

Multithreaded Architectures and Programming on Multithreaded Architectures

Hung-Wei Tseng

Summary: Characteristics of modern processor architectures

- Multiple-issue pipelines with multiple functional units available
 - Multiple ALUs
 - Multiple Load/store units
 - Dynamic OoO scheduling to reorder instructions whenever possible
- Cache — very high hit rate if your code has good locality
 - Very matured data/instruction prefetcher
- Branch predictors — very high accuracy if your code is predictable
 - Perceptron
 - Variable history predictors

Recap: popcount

- Which of the following implementations will perform the best on modern processors?

A

```
inline int __popcount(uint64_t x) {
    int c=0;
    while(x) {
        c += x & 1;
        x = x >> 1;
    }
    return c;
}
```

B

```
inline int __popcount(uint64_t x) {
    int c = 0;
    int table[16] = {0, 1, 1, 2, 1,
                    2, 2, 3, 1, 2, 2, 3, 2, 3, 3, 4};
    while(x) {
        c += table[(x & 0xF)];
        x = x >> 4;
    }
    return c;
}
```

C

```
inline int __popcount(uint64_t x) {
    int c = 0;
    while(x) {
        c += x & 1;
        x = x >> 1;
    }
    return c;
}
```

D

```
inline int __popcount(uint64_t x) {
    int c = 0;
    int table[16] = {0, 1, 1, 2, 1,
                    2, 2, 3, 1, 2, 2, 3, 2, 3, 3, 4};
    for (uint64_t i = 0; i < 16; i++)
    {
        c += table[(x & 0xF)];
        x = x >> 4;
    }
    return c;
}
```

E

```
#include <simmmintrin.h>
inline int __popcount(uint64_t x) {
    int c = _mm_popcnt_u64(x);
    return c;
}
```

B

```
inline int __popcount(uint64_t x) {
    int c = 0;
    while(x) {
        c += x & 1;
        x = x >> 1;
    }
    return c;
}
```

C

```
inline int __popcount(uint64_t x) {
    int c = 0;
    while(x) {
        c += x & 1;
        x = x >> 1;
    }
    return c;
}
```

D

```
inline int __popcount(uint64_t x) {
    int c = 0;
    int table[16] = {0, 1, 1, 2, 1,
                    2, 2, 3, 1, 2, 2, 3, 2, 3, 3, 4};
    for (uint64_t i = 0; i < 16; i++)
    {
        c += table[(x & 0xF)];
        x = x >> 4;
    }
    return c;
}
```

E

```
#include <simmmintrin.h>
inline int __popcount(uint64_t x) {
    int c = _mm_popcnt_u64(x);
    return c;
}
```

W

```
inline int __popcount(uint64_t x) {
    int c = 0;
    for (uint64_t i = 0; i < 16; i++)
    {
        switch((x & 0xF))
        {
            case 1: c+=1; break;
            case 2: c+=1; break;
            case 3: c+=2; break;
            case 4: c+=1; break;
            case 5: c+=2; break;
            case 6: c+=2; break;
            case 7: c+=3; break;
            case 8: c+=1; break;
            case 9: c+=2; break;
            case 10: c+=2; break;
            case 11: c+=3; break;
            case 12: c+=2; break;
            case 13: c+=3; break;
            case 14: c+=3; break;
            case 15: c+=4; break;
            default: break;
        }
        x = x >> 4;
    }
    return c;
}
```

L

```
#include <simmmintrin.h>
inline int __popcount(uint64_t x) {
    int c = _mm_popcnt_u64(x);
    return c;
}
```

Recap: Tips of programming on modern processors

- Minimize the critical path operations
 - Don't forget about optimizing cache/memory locality first!
 - Memory latencies are still way longer than any arithmetic instruction
 - Can we use arrays/hash tables instead of lists?
 - Branch can be expensive as pipeline get deeper
 - Sorting
 - Loop unrolling
 - Still need to carefully avoid long latency operations (e.g., mod)
- Since processors have multiple functional units — code must be able to exploit instruction-level parallelism
 - Hide as many instructions as possible under the "critical path"
 - Try to use as many different functional units simultaneously as possible
- Modern processors also have accelerated instructions
- Compiler can do fairly go optimizations, but with limitations

Recap: Summary of popcounts

	ET	IC	IPC/ILP	# of branches	Branch mis-prediction rate
A	22.21	332 Trillions	2.88	65 Trillions	1.13%
B	12.29	287 Trillions	4.52	17 Trillions	0.04%
C	5.01	102 Trillions	3.95	17 Trillions	0.04%
D	3.73	80 Trillions	4.13	1 Trillions	~0%
E	54.4	173 Trillions	0.61	44 Trillions	18.6%
SSE4.2	1.57	22 Trillions	2.7	1 Trillions	~0%

Best performing one at 2.7

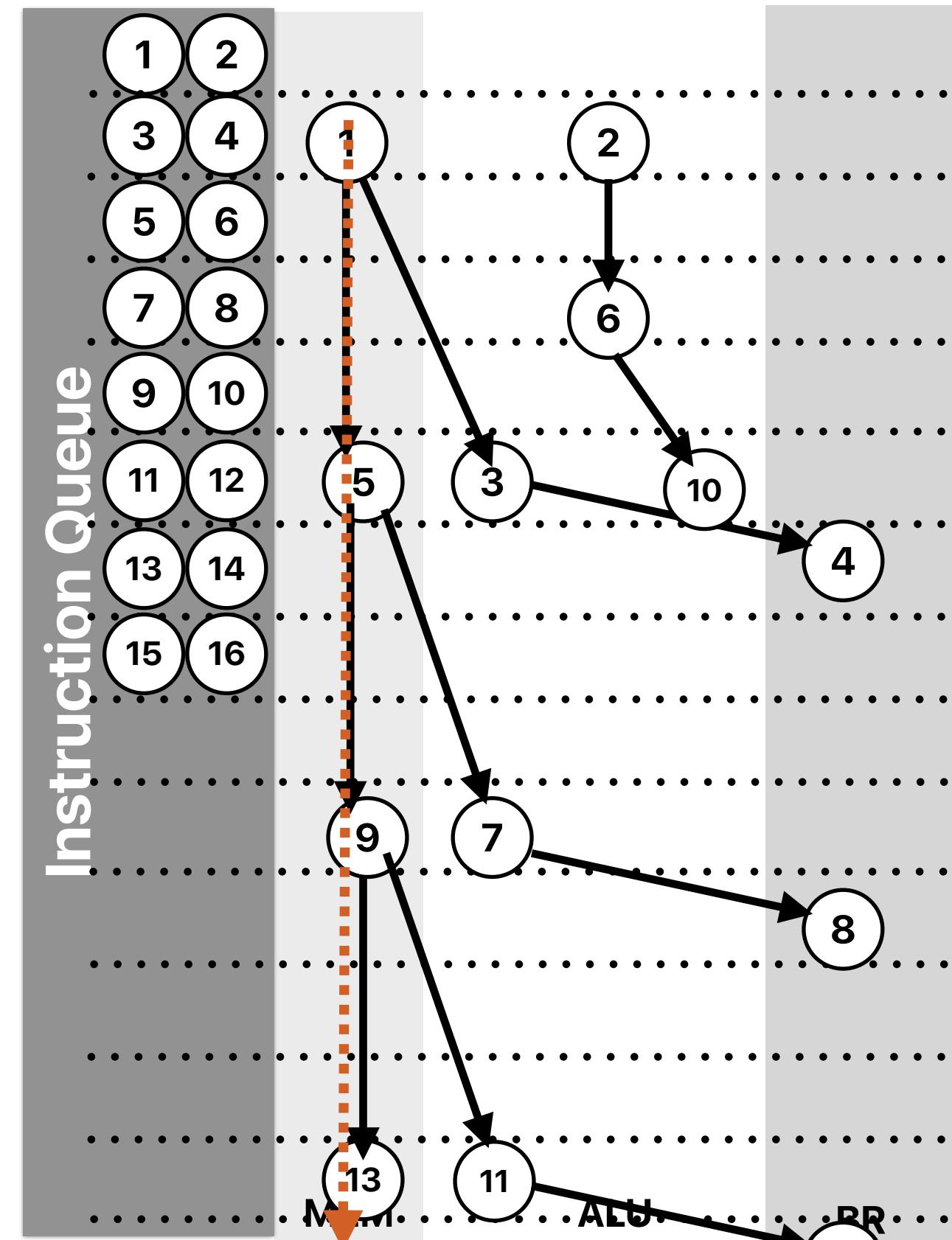
What if we have “unlimited” fetch/issue width — “linked list”

Doesn't help that much!

- It's important that the programmer should write code that can exploit “ILP”
- But — there're always cases we cannot do further in ILP

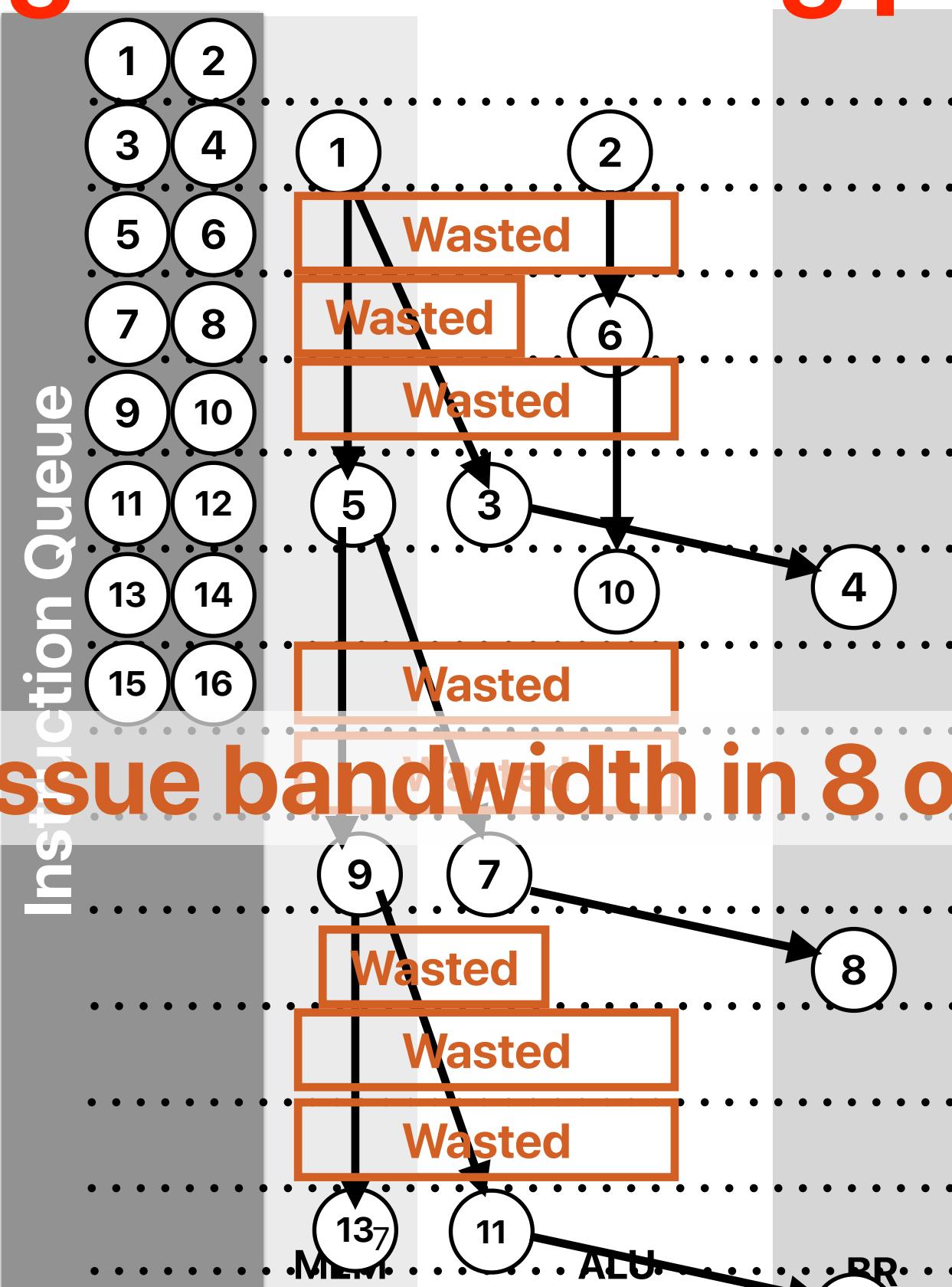
```
do {  
    number_of_nodes++;  
    current = current->next;  
} while ( current != NULL );
```

```
① .L3:    movq    8(%rdi), %rdi  
②          addl    $1, %eax  
③          testq   %rdi, %rdi  
④          jne     .L3
```



2-issue register renaming pipeline

```
① movq    8(%rdi), %rdi  
② addl    $1, %eax  
③ testq   %rdi, %rdi  
④ jne     .L3  
⑤ movq    8(%rdi), %rdi  
⑥ addl    $1, %eax  
⑦ testq   %rdi, %rdi  
⑧ jne     .L3  
⑨ movq    8(%rdi), %rdi  
⑩ addl    $1, %eax  
⑪ testq   %rdi, %rdi  
⑫ jne     .L3
```

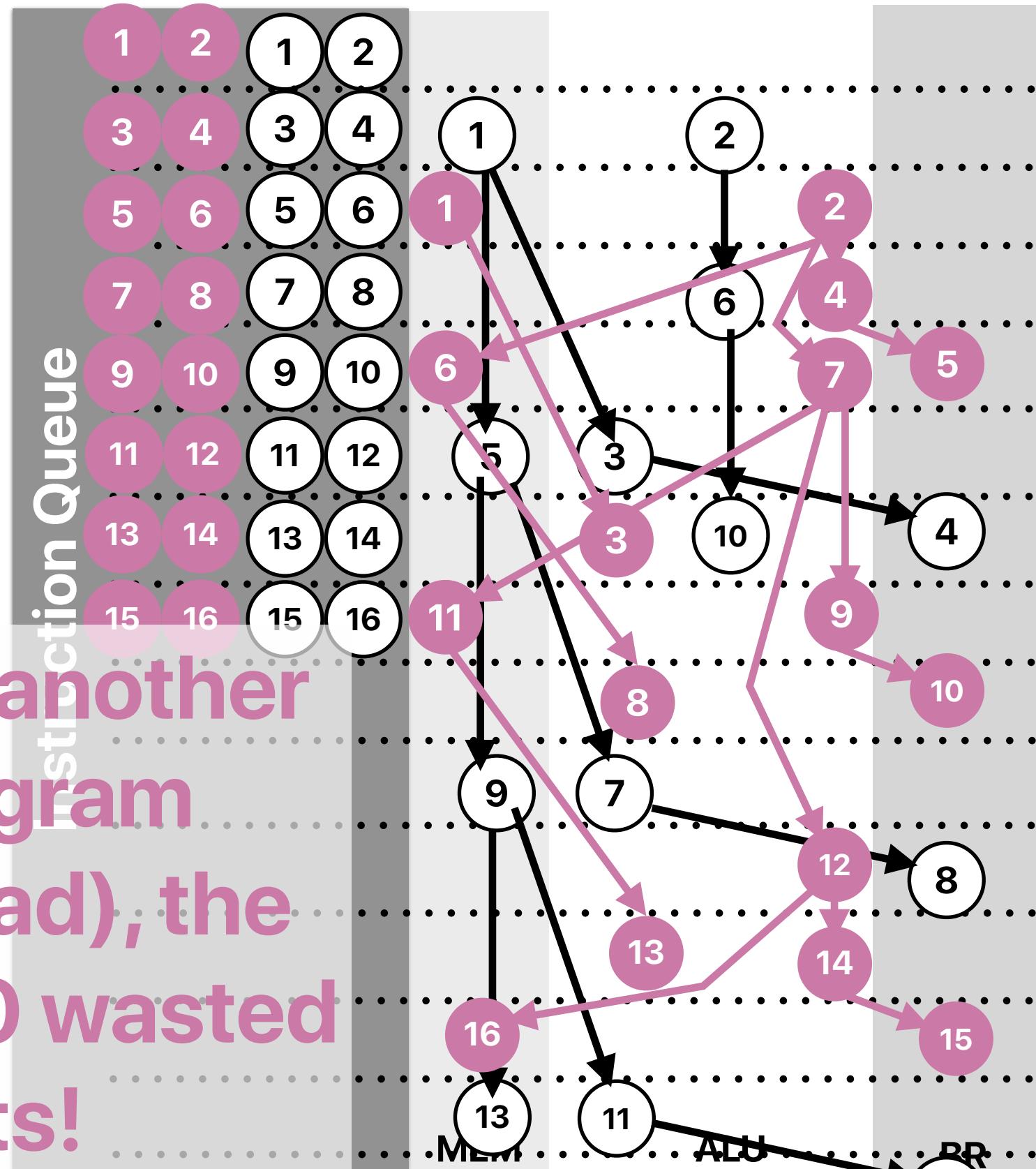


We're wasting the issue bandwidth in 8 out of 12 cycles

Concept: Simultaneous Multithreading (SMT)

① movq 8(%rdi), %rdi
② addl \$1, %eax
③ testq %rdi, %rdi
④ jne .L3
⑤ movq 8(%rdi), %rdi
⑥ addl \$1, %eax
⑦ testq %rdi, %rdi
⑧ jne .L3
⑨ movq 8(%rdi), %rdi
⑩ addl \$1, %eax
⑪ testq %rdi, %rdi
⑫ jne .L3
⑬ movq 8(%rdi), %rdi
⑭ addl \$1, %eax
⑮ testq %rdi, %rdi
⑯ jne .L3
⑰ movl (%rdi), %ecx

By scheduling another running program instance (thread), the processor has 0 wasted issue slots!



- ① movl (%rdi), %ecx
- ② addq \$4, %rdi
- ③ addl %ecx, %eax
- ④ cmpq %rdx, %rdi
- ⑤ jne .L3
- ⑥ movl (%rdi), %ecx
- ⑦ addq \$4, %rdi
- ⑧ addl %ecx, %eax
- ⑨ cmpq %rdx, %rdi
- ⑩ jne .L3
- ⑪ movl (%rdi), %ecx
- ⑫ addq \$4, %rdi
- ⑬ addl %ecx, %eax
- ⑭ cmpq %rdx, %rdi
- ⑮ jne .L3
- ⑯ movl (%rdi), %ecx

SMT from the user/OS' perspective



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!

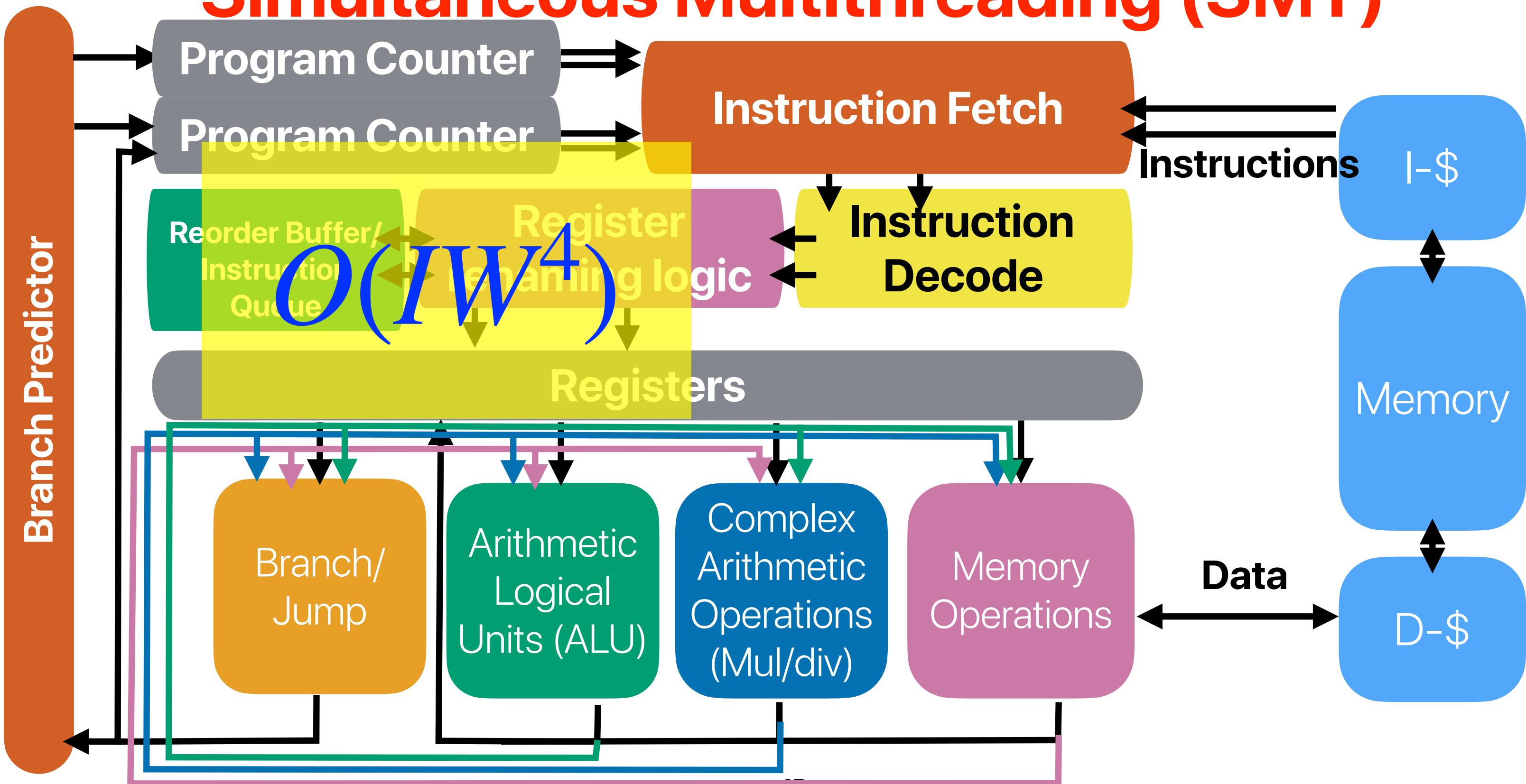
Architecture:	x86_64
CPU op-mode(s):	32-bit, 64-bit
Byte Order:	Little Endian
Address sizes:	39 bits physical, 48 bits virtual
CPU(s):	8
On-line CPU(s) list:	0-7
Thread(s) per core:	2
Core(s) per socket:	4
Socket(s):	1
NUMA node(s):	1
Vendor ID:	GenuineIntel
CPU family:	6
Model:	151
Model name:	12th Gen Intel(R) Core(TM) i3-12100F
Stepping:	5
CPU MHz:	3300.000
CPU max MHz:	5500.0000
CPU min MHz:	800.0000
BogoMIPS:	6604.80
Virtualization:	VT-x
L1d cache:	192 KiB
L1i cache:	128 KiB

Outline

- Multithreaded processors
- Programming multithreaded processors & necessary architectural support

Architecture:	x86_64
CPU op-mode(s):	32-bit, 64-bit
Byte Order:	Little Endian
Address sizes:	39 bits physical, 48 bits virtual
CPU(s):	8
On-line CPU(s) list:	0-7
Thread(s) per core:	2
Core(s) per socket:	4
Socket(s):	1
NUMA node(s):	1
Vendor ID:	GenuineIntel
CPU family:	6
Model:	151
Model name:	12th Gen Intel(R) Core(TM) i3-12100F
Stepping:	5
CPU MHz:	3300.000
CPU max MHz:	5500.0000
CPU min MHz:	800.0000
BogoMIPS:	6604.80
Virtualization:	VT-x
L1d cache:	192 KiB
L1i cache:	128 KiB

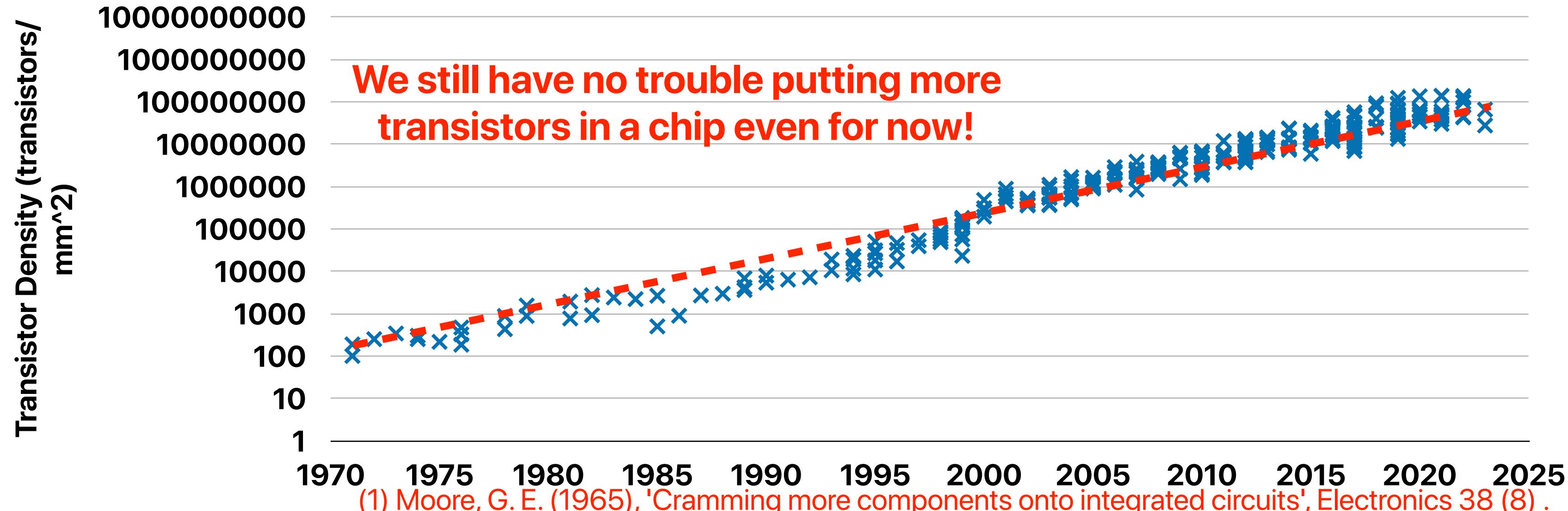
Simultaneous Multithreading (SMT)



Chip-Multiprocessors (CMP) or Multi-core processors

Recap: Moore's Law⁽¹⁾

- The number of transistors we can build in a fixed area of silicon doubles every 12 ~ 24 months.
- Moore's Law "was" the most important driver for historic CPU performance gains



Transistor Counts

Microarchitecture	Transistor Count	Issue-width	Year
Alder Lake	325 M	5x ALU, 7x Memory	2021
Coffee Lake	217 M	4x ALU, 4x Memory	2017
Sandy Bridge	290 M	3x ALU, 3x Memory	2011
Nehalem	182.75 M	3x ALU, 3x Memory	2008



How many transistors per core on Coffee Lake?



The Coffee Lake processor has 217 million transistors per core. It is manufactured using Intel's second 14 nm process. The Coffee Lake processors introduced i5 and i7 CPUs featuring six cores (along with hyper-threading in the case of the later models) and no hyperthreading.



The transistor count per core on Coffee Lake is lower than that of some other modern processors, such as the Ryzen 5 5600X, which has 390 million transistors per core. However, Coffee Lake still offers good performance, thanks to its high clock speeds and efficient design.

Here is a table of the transistor counts per core for some other modern processors:

Processor	Transistors per core
Coffee Lake	217 million
Ryzen 5 5600X	390 million
Core i7 12700K	475.8 million

Nehalem Alder Lake
6-issue 12-issue

Recap: do we really need very wide issue processors?

	ET	IC	IPC/ILP	# of branches	Branch mis-prediction rate
A	22.21	332 Trillions	2.88	65 Trillions	1.13%
B	12.29	287 Trillions	4.52	17 Trillions	0.04%
C	5.01	102 Trillions	3.95	17 Trillions	0.04%
D	3.73	80 Trillions	4.13	1 Trillions	~0%
E	54.4	173 Trillions	0.61	44 Trillions	18.6%
SSE4.2	1.57	22 Trillions	2.7	1 Trillions	~0%

We don't we build processors with more cores?

Microarchitecture	Transistor Count	Issue-width	Year
Alder Lake	325 M	5x ALU, 7x Memory	2021
Coffee Lake	217 M	4x ALU, 4x Memory	2017
Sandy Bridge	290 M	3x ALU, 3x Memory	2011
Nehalem	182.75 M	3x ALU, 3x Memory	2008



How many transistors per core on Coffee Lake?



The Coffee Lake processor has 217 million transistors per core. It is manufactured using Intel's second 14 nm process. Coffee Lake processors introduced i5 and i7 CPUs featuring six cores (along with hyper-threading in the case of the latter) and no hyperthreading.



The transistor count per core on Coffee Lake is lower than that of some other modern processors, such as the Ryzen 5 5600X, which has 390 million transistors per core. However, Coffee Lake still offers good performance, thanks to its high clock speeds and efficient power delivery.

Here is a table of the transistor counts per core for some other modern processors:

Processor	Transistors per core
Coffee Lake	217 million
Ryzen 5 5600X	390 million
Core i7 10700K	425.8 million

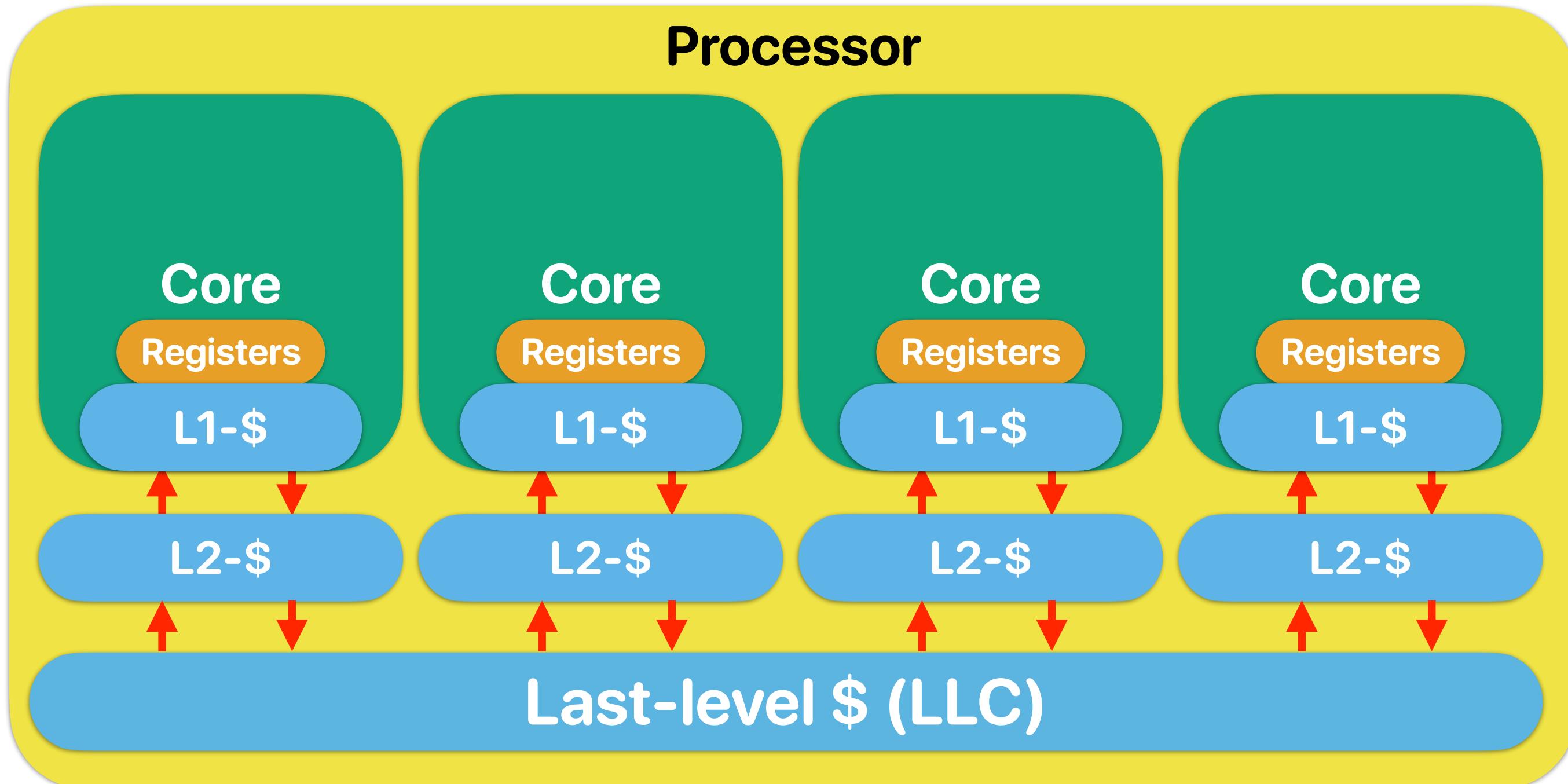
2x 3-issue ALUs Nehalem

Nehalem Alder Lake Nehalem
6-issue 12-issue 6-issue

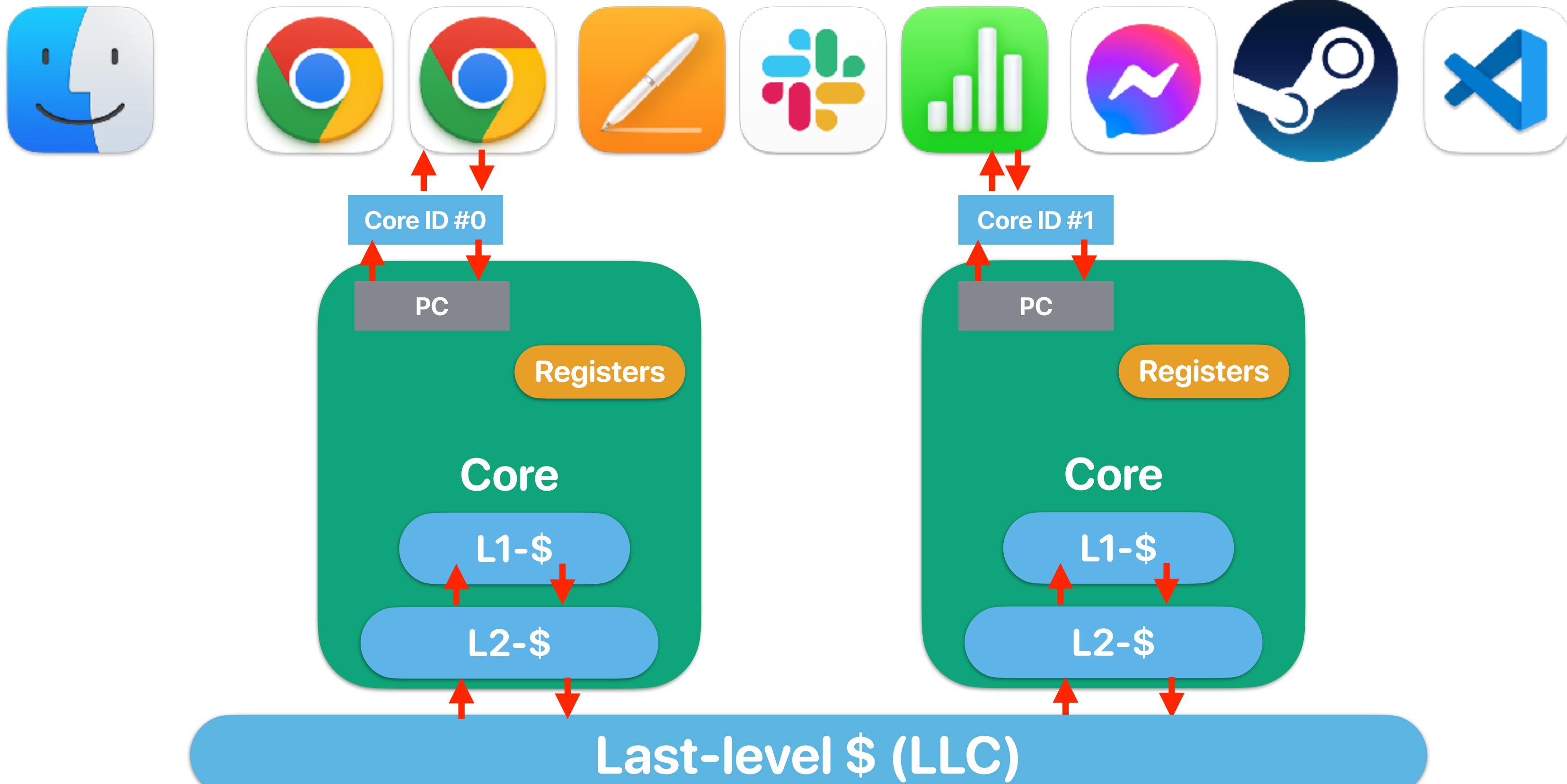
1x 5-issue ALUs Alder Lake

Based on https://en.wikipedia.org/wiki/Transistor_count

Concept of CMP



CMP from the user/OS' perspective





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 **There is no clear win on each — why not having both?**

B. 2

C. 3

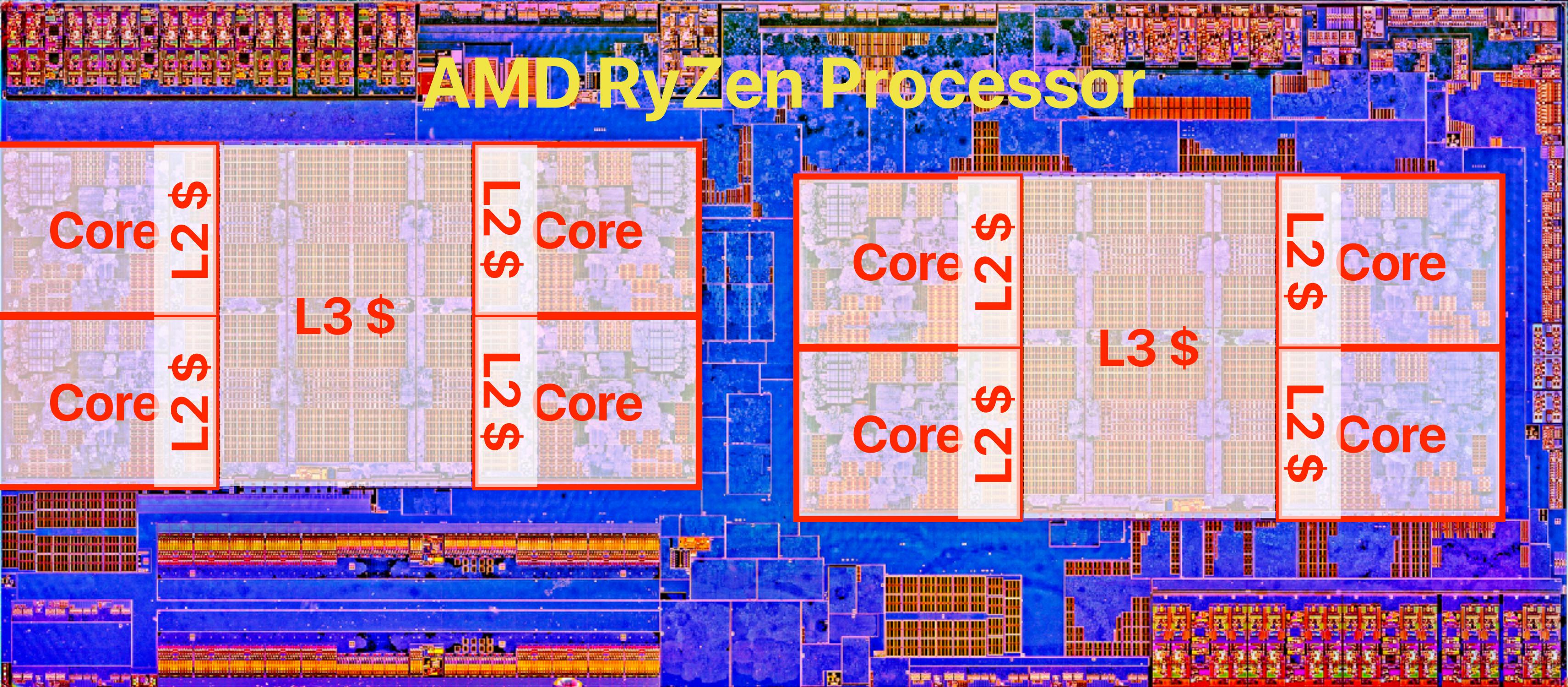
D. 4

E. 5

Modern processors have both CMP/SMT



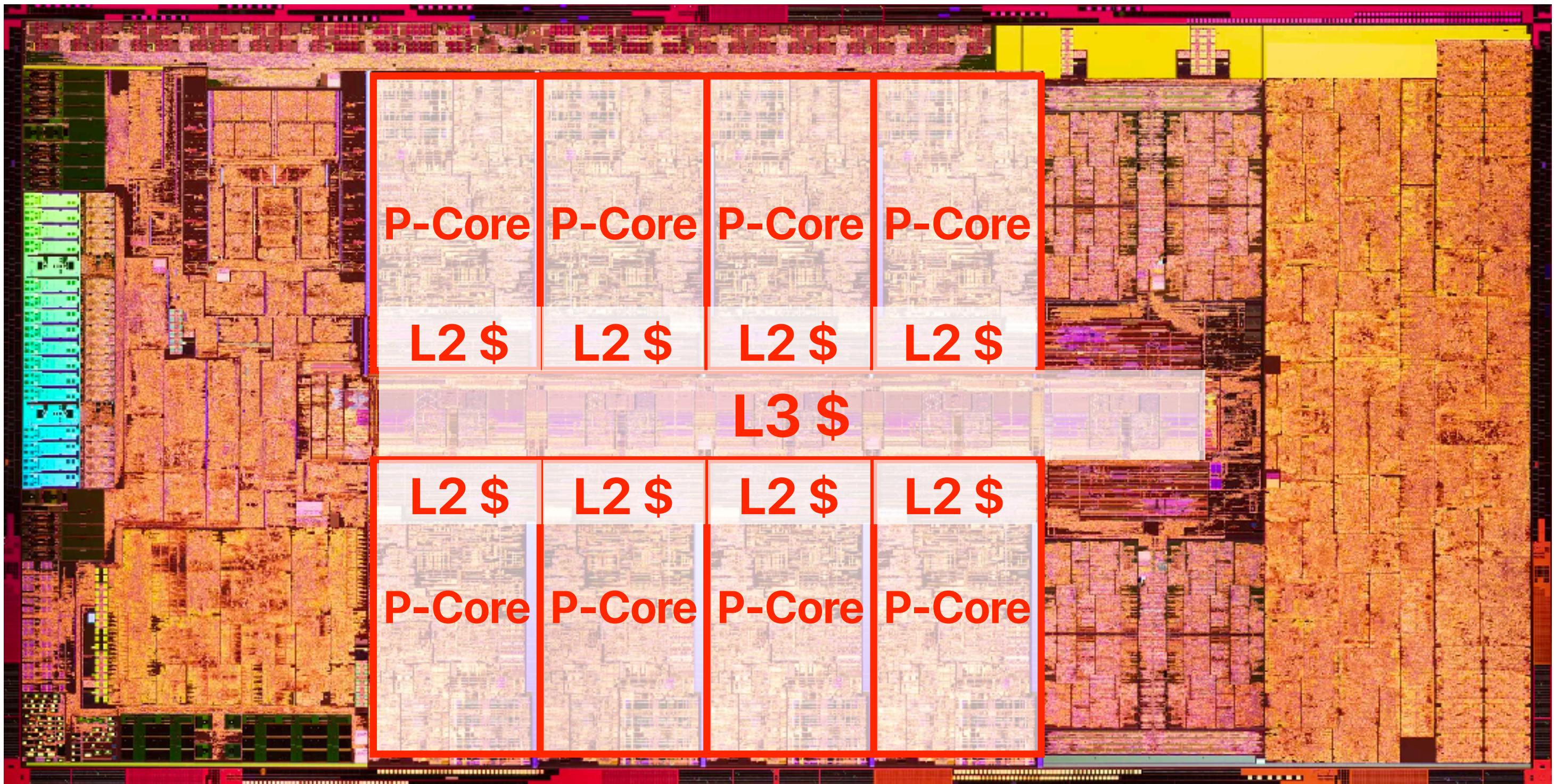
AMD Ryzen Processor



AMD

RYZEN

Intel Alder Lake



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 **There is no clear win on each — why not having both?**

B. 2

C. 3

D. 4 **The only thing we know for sure — if we don't parallel the program, it won't get any faster on CMP**

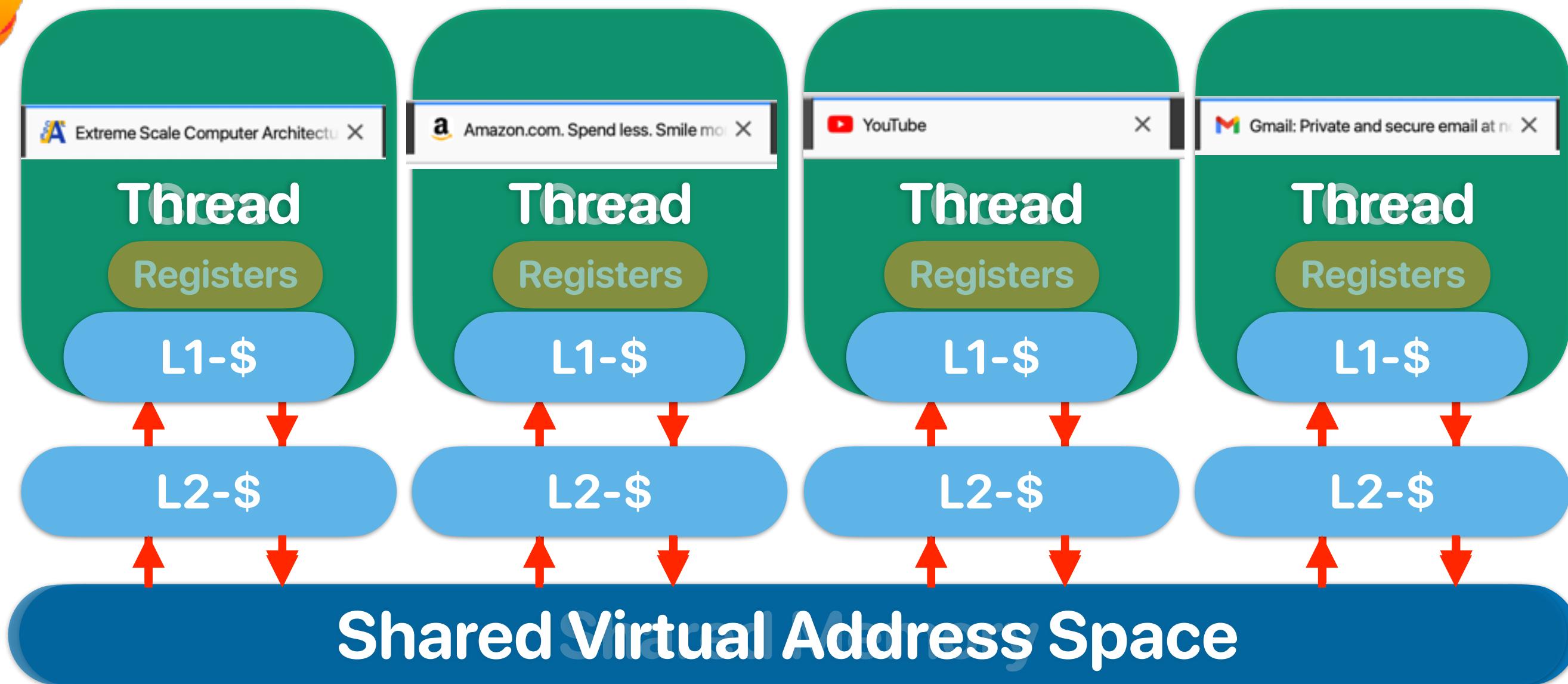
E. 5

Parallel Programming & Architectural Supports for Parallel Programming

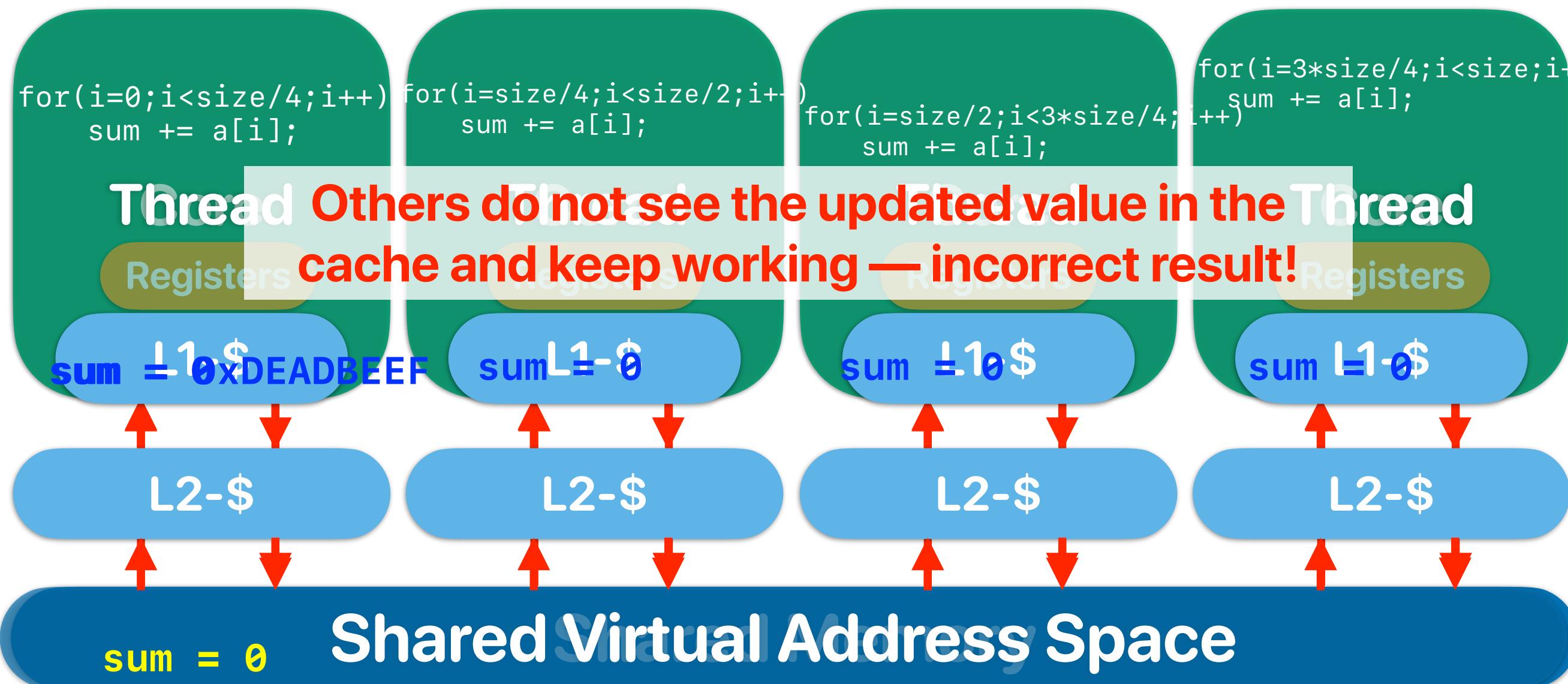
Parallel programming

- 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

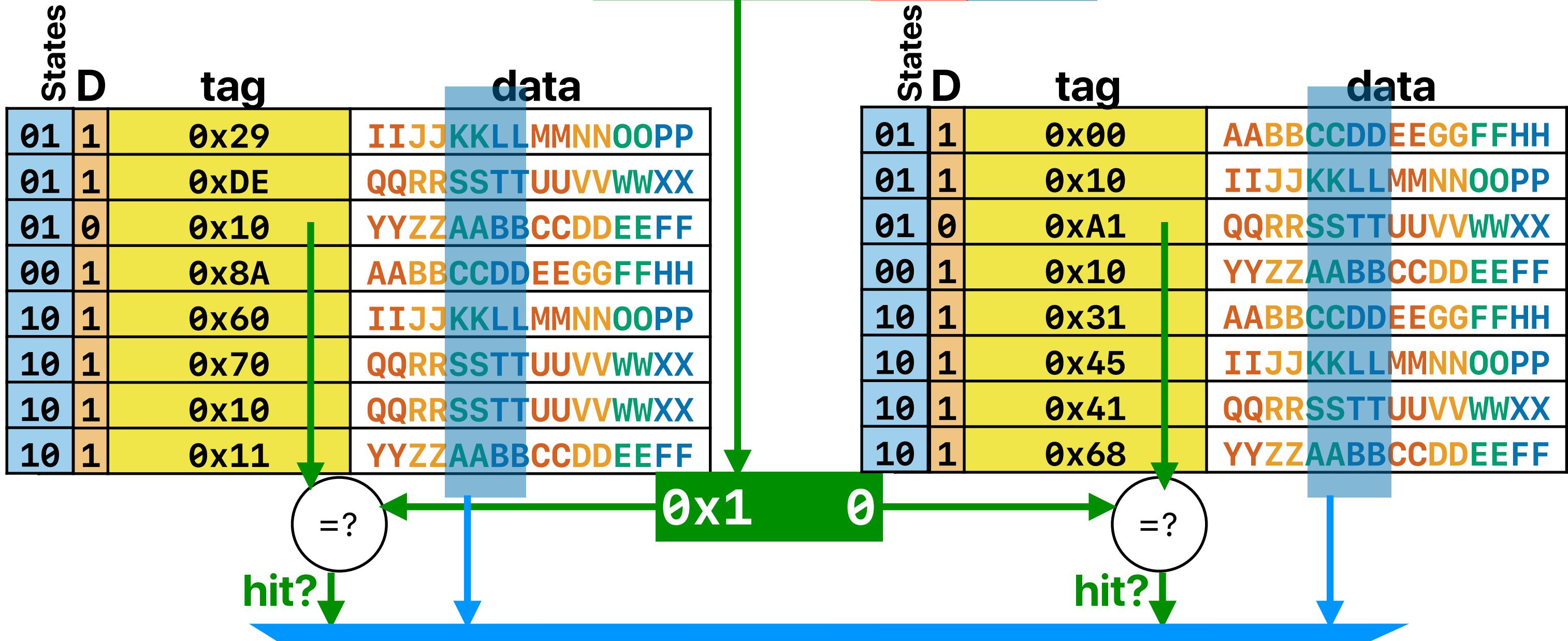
Coherent way-associative cache

memory address:

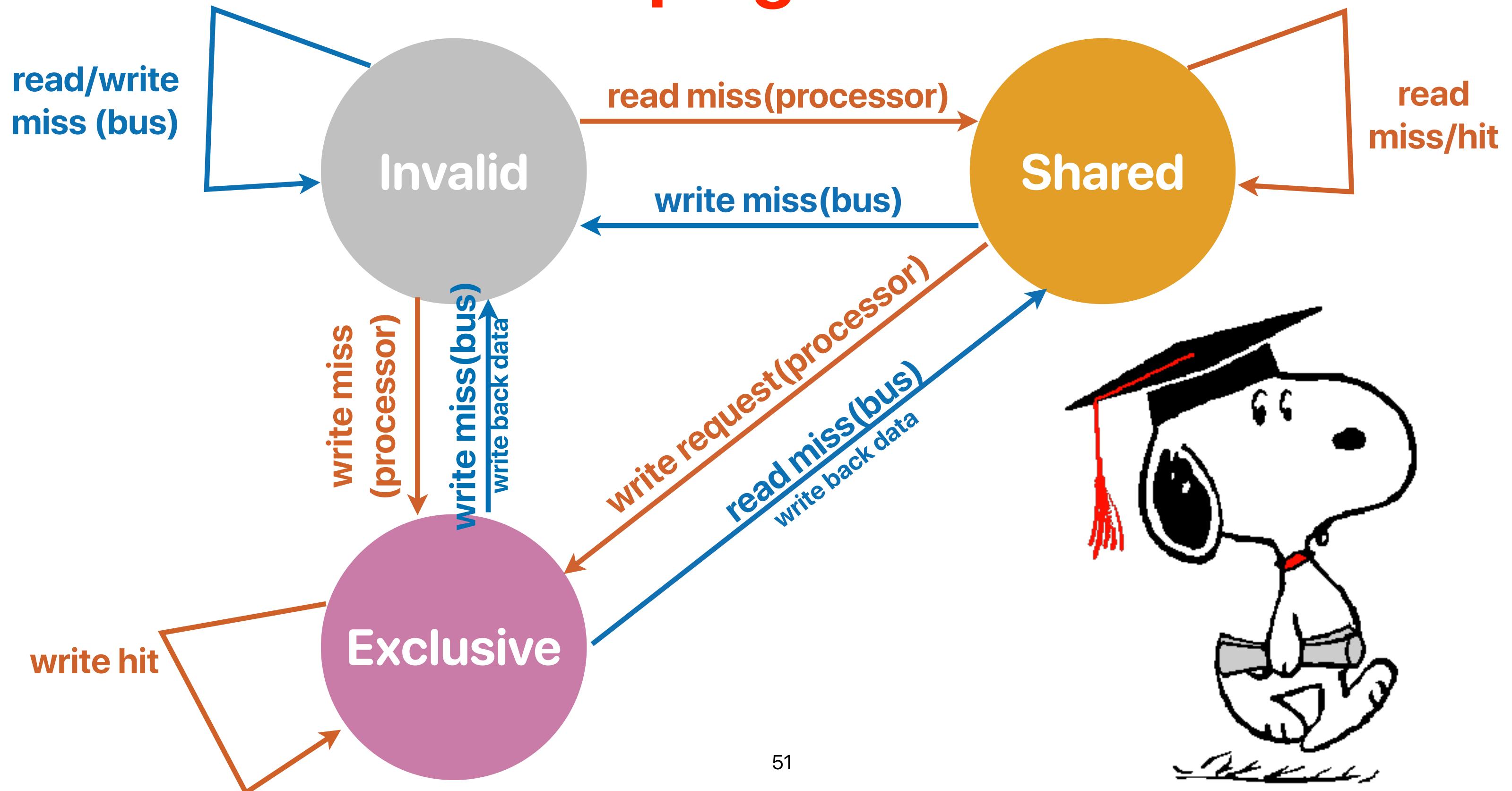
$0x0$ 8 tag 2 set 4 block
index offset

memory address:

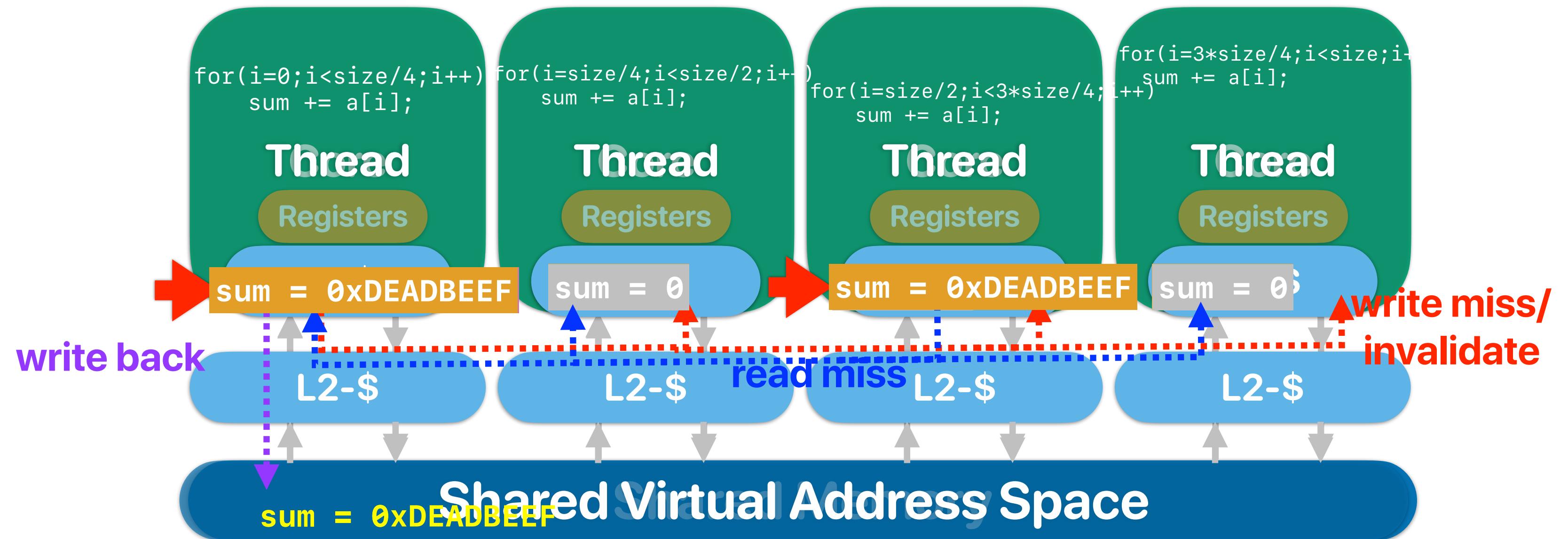
0b0000100000100100



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",</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",
```

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 this Thursday
- iEVAL — if you fill out iEVAL, please submit a screenshot
 - You will receive a full credit reading quiz if you submit the screenshot before the in-person final examine
 - We will drop your **lowest two** reading quizzes — so that if you don't want to fill the form, you still have the lowest one dropped
- Final Exam — will release some sample questions next Tuesday at the end of the lecture
 - 12/7 3:30p—4:50p
 - Including CSMS comprehensive examine questions
 - In-person, no outside materials
 - 12/11-12/14 6pm—6pm — any consecutive three-hour slot you pick
 - Online examine
 - Including coding assessments (similar to company interviews)
 - Essay-style questions that require thorough understandings of course materials

Computer Science & Engineering

203

つづく

