

Programming on Parallel Architectures

Hung-Wei Tseng

AMD Zen 2 (RyZen 3000 Series)

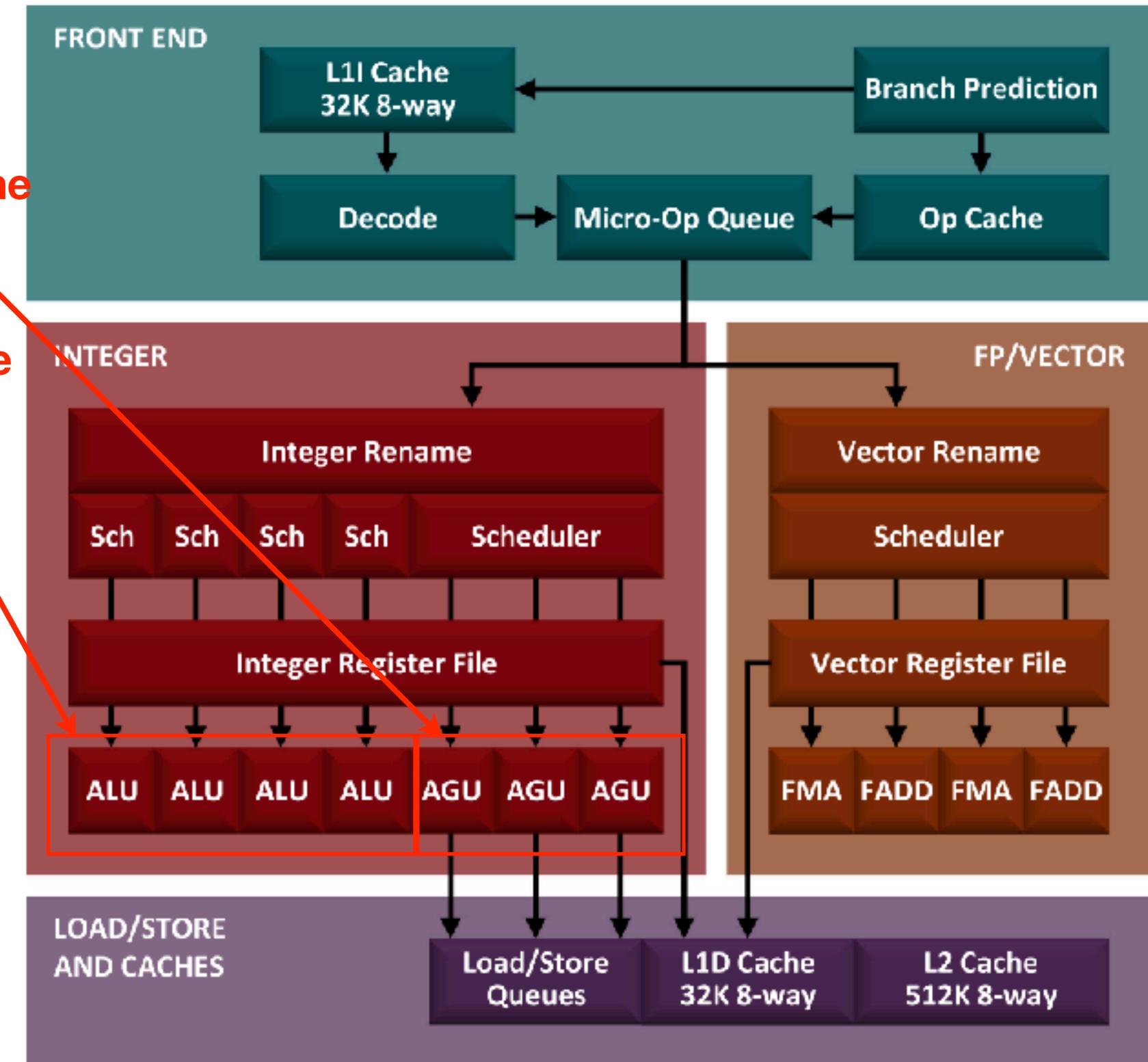
3-issue memory pipeline

4-issue integer pipeline

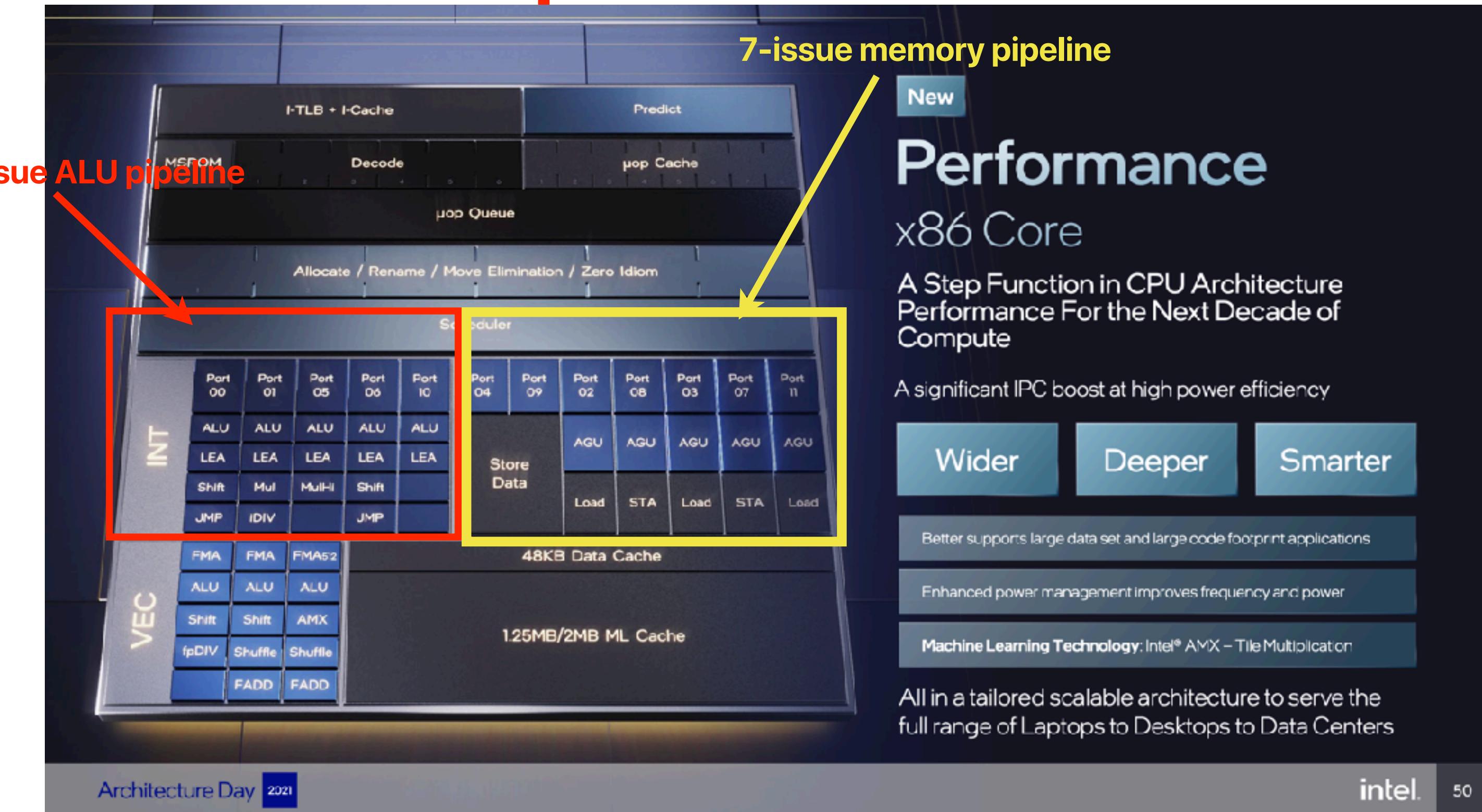
$$MinCPI = \frac{1}{7}$$

$$MinINTInst . CPI = \frac{1}{4}$$

$$MinMEMInst . CPI = \frac{1}{3}$$

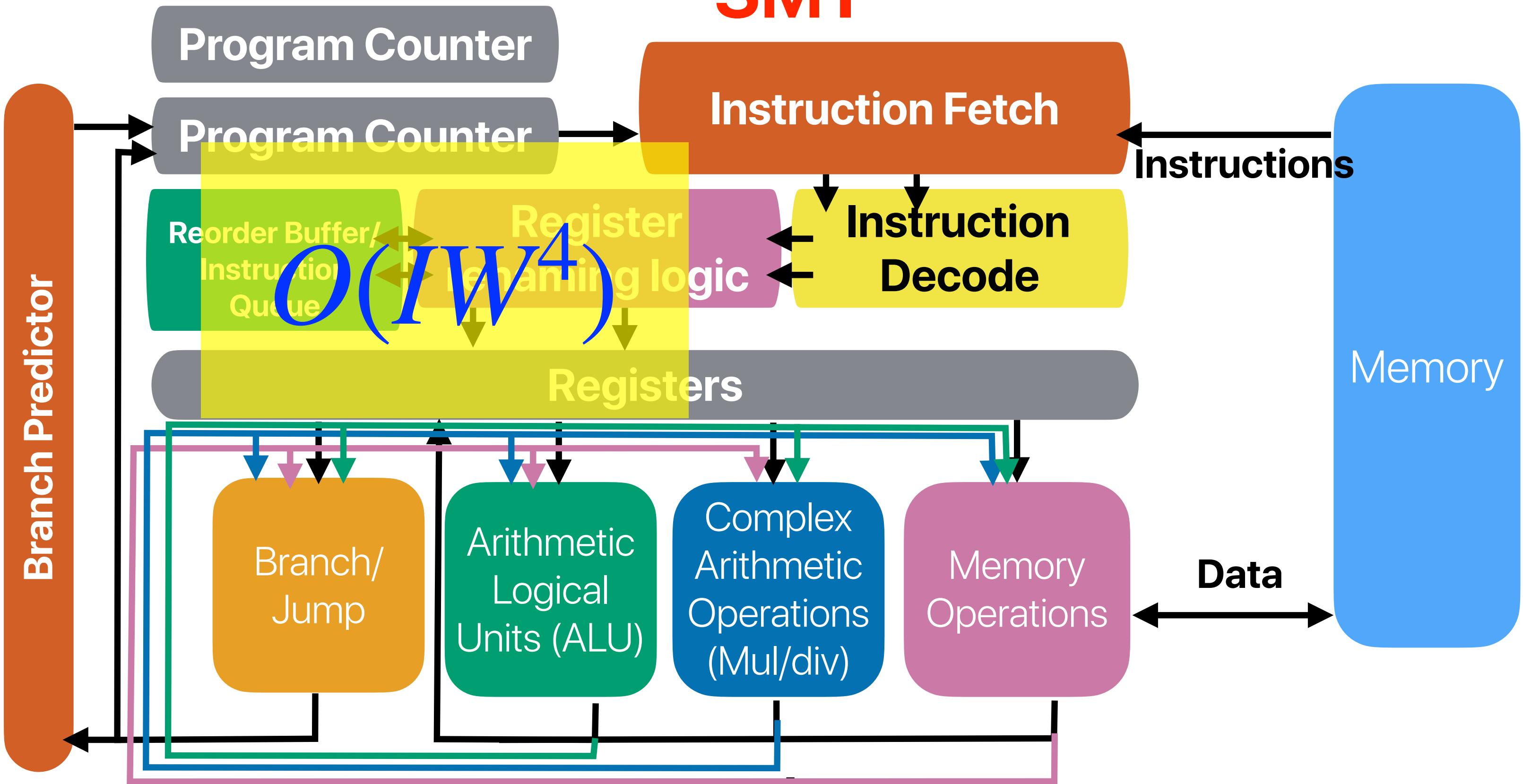


Recap: Intel Alder Lake



Architecture:	x86_64
CPU op-mode(s):	32-bit, 64-bit
Byte Order:	Little Endian
Address sizes:	48 bits physical, 48 bits virtual
CPU(s):	16
On-line CPU(s) list:	0-15
Thread(s) per core:	2
Core(s) per socket:	8
Socket(s):	1
NUMA node(s):	1
Vendor ID:	AuthenticAMD
CPU family:	25
Model:	80
Model name:	AMD Ryzen 7 5700G with Radeon Graphics
Stepping:	0

SMT



SMT

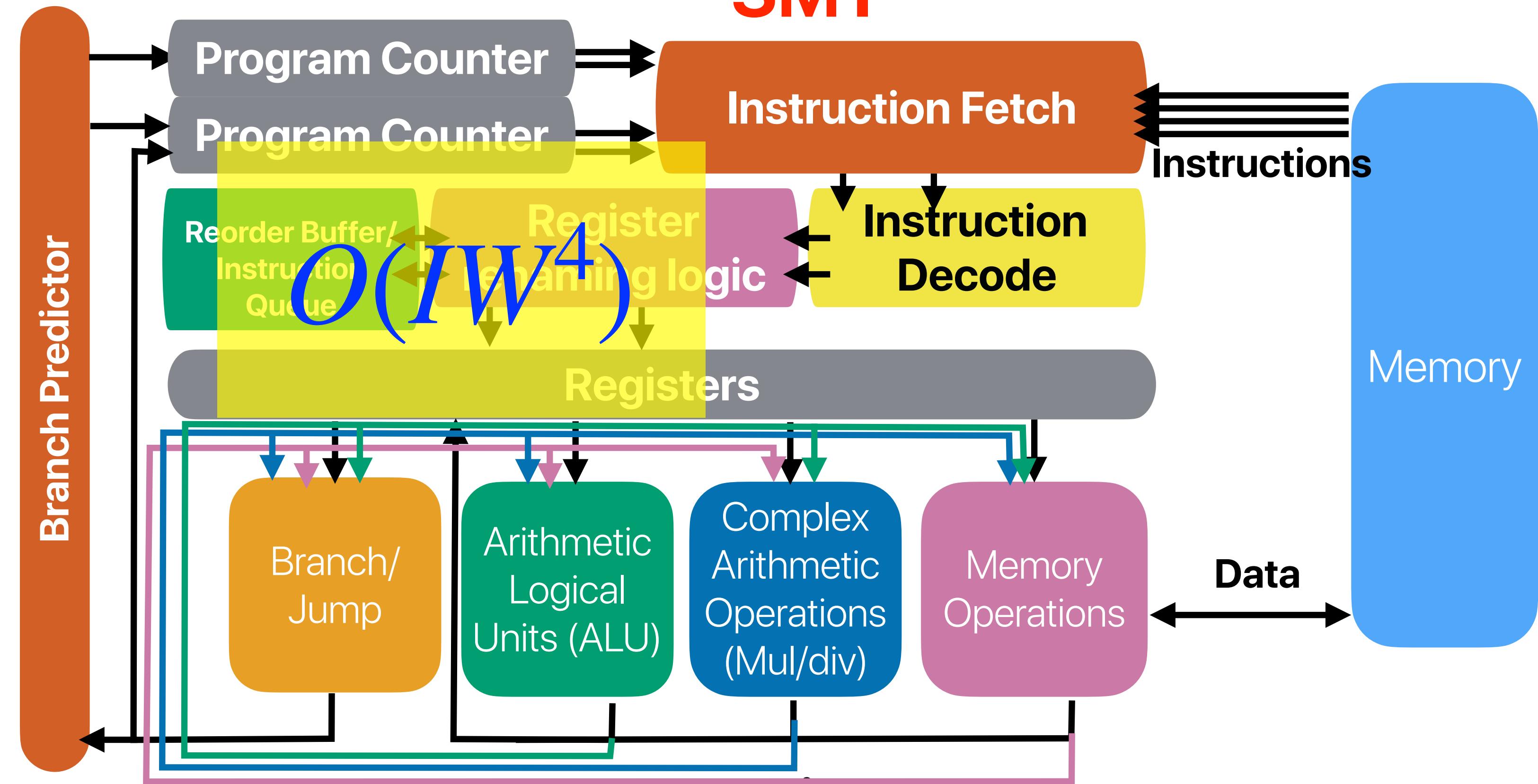
- Improve the throughput of execution
 - May increase the latency of a single thread
- Less branch penalty per thread
- Increase hardware utilization
- Simple hardware design: Only need to duplicate PC/Register Files
- Real Case:
 - Intel HyperThreading (supports up to two threads per core)
 - Intel Pentium 4, Intel Atom, Intel Core i7
 - AMD Ryzen (Zen microarchitecture)
 - If you see a processor with “threads” more than “cores”, that must be because of SMT!

Outline

- Multithreaded architectures
- Programming on Multithreaded Processors

Multithreaded Processors (cont.)

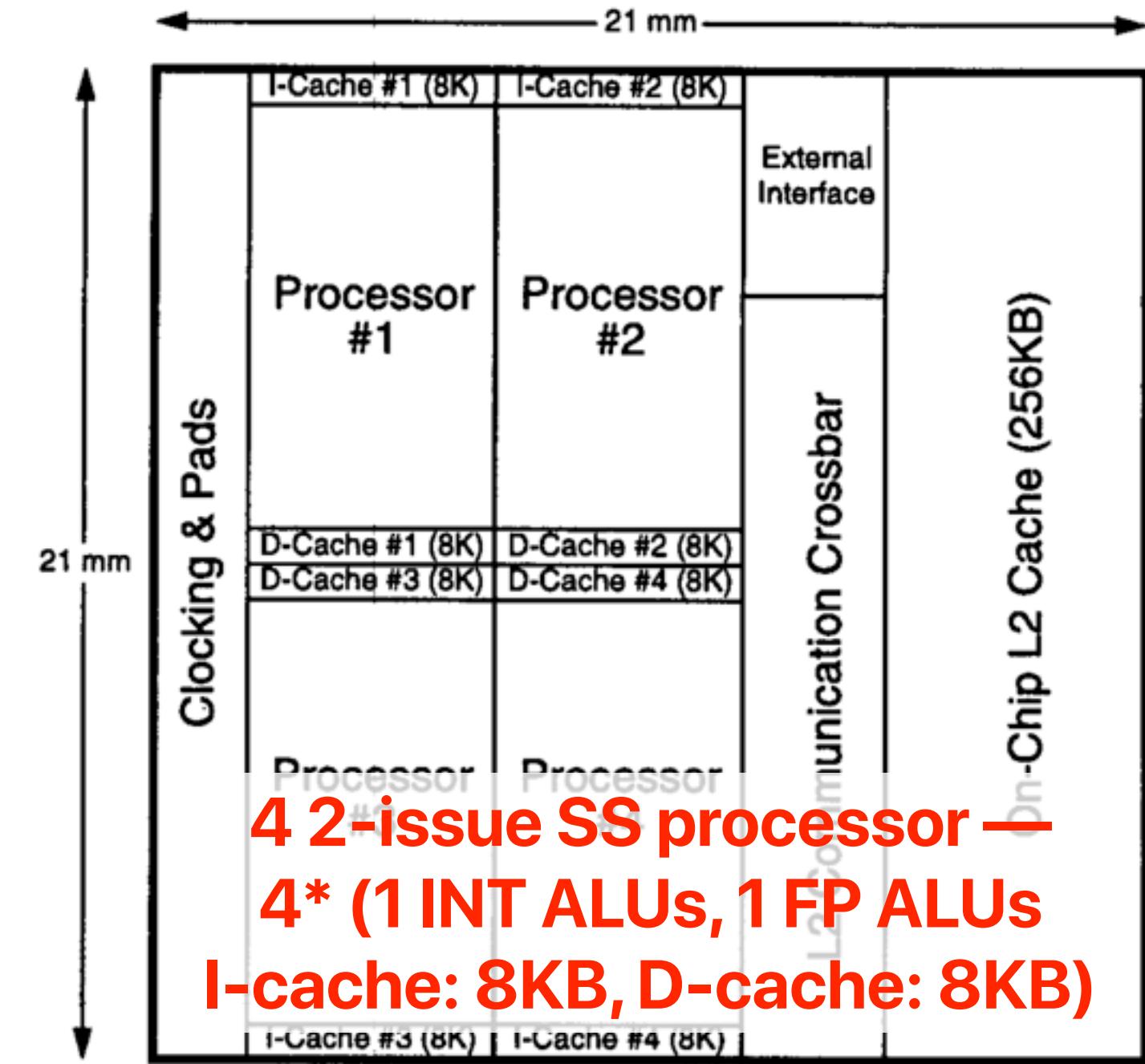
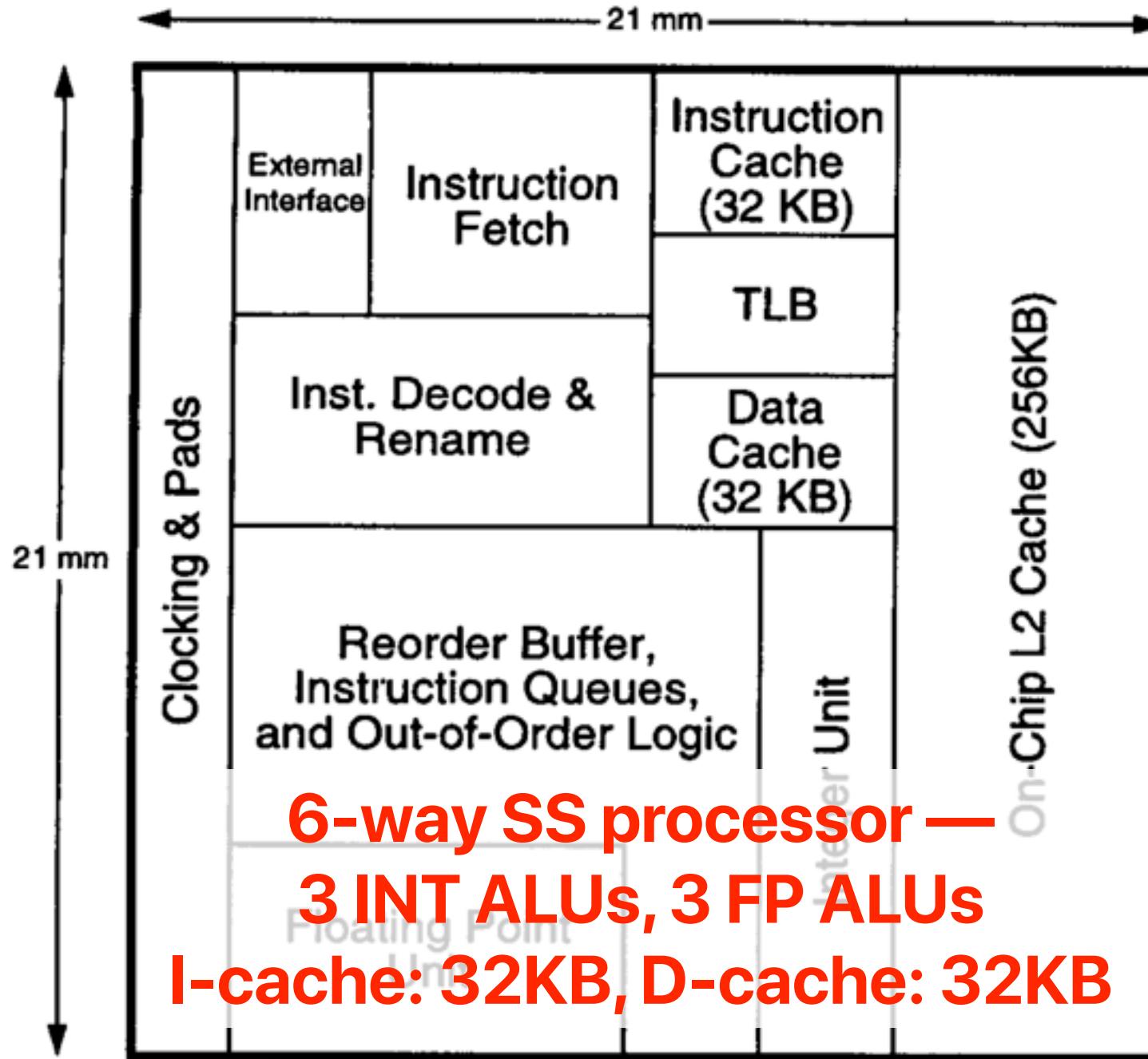
SMT



The case for a Single-Chip Multiprocessor

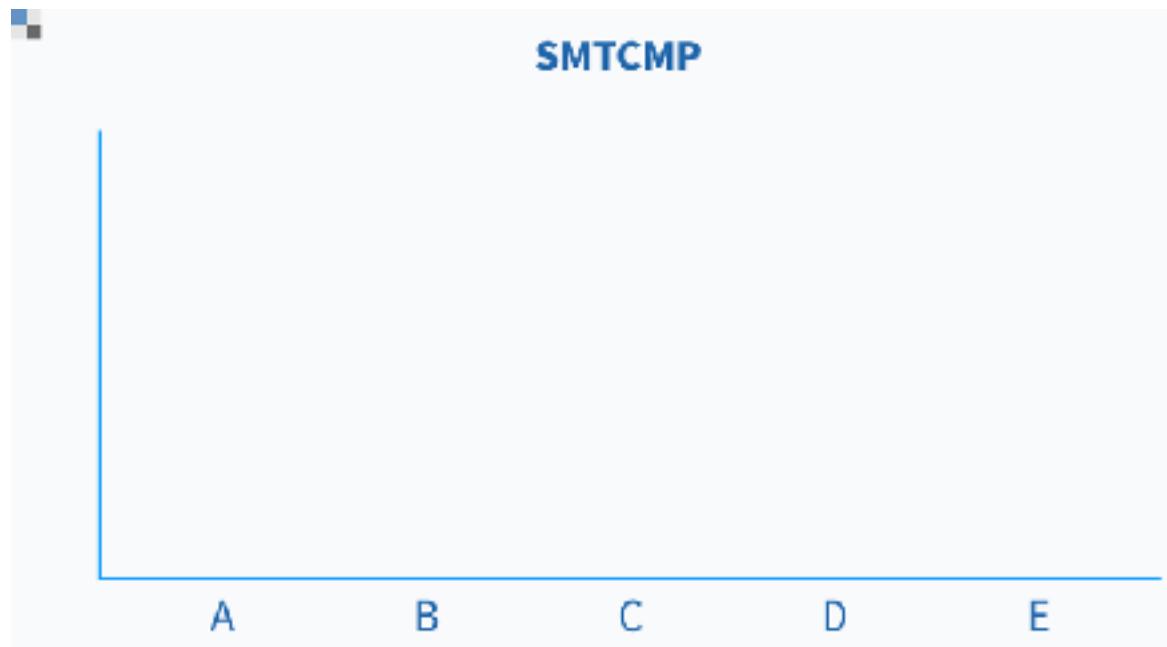
**Kunle Olukotun, Basem A. Nayfeh, Lance Hammond, Ken Wilson, and Kunyung
Chang
Stanford University**

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



SMT v.s. CMP

- An SMT processor is basically a SuperScalar processor with multiple instruction front-end. Assume within the same chip area, we can build an SMT processor supporting 4 threads, with 6-issue pipeline, 64KB cache or — a CMP with 4x 2-issue pipeline & 16KB cache in each core. Please identify how many of the following statements are/is correct when running programs on these processors.
 - ① If we are just running one program in the system, the program will perform better on an SMT processor
 - ② If we are running 4 applications simultaneously, the cache miss rates will be higher in the SMT processor
 - ③ If we are running 4 applications simultaneously, the branch mis-prediction will be higher in the SMT processor
 - ④ If we are running one program with 4 parallel threads, the cache miss rates will be higher in the SMT processor
 - ⑤ If we are running one program with 4 parallel threads simultaneously, the branch mis-prediction will be longer in the SMT processor
- A. 1
B. 2
C. 3
D. 4
E. 5

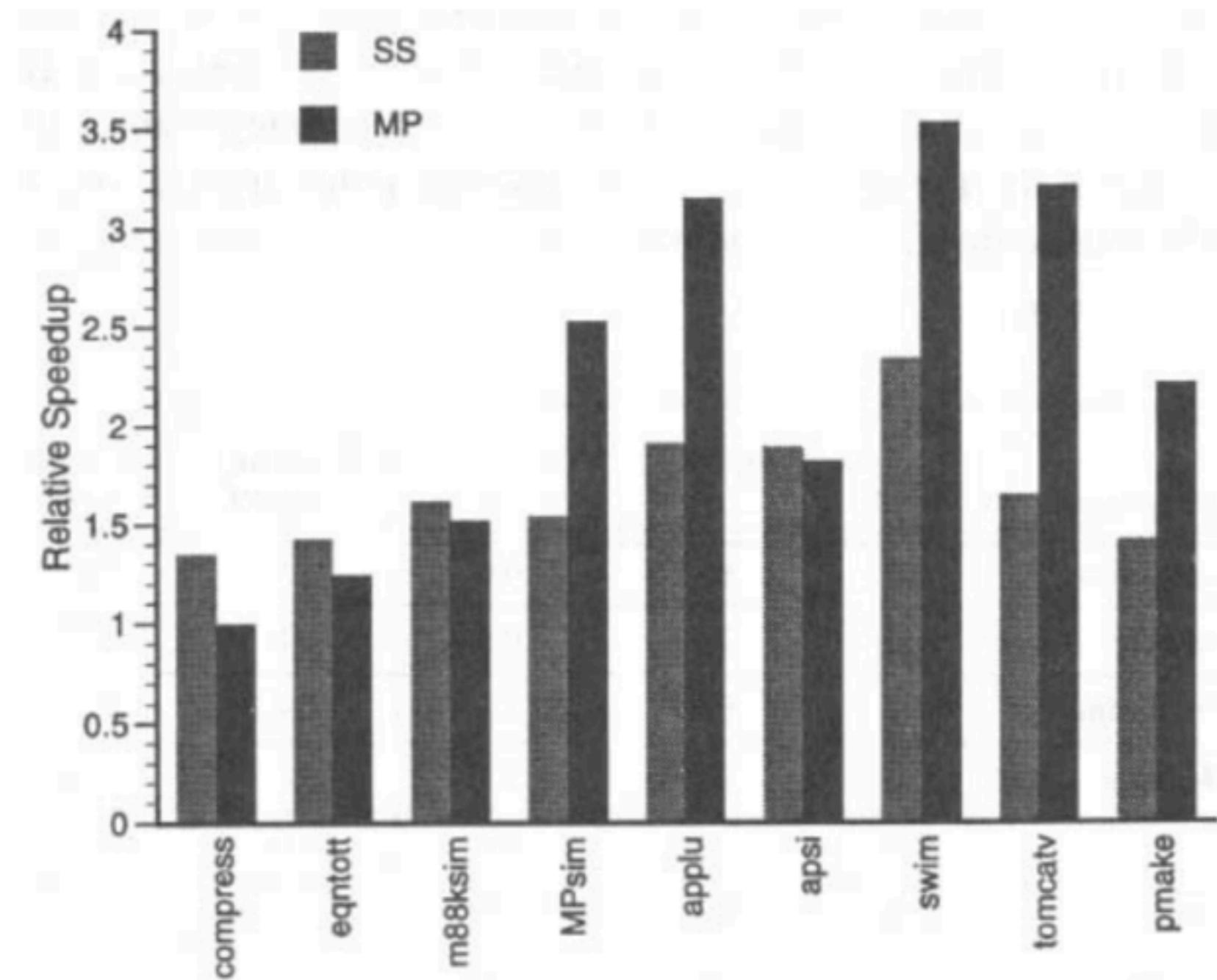


SMT v.s. CMP

- An SMT processor is basically a SuperScalar processor with multiple instruction front-end. Assume within the same chip area, we can build an SMT processor supporting 4 threads, with 6-issue pipeline, 64KB cache or — a CMP with 4x 2-issue pipeline & 16KB cache in each core. Please identify how many of the following statements are/is correct when running programs on these processors.
 - ① If we are just running one program in the system, the program will perform better on an SMT processor — **you have more resources for the program**
 - ② If we are running 4 applications simultaneously, the cache miss rates will be higher in the SMT processor
 - ③ If we are running 4 applications simultaneously, the branch mis-prediction will be higher in the SMT processor — **it depends!**
 - ④ If we are running one program with 4 parallel threads, the cache miss rates will be higher in the SMT processor — **it depends!**
 - ⑤ If we are running one program with 4 parallel threads simultaneously, the branch mis-prediction will be longer in the SMT processor — **it depends!**

A. 1
B. 2
C. 3 **The only thing we know for sure — if we don't parallel the program, it won't get any faster on CMP**
D. 4
E. 5

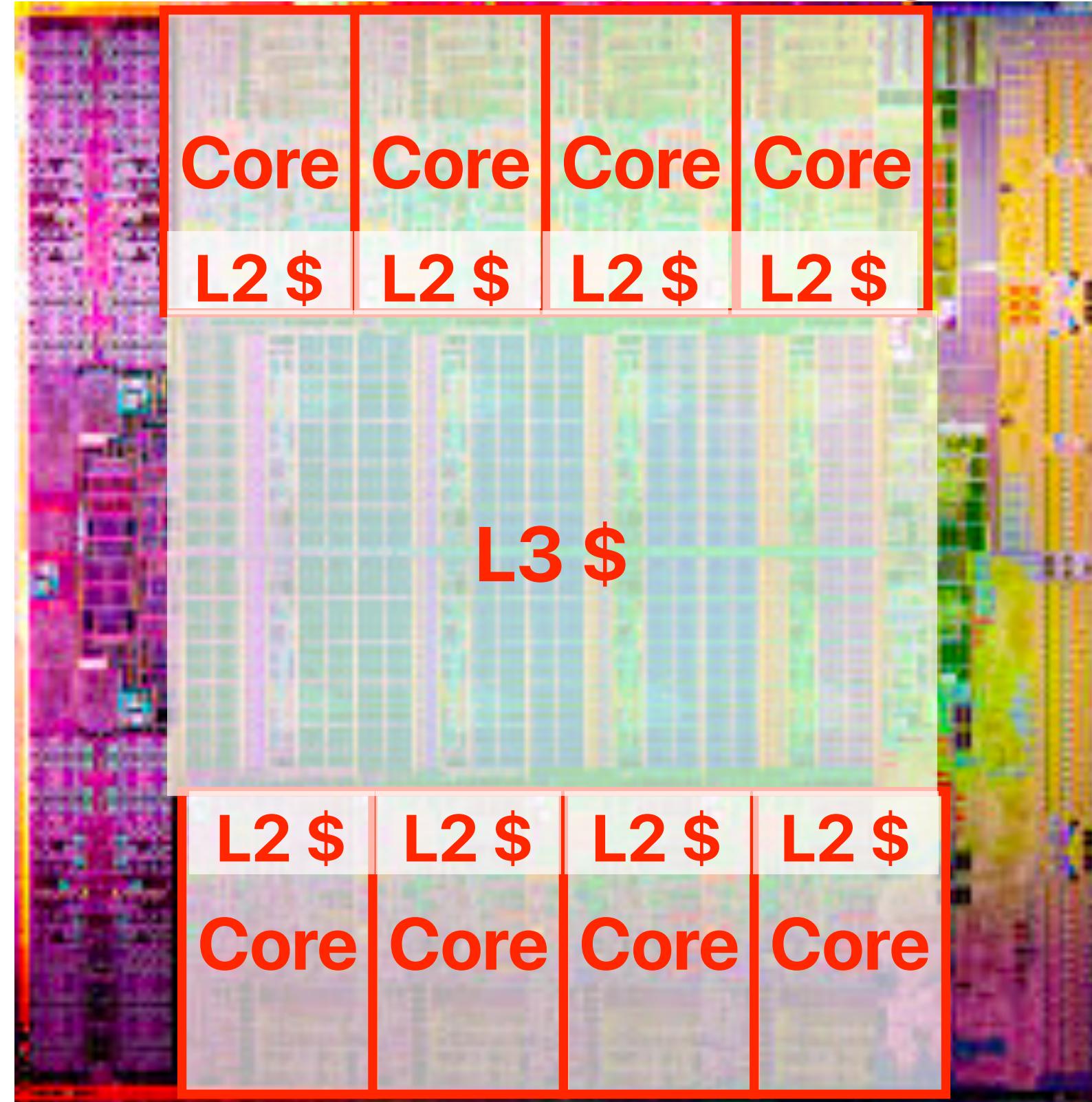
6-way SuperScalar v.s. quad-core CMP

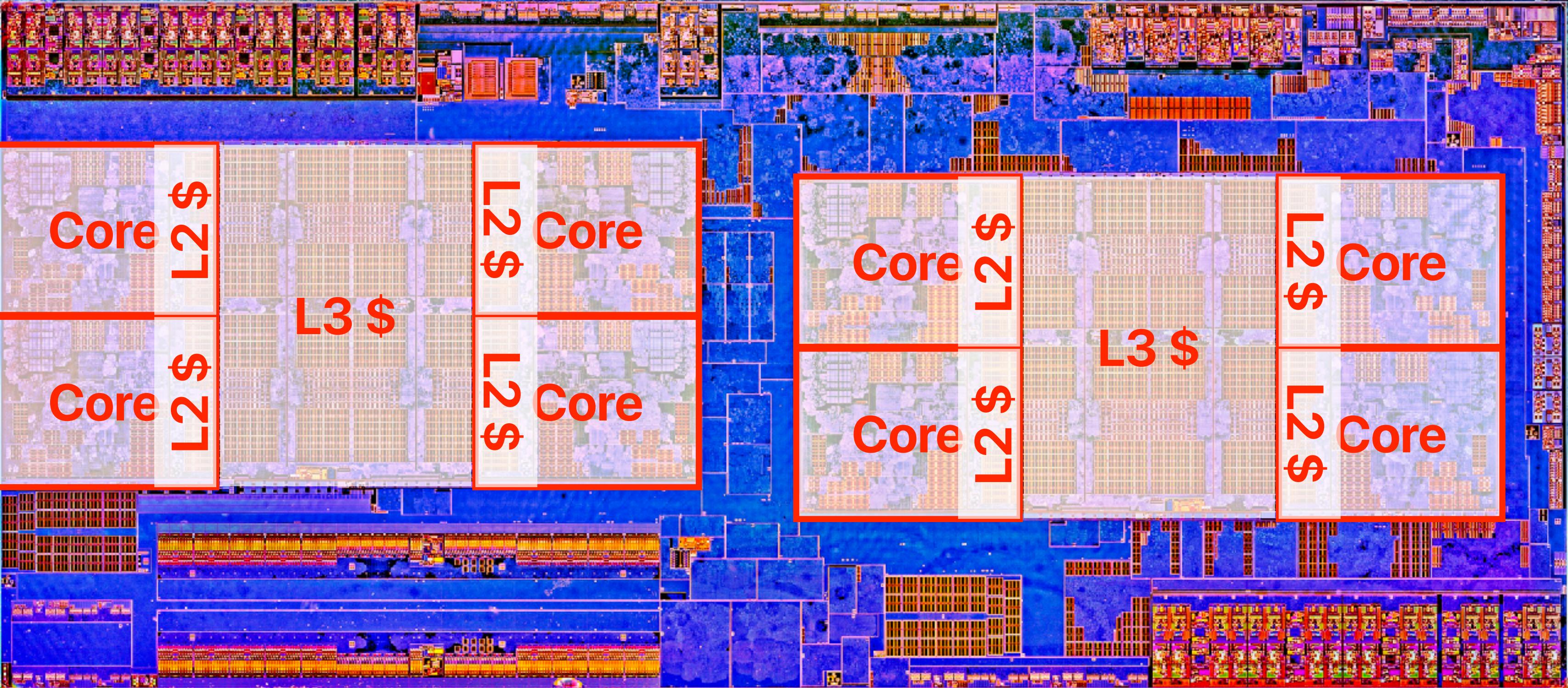


The applications are parallelized in different ways to run on the MP microarchitecture. Compress is run unmodified on both the SS and MP microarchitectures; using only one processor of the MP architecture. Eqntott is parallelized manually by modifying a single bit vector comparison routine that is responsible for 90% of the execution time of the application [16]. The CPU simulator m88ksim is also parallelized manually into three threads using the SUIF compiler runtime system. Each of the three threads is allowed to be in a different phase of simulating a different instruction at the same time. This style of parallelization is very similar to the overlap of instruction execution that occurs in hardware pipelining. The

Figure 6. Performance comparison of SS and MP.

Intel Sandy Bridge

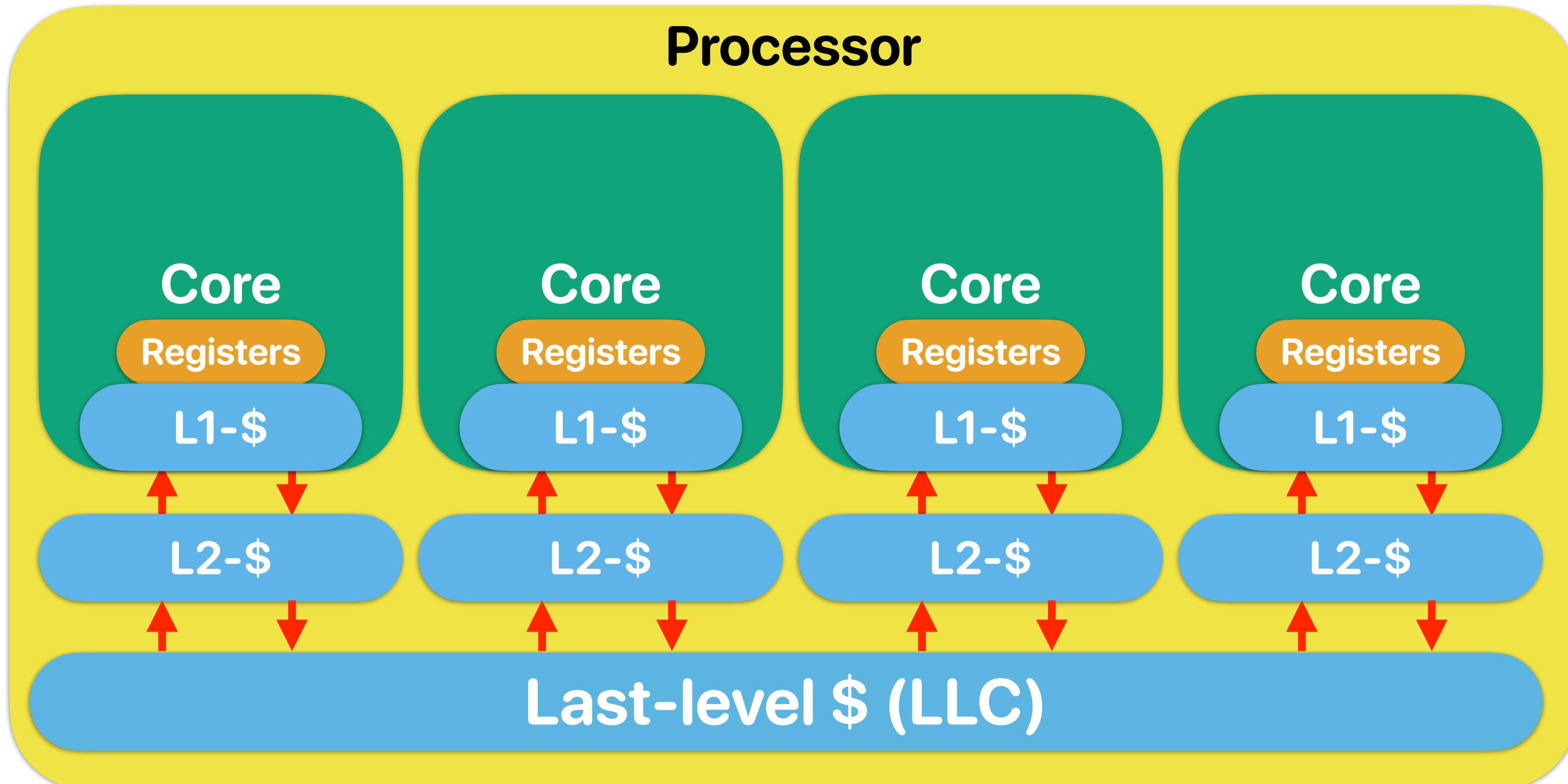




AMD

RYZEN

Concept of CMP

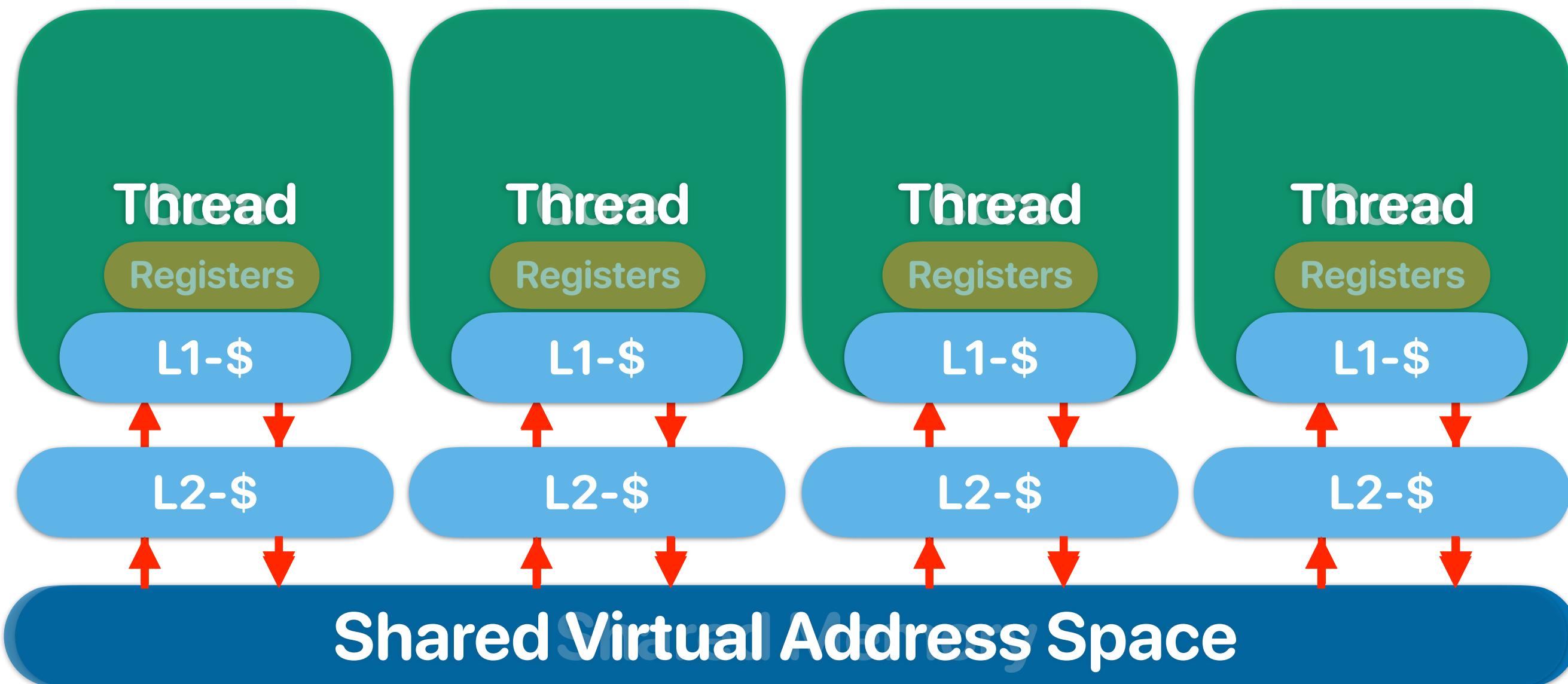


Architectural Support for Parallel Programming

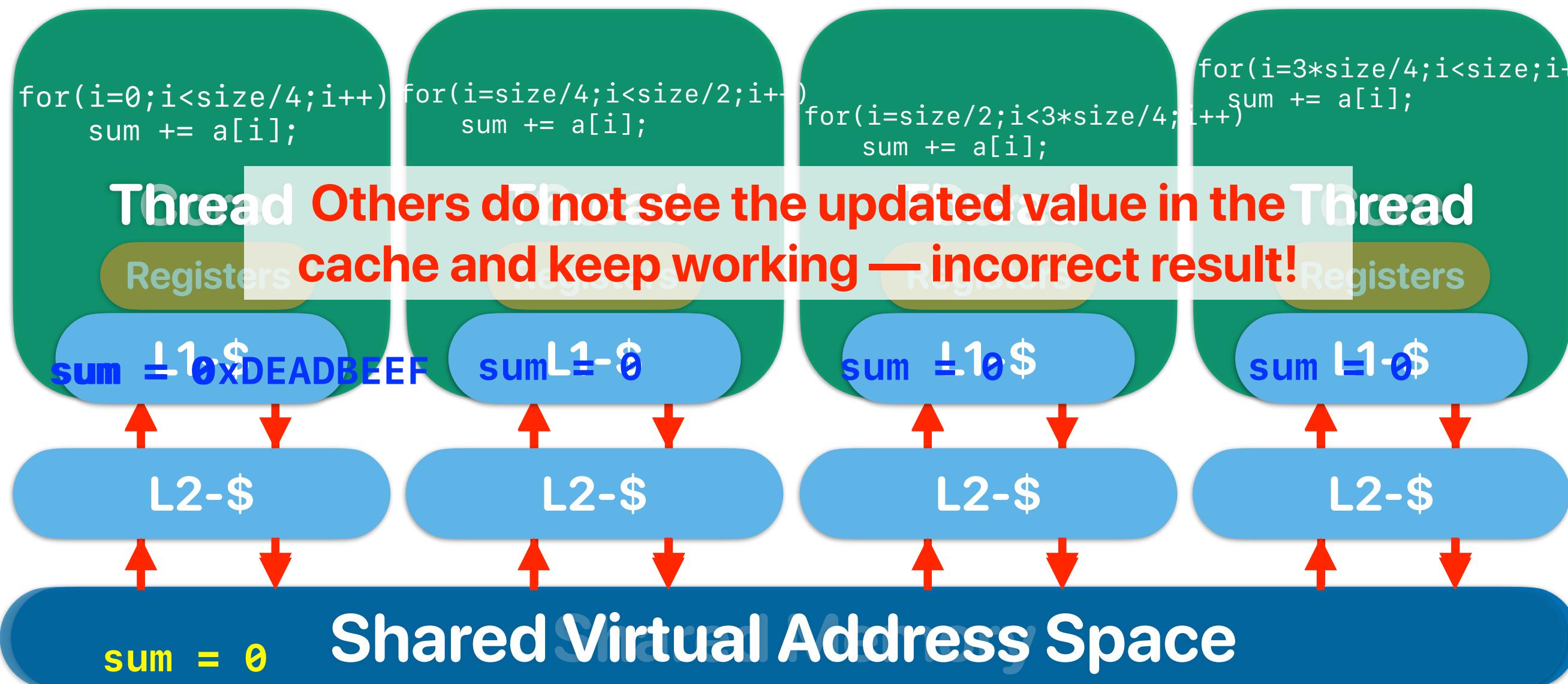
Parallel programming

- 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

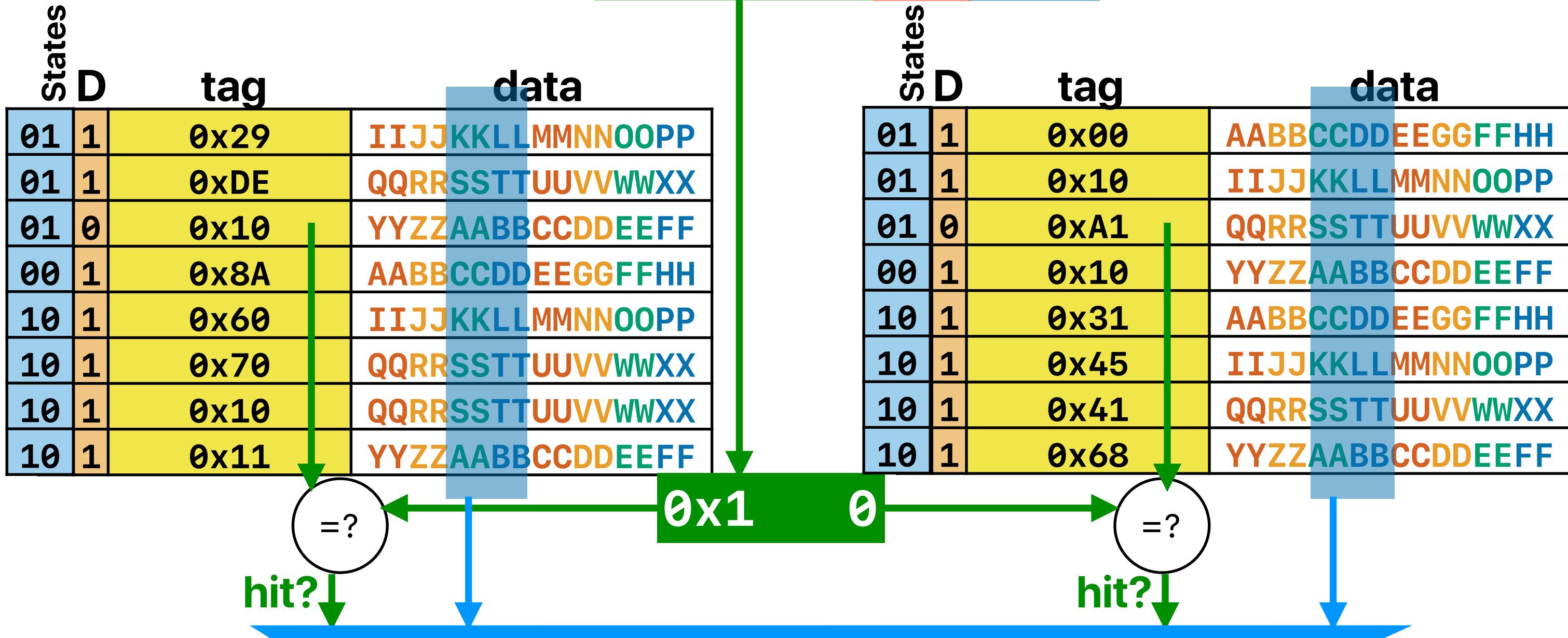
Coherent way-associative cache

memory address:

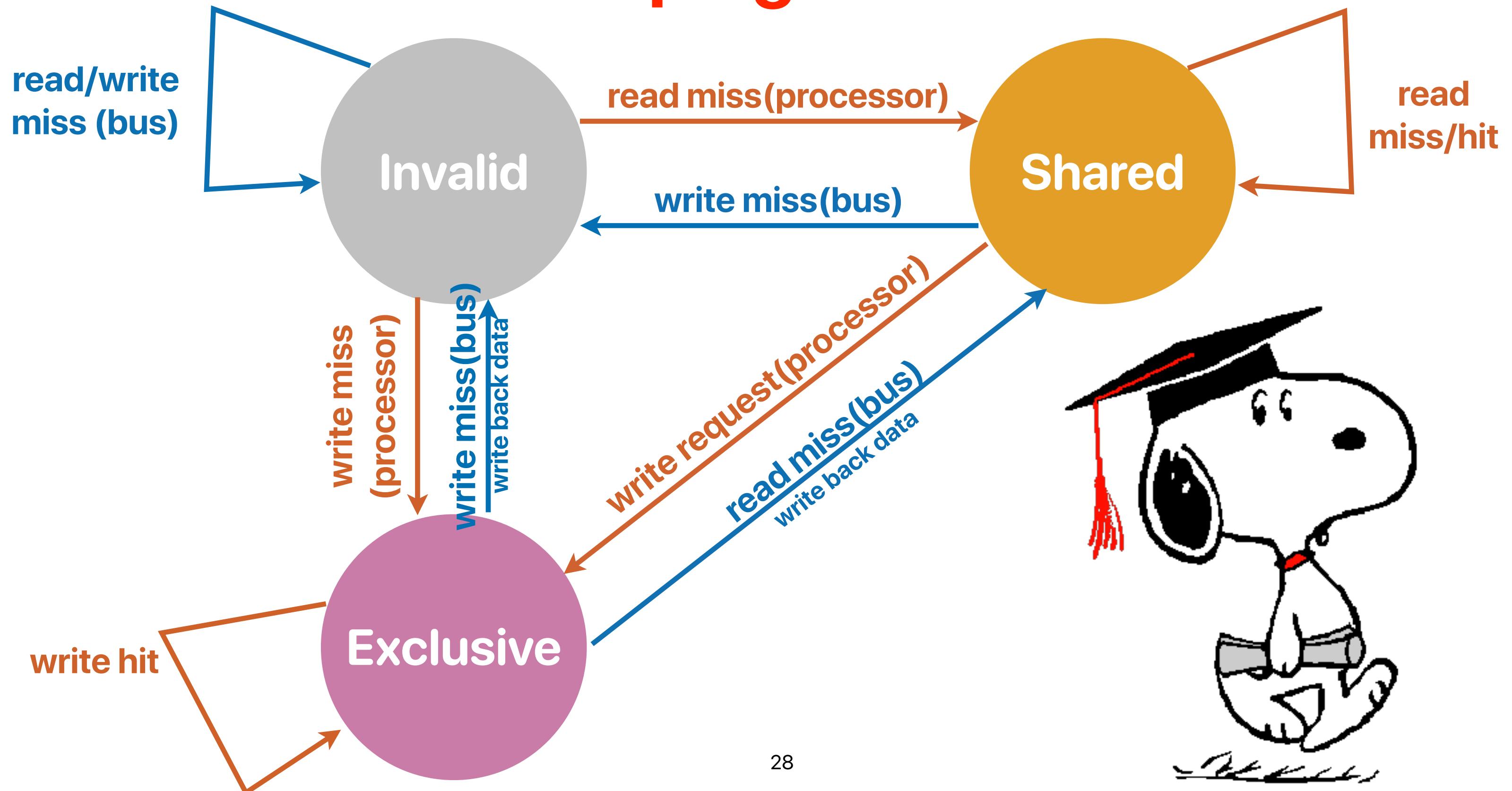
$0x0$ 8 tag 2 set 4 block
index offset

memory address:

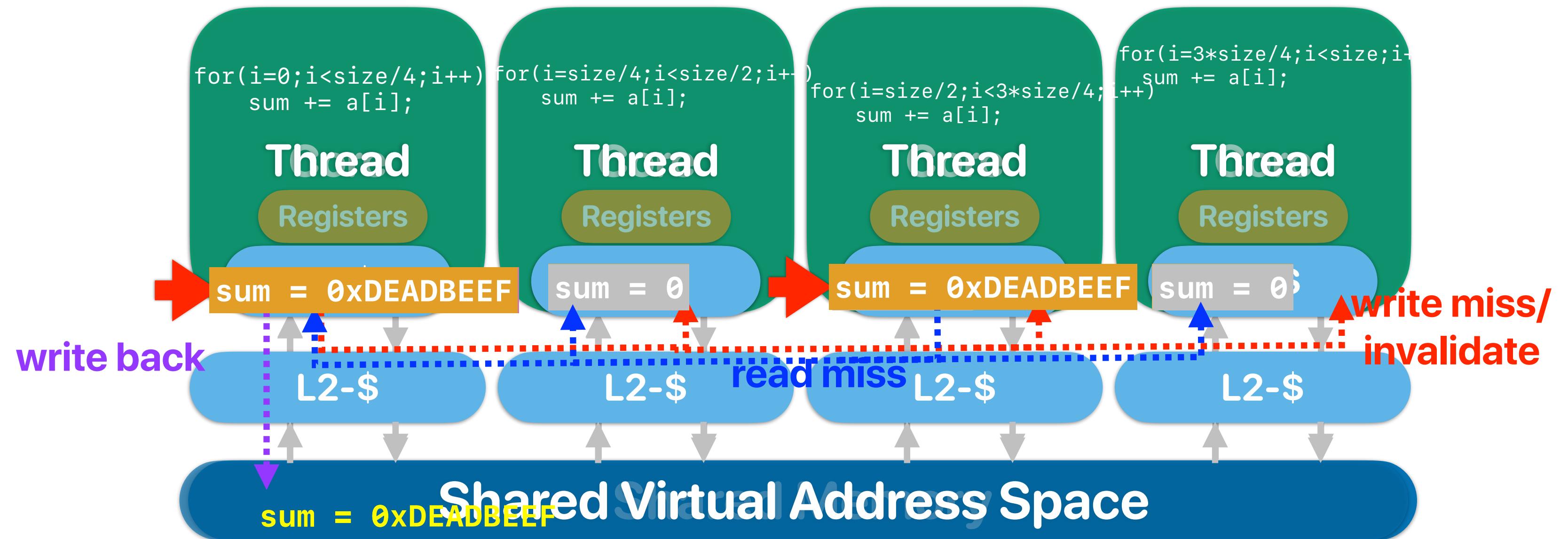
0b0000100000100100



Snooping Protocol



What happens when we write in coherent caches?



Observer

thread 1	thread 2
<pre>int loop; int main() { pthread_t thread; loop = 1; pthread_create(&thread, NULL, modifyloop, NULL); while(loop == 1) { continue; } pthread_join(thread, NULL); fprintf(stderr,"User input: %d\n", loop); return 0; }</pre>	<pre>void* modifyloop(void *x) { sleep(1); printf("Please input a number:\n"); scanf("%d",&loop); return NULL; }</pre>

Observer

prevents the compiler from putting the variable "loop" in the "register"

thread 1

```
volatile int loop;  
  
int main()  
{  
    pthread_t thread;  
    loop = 1;  
  
    pthread_create(&thread, NULL, modifyloop,  
NULL);  
    while(loop == 1)  
    {  
        continue;  
    }  
    pthread_join(thread, NULL);  
    fprintf(stderr, "User input: %d\n", loop);  
    return 0;  
}
```

thread 2

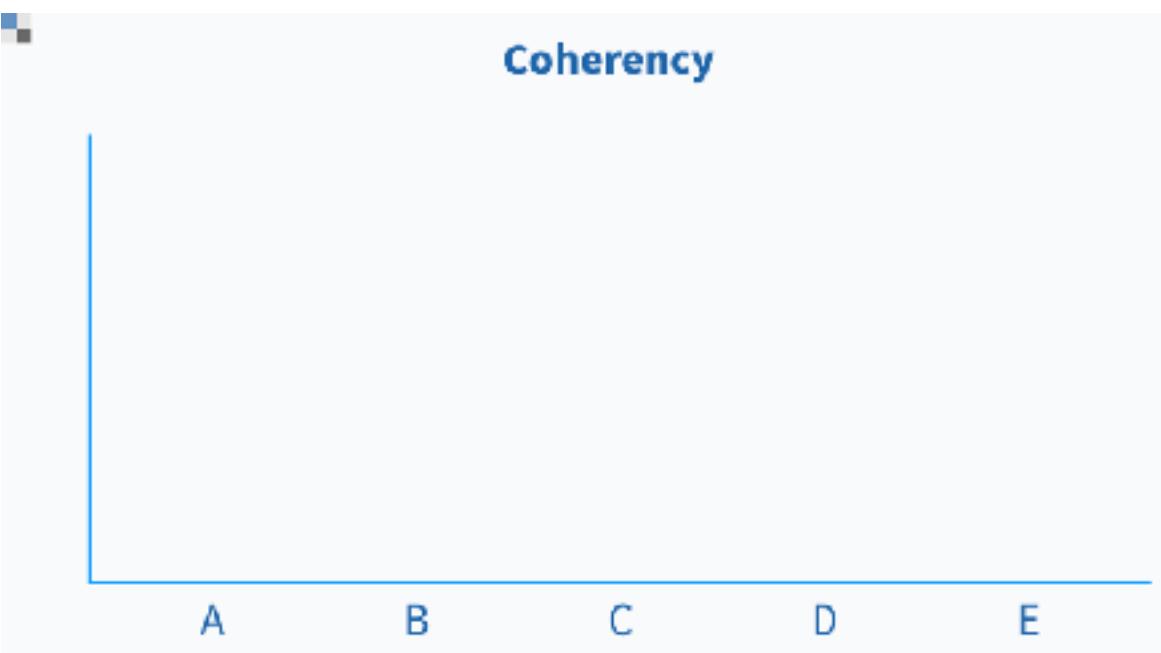
```
void* modifyloop(void *x)  
{  
    sleep(1);  
    printf("Please input a number:\n");  
    scanf("%d", &loop);  
    return NULL;  
}
```

Cache coherency

- Assuming that we are running the following code on a CMP with a cache coherency protocol, how many of the following outputs are possible? (a is initialized to 0 as assume we will output more than 10 numbers)

thread 1	thread 2
while(1) printf("%d ", a);	while(1) a++;

- ① 0123456789
- ② 1259368101213
- ③ 1111111164100
- ④ 111111111100
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4



Cache coherency

- Assuming that we are running the following code on a CMP with a cache coherency protocol, how many of the following outputs are possible? (a is initialized to 0 as assume we will output more than 10 numbers)

thread 1	thread 2
while(1) printf("%d ", a);	while(1) a++;

- ① 0123456789
- ② 1259368101213
- ③ 1111111164100
- ④ 111111111100
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

Announcements

- Last Reading Quiz due next Monday
 - We drop two of your lowest ones
- Assignment #4 due 12/2
- If you submit iEVAL and submit the screenshot through eLearn, it counts as a “full-credit” notebook assignment
 - We drop two notebook assignments with this one included
 - In other words, if you submit iEVAL and the screenshot, you got two lowest assignments dropped.

Computer Science & Engineering

203



づづく