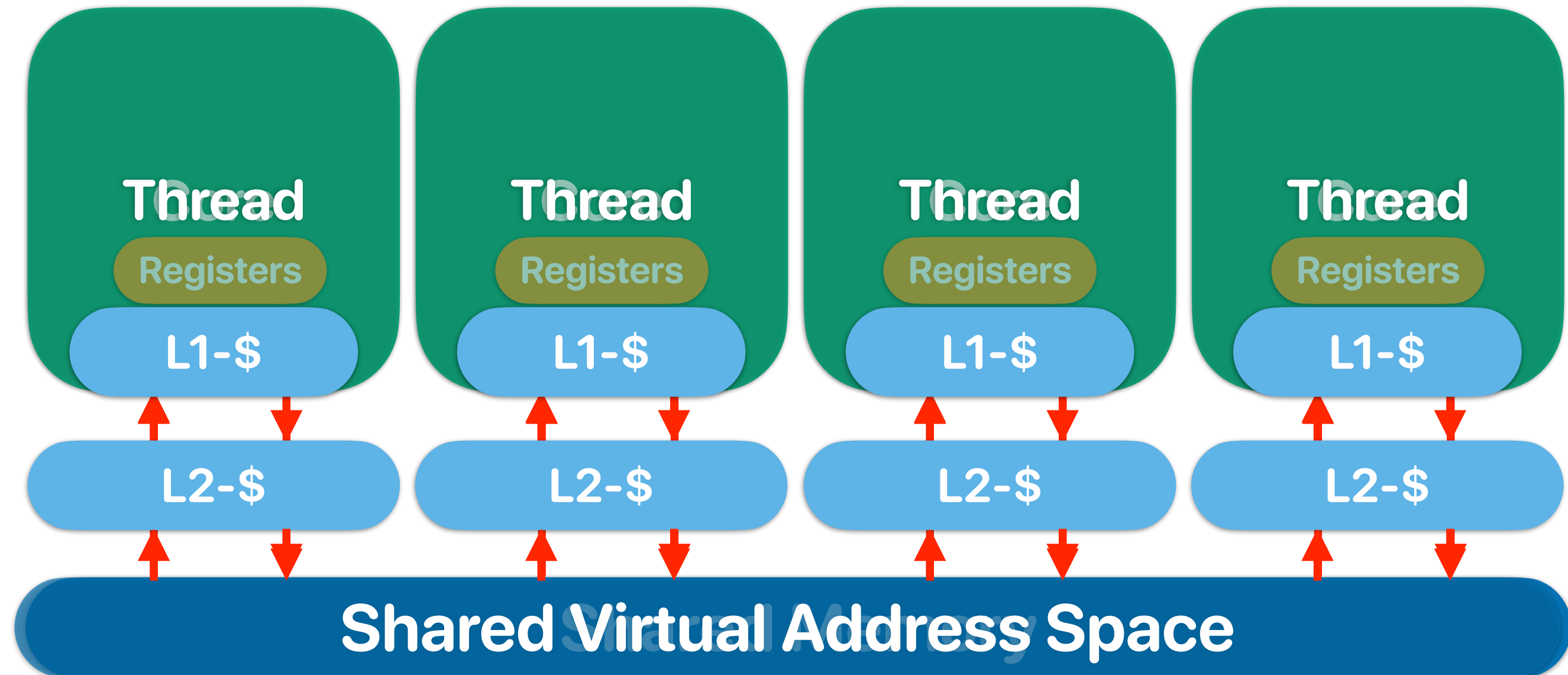


Architectural Support for Parallel Programming

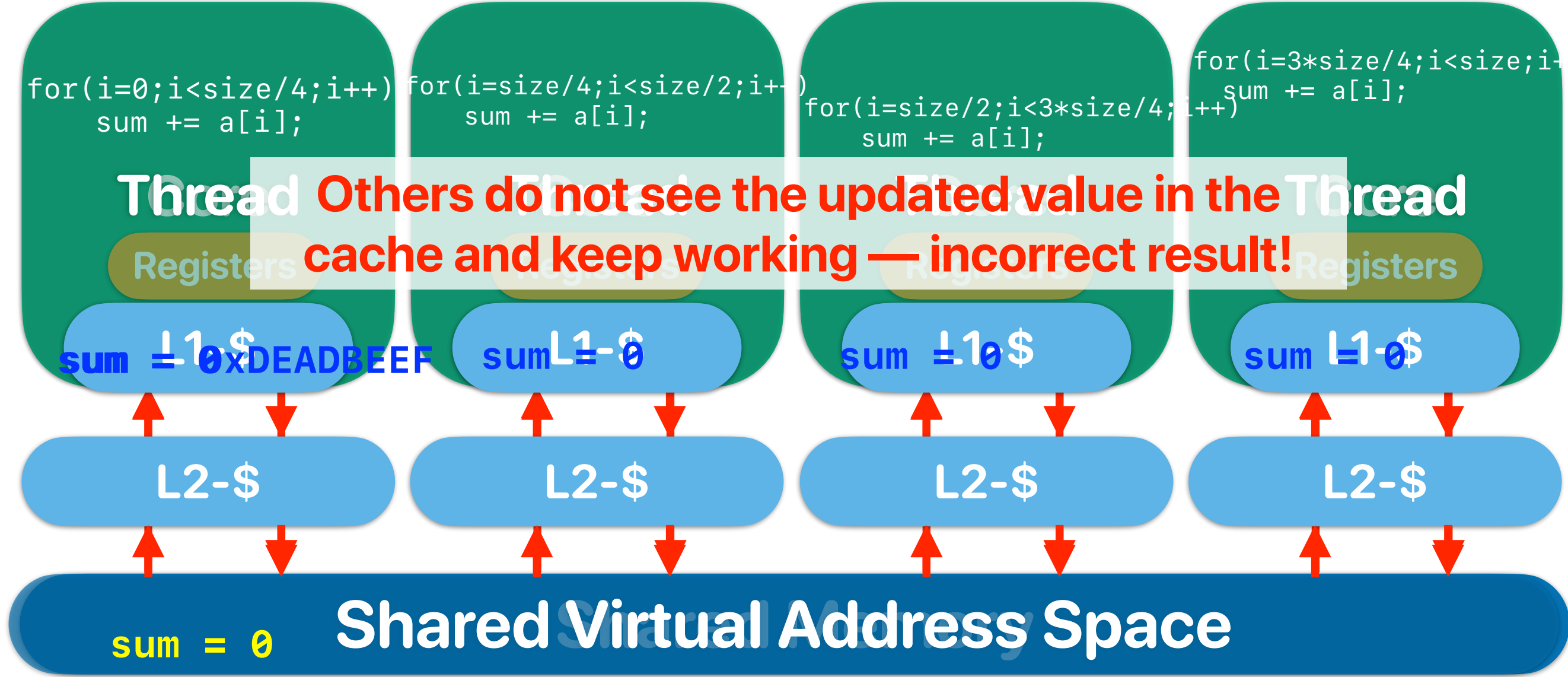
Parallel programming

- To exploit parallelism you need to break your computation into multiple "processes" or multiple "threads"
- Processes (in OS/software systems)
 - Separate programs actually running (not sitting idle) on your computer at the same time.
 - Each process will have its own virtual memory space and you need explicitly exchange data using inter-process communication APIs
- Threads (in OS/software systems)
 - Independent portions of your program that can run in parallel
 - All threads share the same virtual memory space
- We will refer to these collectively as "threads"
 - A typical user system might have 1-8 actively running threads.
 - Servers can have more if needed (the sysadmins will hopefully configure it that way)

What software thinks about "multiprogramming" hardware



What software thinks about “multiprogramming” hardware



Coherency & Consistency

- Coherency — Guarantees all processors see the same value for a variable/memory address in the system when the processors need the value at the same time
 - What value should be seen
- Consistency — All threads see the change of data in the same order
 - When the memory operation should be done

Simple cache coherency protocol

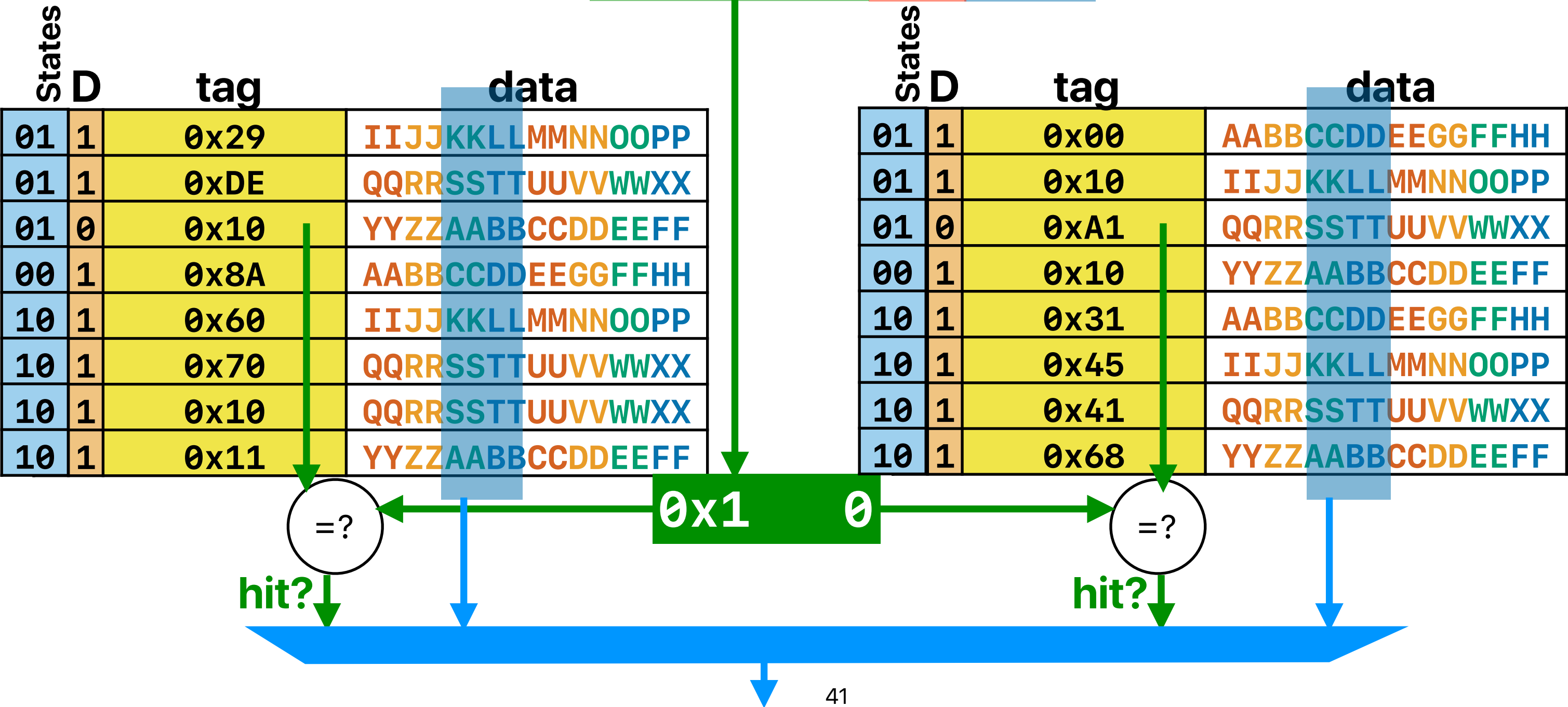
- Snooping protocol
 - Each processor broadcasts / listens to cache misses
- State associate with each block (cacheline)
 - Invalid
 - The data in the current block is invalid
 - Shared
 - The processor can read the data
 - The data may also exist on other processors
 - Exclusive
 - The processor has full permission on the data
 - The processor is the only one that has up-to-date data

Coherent way-associative cache

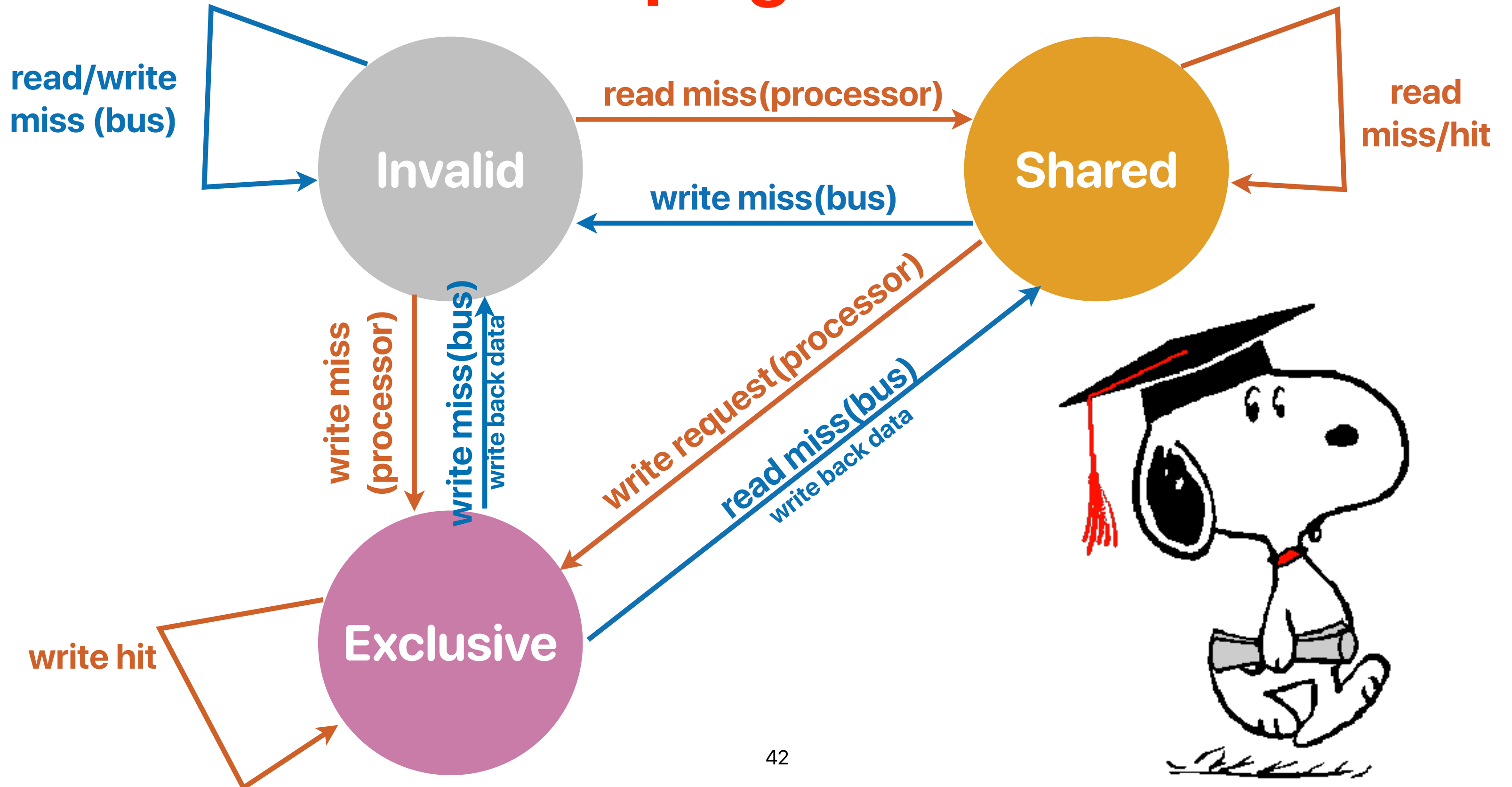
memory address: 0x0

memory address: 0b00001000000100100

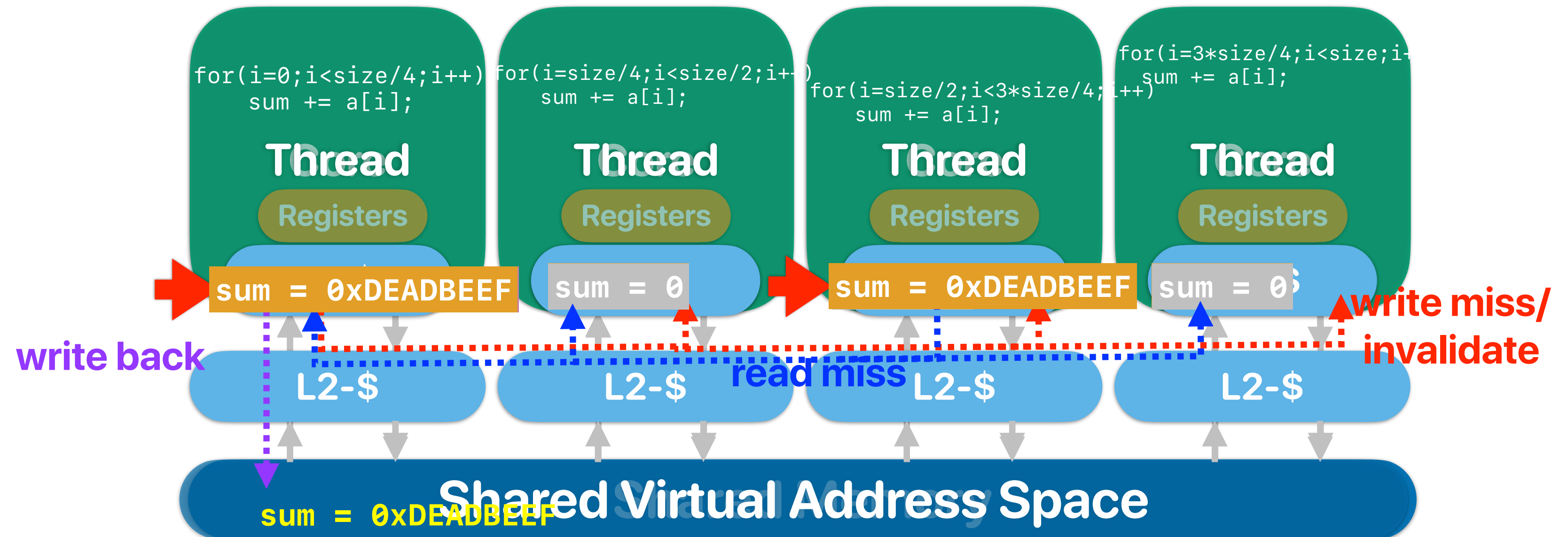
8 tag 2 set 4 block
index offset



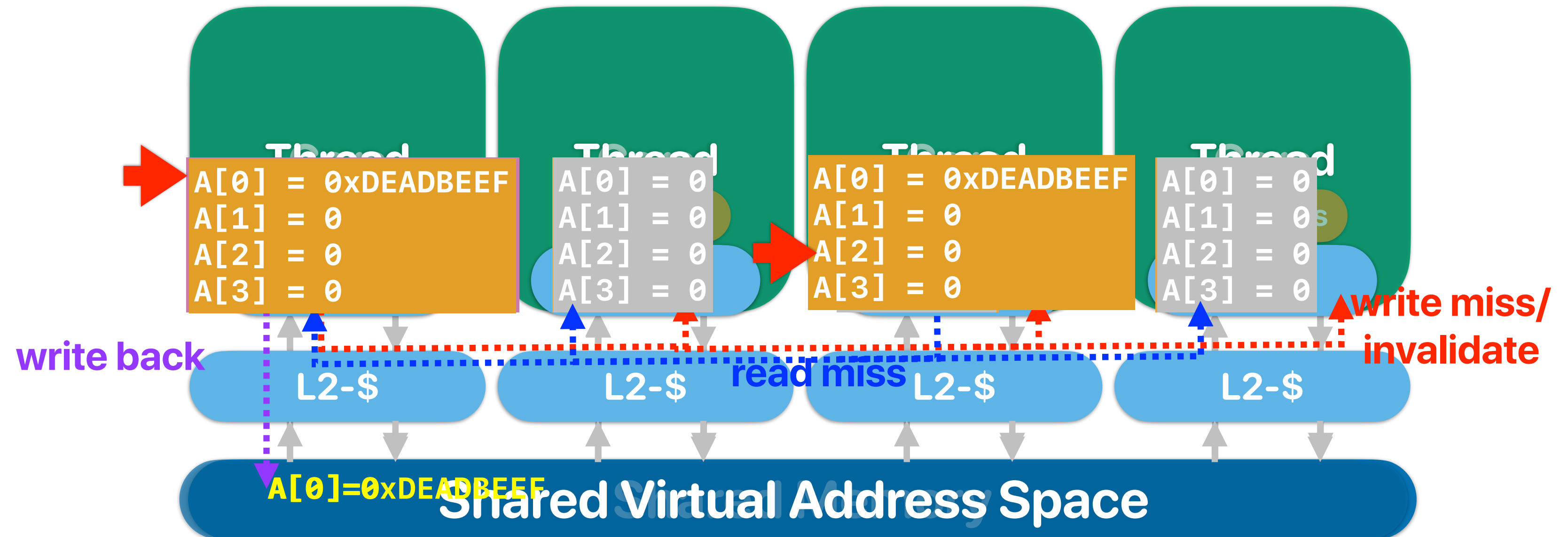
Snooping Protocol



What happens when we write in coherent caches?



Cache coherency



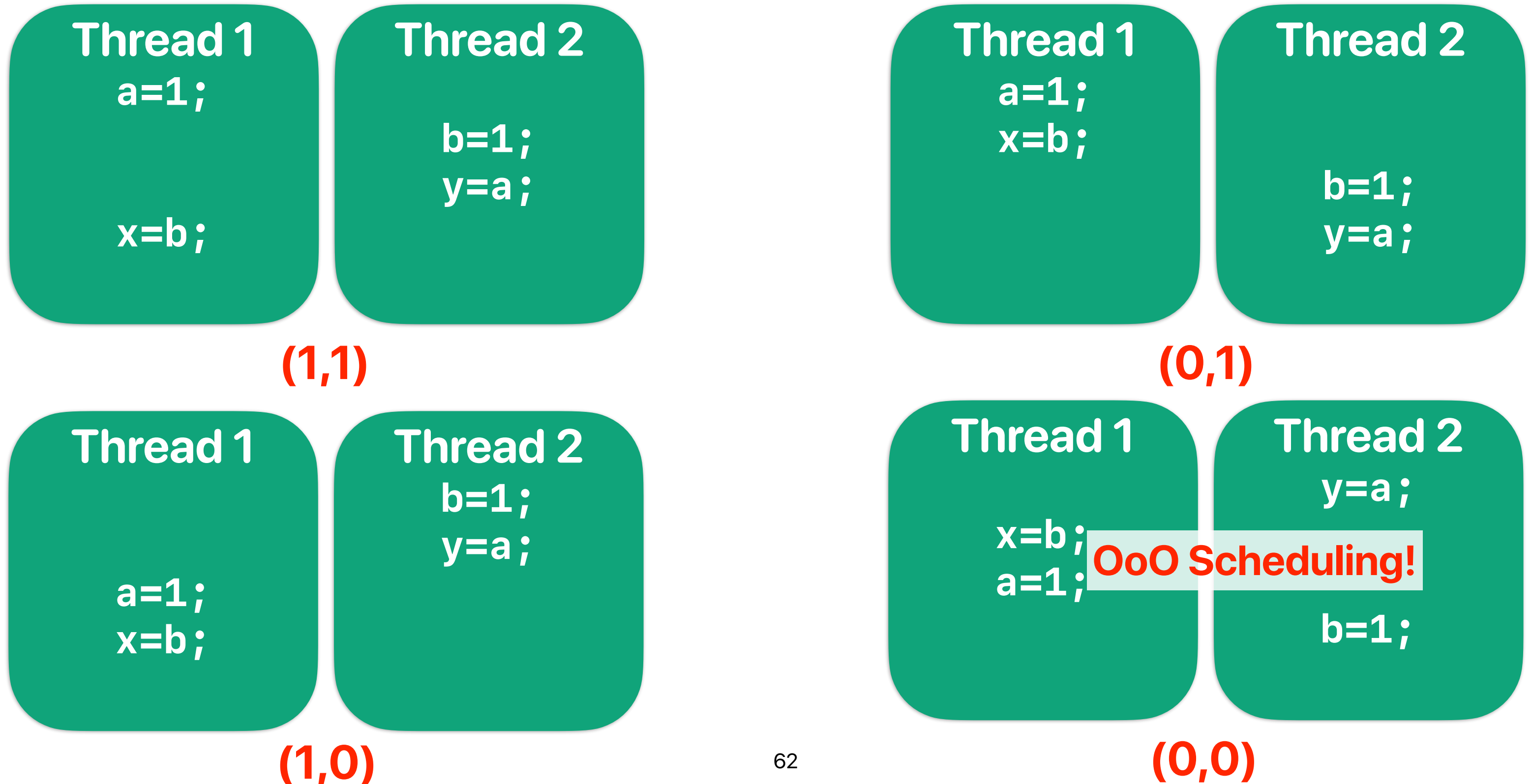
4Cs of cache misses

- 3Cs:
 - Compulsory, Conflict, Capacity
- Coherency miss:
 - A "block" invalidated because of the sharing among processors.

False sharing

- True sharing
 - Processor A modifies X, processor B also want to access X.
- False sharing
 - Processor A modifies X, processor B also want to access Y.
However, Y is invalidated because X and Y are in the same block!

Possible scenarios



Why (0,0)?

- Processor/compiler may reorder your memory operations/instructions
 - Coherence protocol can only guarantee the update of the same memory address
 - Processor can serve memory requests without cache miss first
 - Compiler may store values in registers and perform memory operations later
- Each processor core may not run at the same speed (cache misses, branch mis-prediction, I/O, voltage scaling and etc..)
- Threads may not be executed/scheduled right after it's spawned

Again — how many values are possible?

- Consider the given program. You can safely assume the caches are coherent. How many of the following outputs will you see?

① (0, 0)

② (0, 1)

③ (1, 0)

④ (1, 1)

A. 0

B. 1

C. 2

D. 3

E. 4

```
#include <stdio.h>
#include <stdlib.h>
#include <pthread.h>
#include <unistd.h>

volatile int a,b;
volatile int x,y;
volatile int f;
void* modifya(void *z) {
    a=1;
    x=b;
    return NULL;
}
void* modifyb(void *z) {
    b=1;
    y=a;
    return NULL;
}
```

```
int main() {
    int i;
    pthread_t thread[2];
    pthread_create(&thread[0], NULL, modifya, NULL);
    pthread_create(&thread[1], NULL, modifyb, NULL);
    pthread_join(thread[0], NULL);
    pthread_join(thread[1], NULL);
    fprintf(stderr, "(%d, %d)\n", x, y);
    return 0;
}
```

fence instructions

- x86 provides an "mfence" instruction to prevent reordering across the fence instruction
- x86 only supports this kind of "relaxed consistency" model. You still have to be careful enough to make sure that your code behaves as you expected

thread 1	thread 2
<pre>a=1; mfence a=1 must occur/update before mfence x=b;</pre>	<pre>b=1; mfence b=1 must occur/update before mfence y=a;</pre>

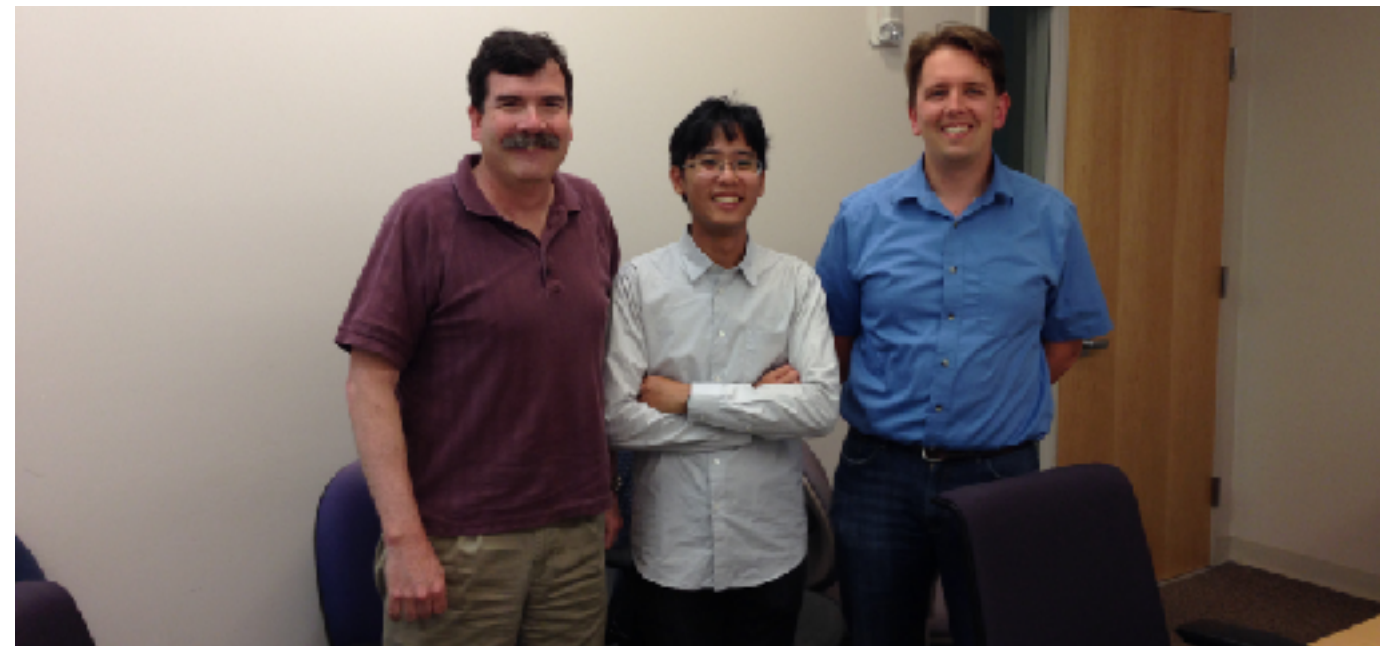
Take-aways of parallel programming

- Processor behaviors are non-deterministic
 - You cannot predict which processor is going faster
 - You cannot predict when OS is going to schedule your thread
- Cache coherency only guarantees that everyone would eventually have a coherent view of data, but not when
- Cache consistency is hard to support

Simultaneous multithreading: maximizing on-chip parallelism

Dean M. Tullsen, Susan J. Eggers, Henry M. Levy

Department of Computer Science and Engineering, University of Washington

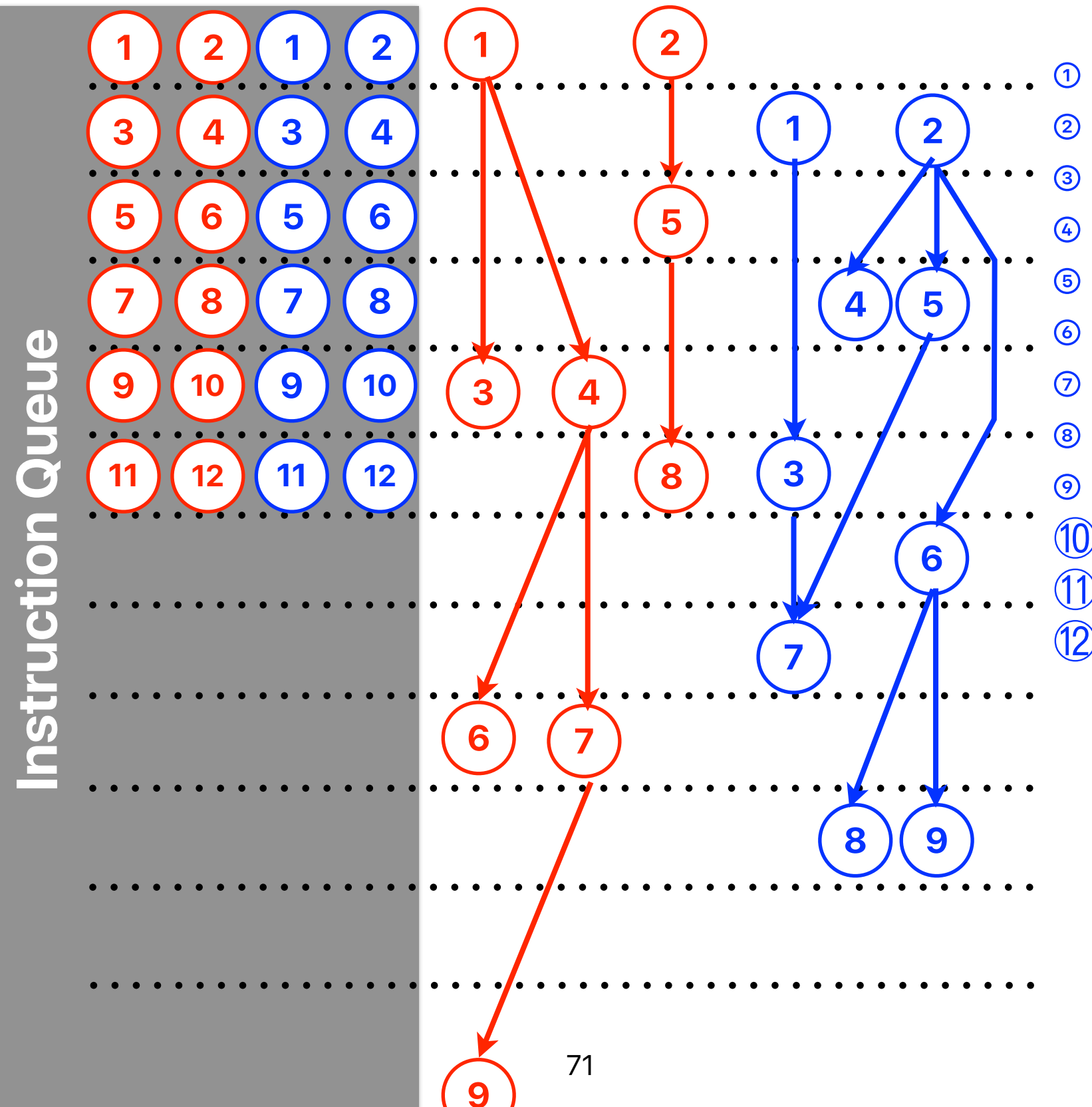


Simultaneous multithreading

- The processor can schedule instructions from different threads/processes/programs
- Fetch instructions from different threads/processes to fill the not utilized part of pipeline
 - Exploit “thread level parallelism” (TLP) to solve the problem of insufficient ILP in a single thread
 - You need to create an illusion of multiple processors for OSs

Simultaneous multithreading

① ld X10, 8(X10)
② addi X7, X7, 1
③ bne X10, X0, LOOP
④ ld X10, 8(X10)
⑤ addi X7, X7, 1
⑥ bne X10, X0, LOOP
⑦ ld X10, 8(X10)
⑧ addi X7, X7, 1
⑨ bne X10, X0, LOOP



① ld X1, 0(X10)
② addi X10, X10, 8
③ add X20, X20, X1
④ bne X10, X2, LOOP
⑤ ld X1, 0(X10)
⑥ addi X10, X10, 8
⑦ add X20, X20, X1
⑧ bne X10, X2, LOOP
⑨ ld X1, 0(X10)
⑩ addi X10, X10, 8
⑪ add X20, X20, X1
⑫ bne X10, X2, LOOP

Architectural support for simultaneous multithreading

- To create an illusion of a multi-core processor and allow the core to run instructions from multiple threads concurrently, how many of the following units in the processor must be duplicated/extended?
 - ① Program counter
 - ② Register mapping tables
 - ③ Physical registers
 - ④ ALUs
 - ⑤ Data cache
 - ⑥ Reorder buffer/Instruction Queue
 - A. 2
 - B. 3
 - C. 4
 - D. 5
 - E. 6

Architectural support for simultaneous multithreading



- To create an illusion of a multi-core processor and allow the core to run instructions from multiple threads concurrently, how many of the following units in the processor must be duplicated/extended?
 - ① Program counter
 - ② Register mapping tables
 - ③ Physical registers
 - ④ ALUs
 - ⑤ Data cache
 - ⑥ Reorder buffer/Instruction Queue
- A. 2
- B. 3
- C. 4
- D. 5
- E. 6

Architectural support for simultaneous multithreading

- To create an illusion of a multi-core processor and allow the core to run instructions from multiple threads concurrently, how many of the following units in the processor must be duplicated/extended?

- ① Program counter — **you need to have one for each context**
- ② Register mapping tables — **you need to have one for each context**
- ③ Physical registers — **you can share**
- ④ ALUs — **you can share**
- ⑤ Data cache — **you can share**
- ⑥ Reorder buffer/Instruction Queue
A. 2 — **you need to indicate which context the instruction is from**

B. 3

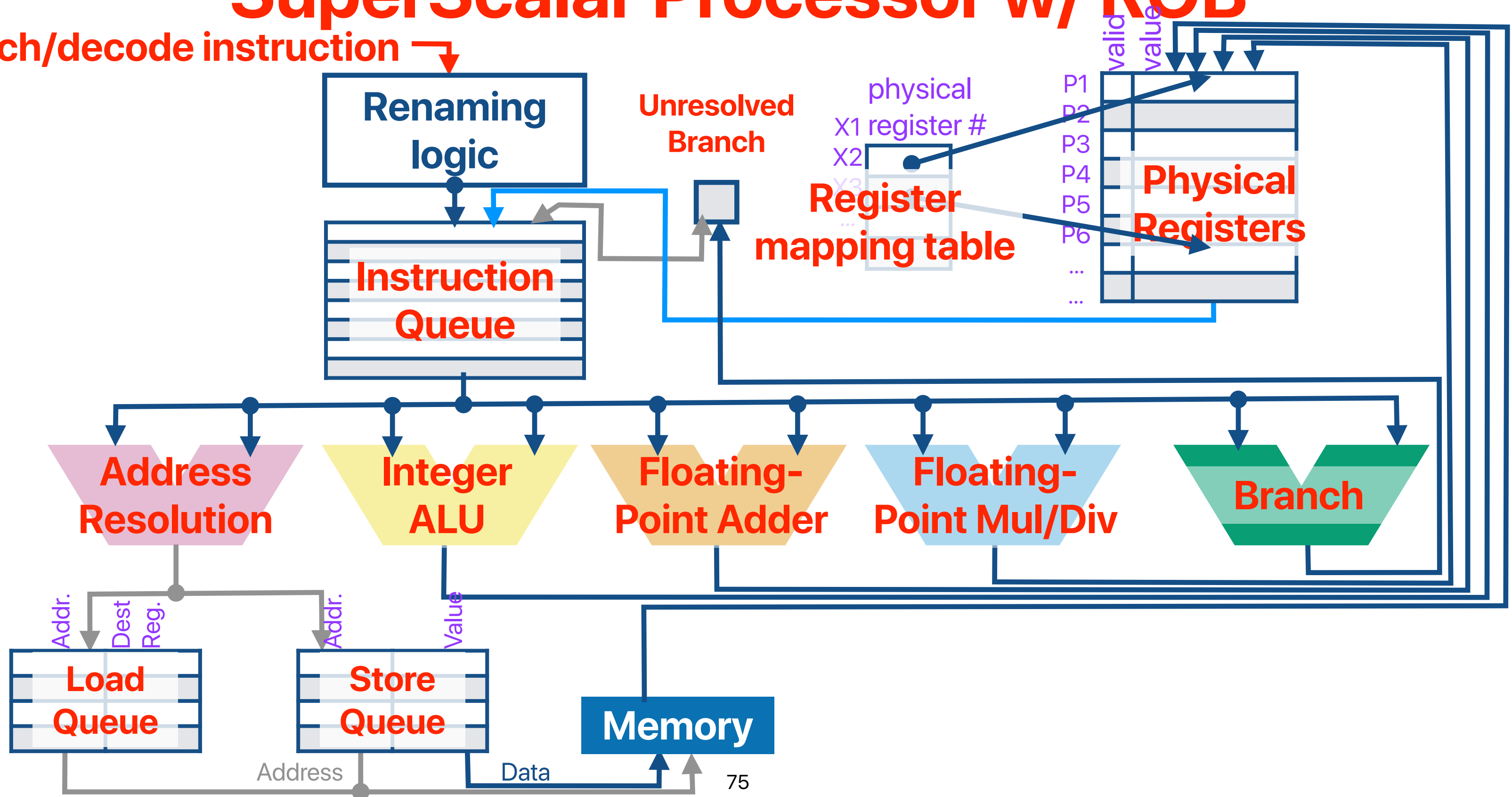
C. 4

D. 5

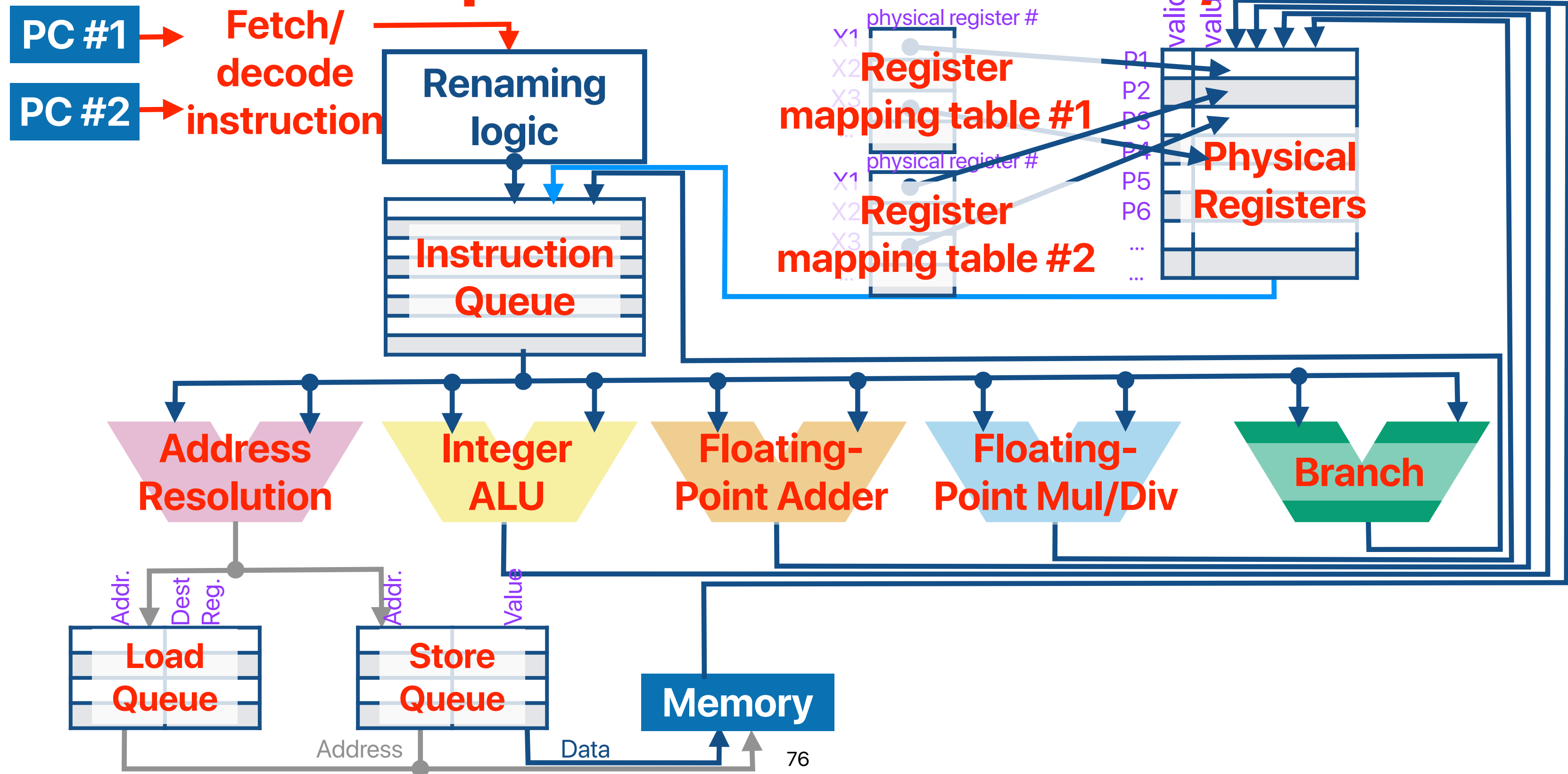
E. 6

SuperScalar Processor w/ ROB

Fetch/decode instruction →



SMT SuperScalar Processor w/ ROB



SMT

- How many of the following about SMT are correct?
 - ① SMT makes processors with deep pipelines more tolerable to mis-predicted branches
 - ② SMT can improve the throughput of a single-threaded application
 - ③ SMT processors can better utilize hardware during cache misses comparing with superscalar processors with the same issue width
 - ④ SMT processors can have higher cache miss rates comparing with superscalar processors with the same cache sizes when executing the same set of applications.
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

SMT



- How many of the following about SMT are correct?
 - ① SMT makes processors with deep pipelines more tolerable to mis-predicted branches
 - ② SMT can improve the throughput of a single-threaded application
 - ③ SMT processors can better utilize hardware during cache misses comparing with superscalar processors with the same issue width
 - ④ SMT processors can have higher cache miss rates comparing with superscalar processors with the same cache sizes when executing the same set of applications.
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

SMT

- How many of the following about SMT are correct?
 - ① SMT makes processors with deep pipelines more tolerable to mis-predicted branches **We can execute from other threads/contexts instead of the current one**
hurt, b/c you are sharing resource with other threads.
 - ② SMT can ~~improve~~ the throughput of a single-threaded application
 - ③ SMT processors can better utilize hardware during cache misses comparing with superscalar processors with the same issue width **We can execute from other threads/contexts instead of the current one**
 - ④ SMT processors can have higher cache miss rates comparing with superscalar processors with the same cache sizes when executing the same set of applications.

A. 0 **b/c we're sharing the cache**

B. 1

C. 2

D. 3

E. 4