

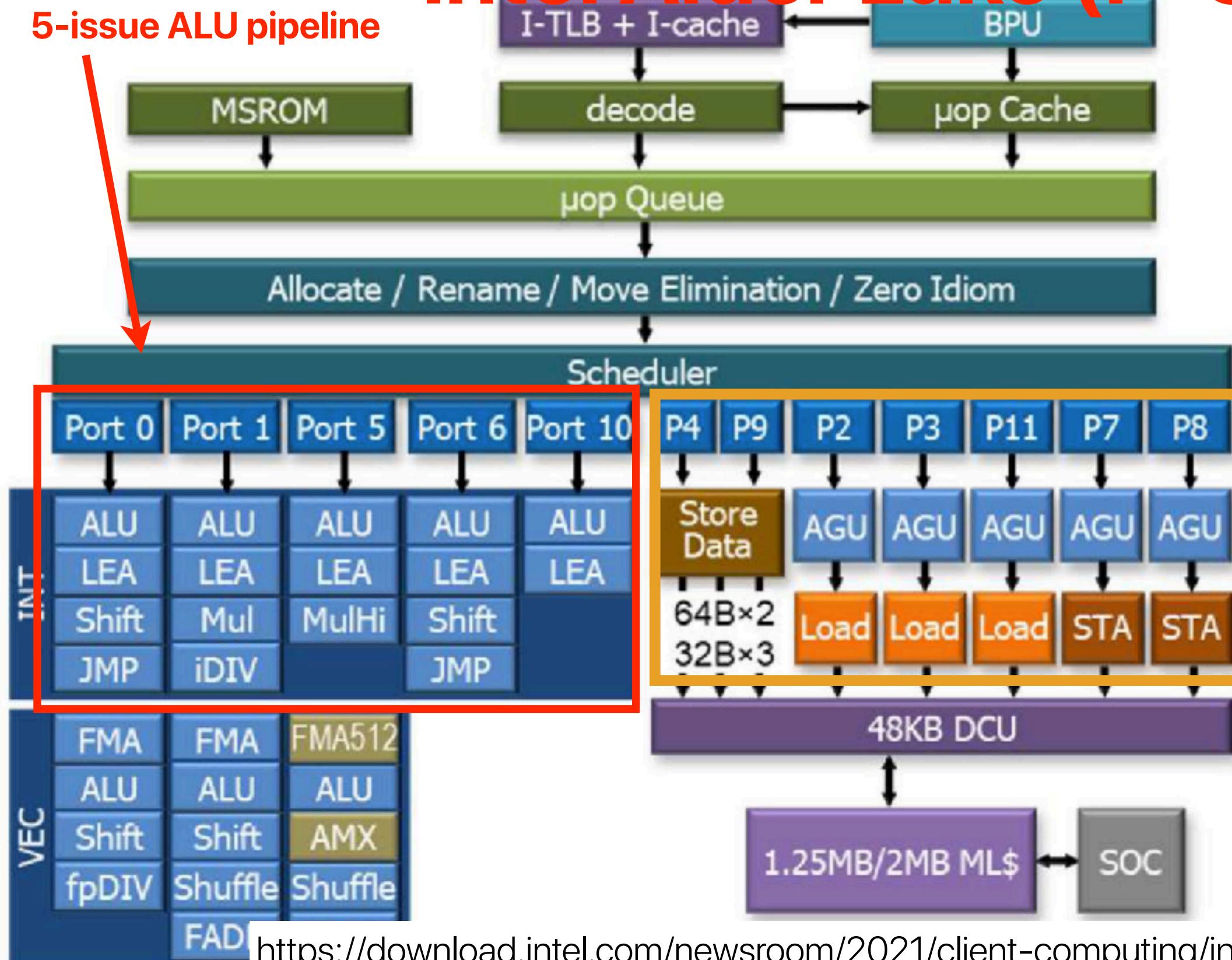
# Programming on Modern Processors

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
  - Tournament predictors

# Intel Alder Lake (P-Core)



$$MinCPI = \frac{1}{12}$$

$$MinINTInst . CPI = \frac{1}{5}$$

$$MinMEMInst . CPI = \frac{1}{7}$$

$$MinBRInst . CPI = \frac{1}{2}$$

# Outline

- Programming on modern processors — exploiting instruction-level parallelism
- Simultaneous multithreading

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

If we cannot improve the performance of executing

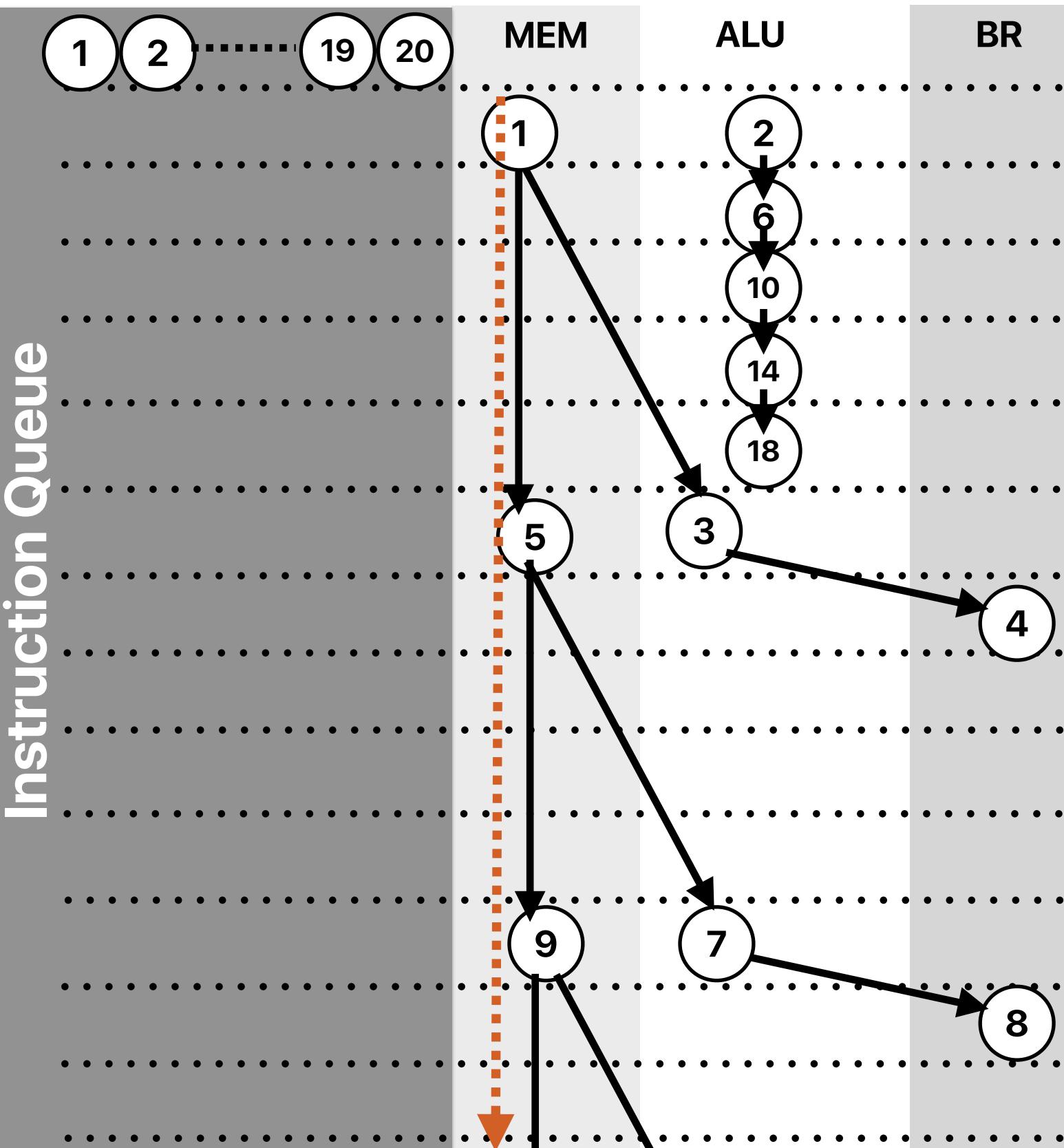
`movq 8(%rdi), %rdi`

we cannot improve the execution time.

That's the “critical path”!

```
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`



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

If we cannot improve the performance of executing

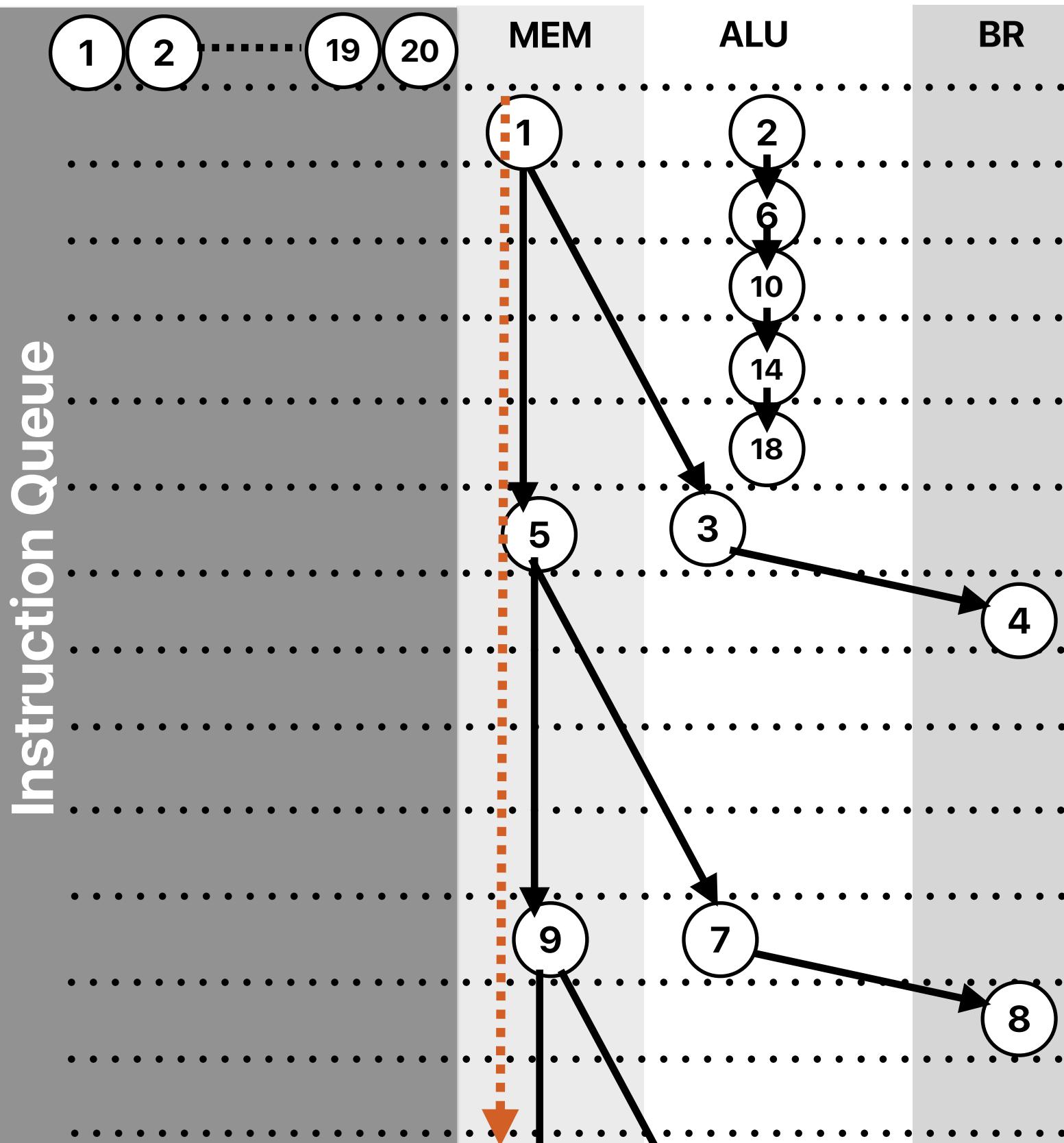
`movq 8(%rdi), %rdi`

we cannot improve the execution time.

That's the “critical path”!

```
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`



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

If we cannot improve the performance of executing

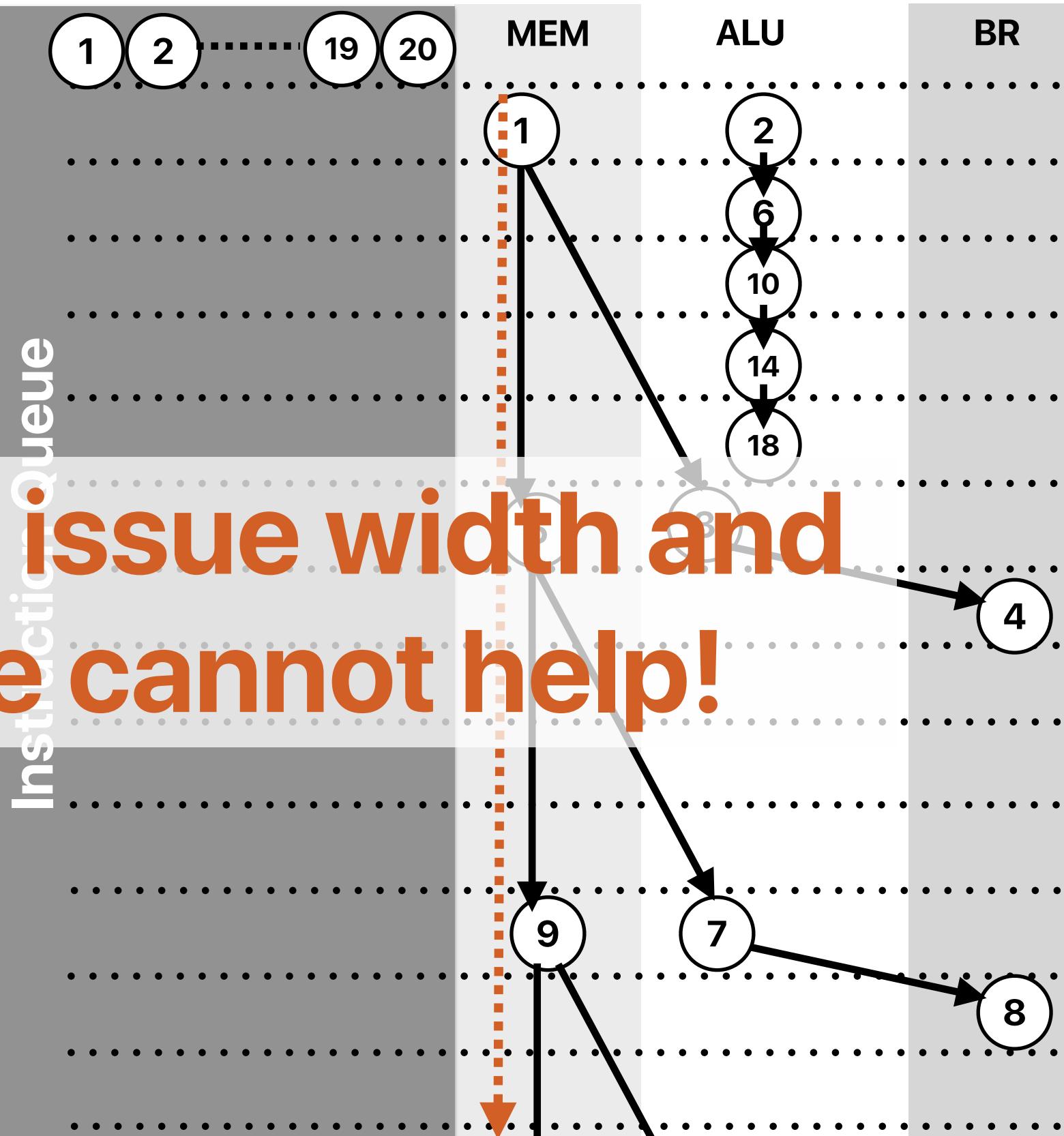
`movq 8(%rdi), %rdi`

we cannot improve the execution time.

That's the “critical path”!

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

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



# Takeaways: programming modern processors

- The key to efficient code is exploiting as much instruction-level parallelism (ILP) or say instructions per cycle (IPC) as possible

# Problem: Popcount

- The population count (or popcount) of a specific value is the number of set bits (i.e., bits in 1s) in that value.
- Applications
  - Parity bits in error correction/detection code
  - Cryptography
  - Sparse matrix
  - Molecular Fingerprinting
  - Implementation of some succinct data structures like bit vectors and wavelet trees.

# Problem: Popcount

- Given a 64-bit integer number, find the number of 1s in its binary representation.
- Example 1:

Input: 59487

Output: 9

Explanation: 59487's binary representation is  
0b10110010100001111

```
int main(int argc, char *argv[]) {  
    uint64_t key = 0xdeadbeef;  
  
    int count = 1000000000;  
    uint64_t sum = 0;  
  
    for (int i=0; i < count; i++)  
    {  
        sum += __builtin_popcount(key);  
    }  
    printf("Result: %lu\n", sum);  
    return sum;  
}
```

# Takeaways: programming modern processors

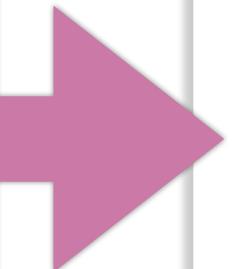
- The key to efficient code is exploiting as much instruction-level parallelism (ILP) or say instructions per cycle (IPC) as possible
- Loop unrolling is effective as control overhead is still significant despite we have branch predictors and OoO.

# Takeaways: programming modern processors

- The key to efficient code is exploiting as much instruction-level parallelism (ILP) or say instructions per cycle (IPC) as possible
- Loop unrolling is effective as control overhead is still significant despite we have branch predictors and OoO.
- With caches, we can potentially use small lookup tables to replace more expensive data dependent operations

# Loop unrolling eliminates all branches!

```
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;
}
```

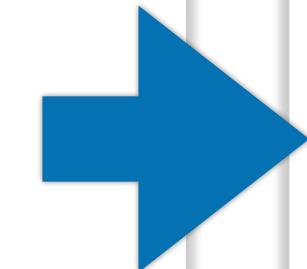


# Takeaways: programming modern processors

- The key to efficient code is exploiting as much instruction-level parallelism (ILP) or say instructions per cycle (IPC) as possible
- Loop unrolling is effective as control overhead is still significant despite we have branch predictors and OoO.
- With caches, we can potentially use small lookup tables to replace more expensive data dependent operations
- Making your code more predictable is the key!
  - Compilers can confidently perform aggressive optimizations
  - Branch predictors can be more accurate
  - Cache miss rate can be really low

# Why is E the slowest?

```
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;
}
```



.L11:

```
    movq    %r9, %rcx
    andl    $15, %ecx
    movslq  (%r8,%rcx,4), %rcx
    addq    %r8, %rcx
    notrack jmp     *%rcx
```

.L7:

```
    .long   .L5-.L7
    .long   .L10-.L7
    .long   .L10-.L7
    .long   .L9-.L7
    .long   .L10-.L7
    .long   .L9-.L7
    .long   .L9-.L7
    .long   .L8-.L7
    .long   .L10-.L7
    .long   .L9-.L7
    .long   .L9-.L7
    .long   .L8-.L7
    .long   .L9-.L7
    .long   .L8-.L7
    .long   .L9-.L7
    .long   .L8-.L7
    .long   .L6-.L7
```

.L8:

```
    addl    $3, %eax
```

.L5:

```
    shrq    $4, %r9
    subq    $1, %rsi
    jne     .L11
```

```
    cltq
    addq    %rax, %rbx
    subl    $1, %edi
    jne     .L12
```

.L9:

```
.cfi_restore_state
    addl    $2, %eax
    jmp    .L5
    .p2align 4,,10
    .p2align 3
```

.L10:

```
    addl    $1, %eax
    jmp    .L5
    .p2align 4,,10
    .p2align 3
```

.L6:

```
    addl    $4, %eax
    jmp    .L5
```

# Hardware acceleration

- Because `popcount` is important, both intel and AMD added a `POPCNT` instruction in their processors with SSE4.2 and SSE4a
- In C/C++, you may use the intrinsic “`_mm_popcnt_u64`” to get # of “1”s in an unsigned 64-bit number
  - You need to compile the program with `-m64 -msse4.2` flags to enable these new features

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

# 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%

# Takeaways: programming modern processors

- The key to efficient code is exploiting as much instruction-level parallelism (ILP) or say instructions per cycle (IPC) as possible
- Loop unrolling is effective as control overhead is still significant despite we have branch predictors and OoO.
- With caches, we can potentially use small lookup tables to replace more expensive data dependent operations
- Making your code more predictable is the key!
  - Compilers can confidently perform aggressive optimizations
  - Branch predictors can be more accurate
  - Cache miss rate can be really low
- If there is a hardware feature supporting the desire computation — we should try it!

# 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

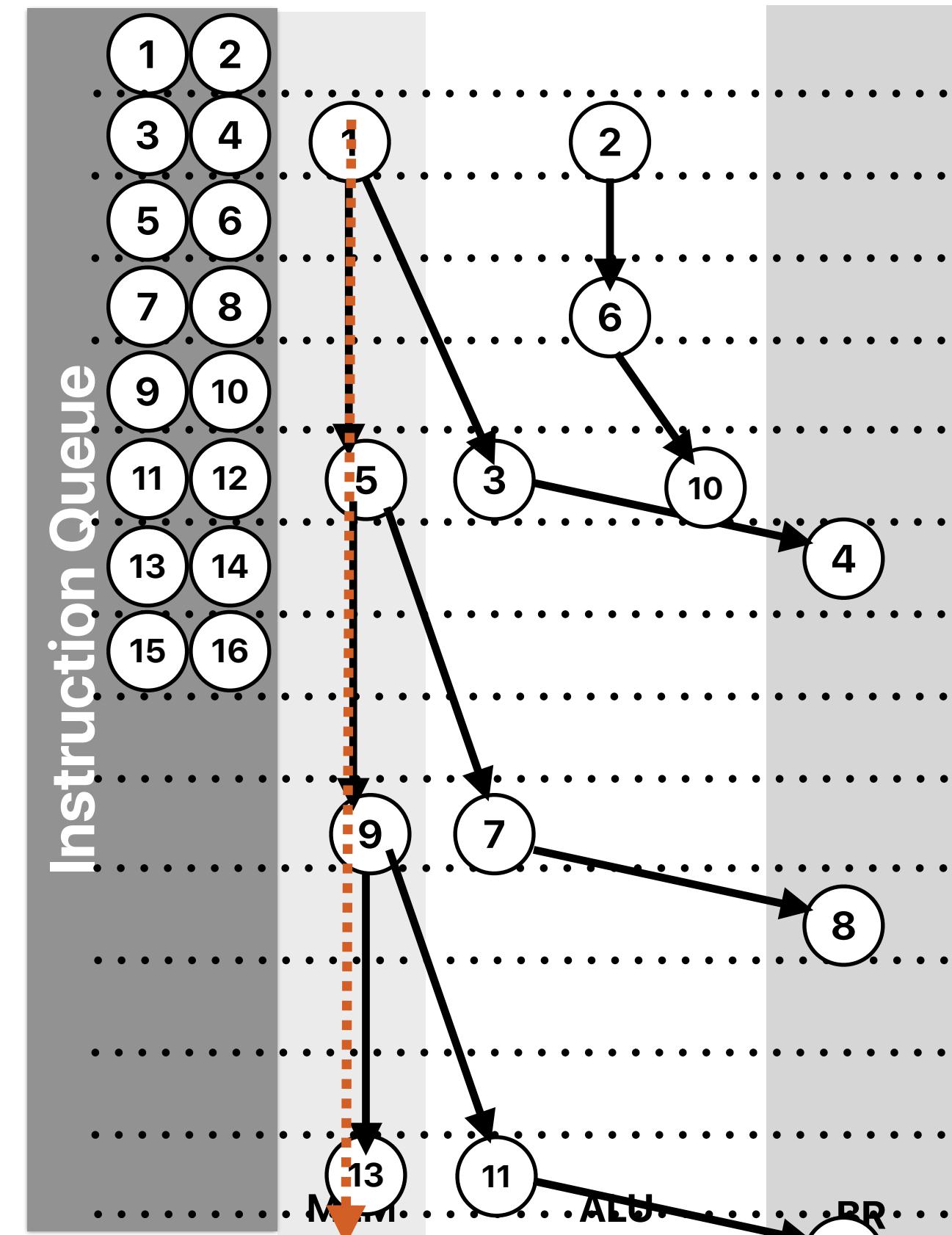
# What if we have “unlimited” fetch/issue width — “linked list”

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
```

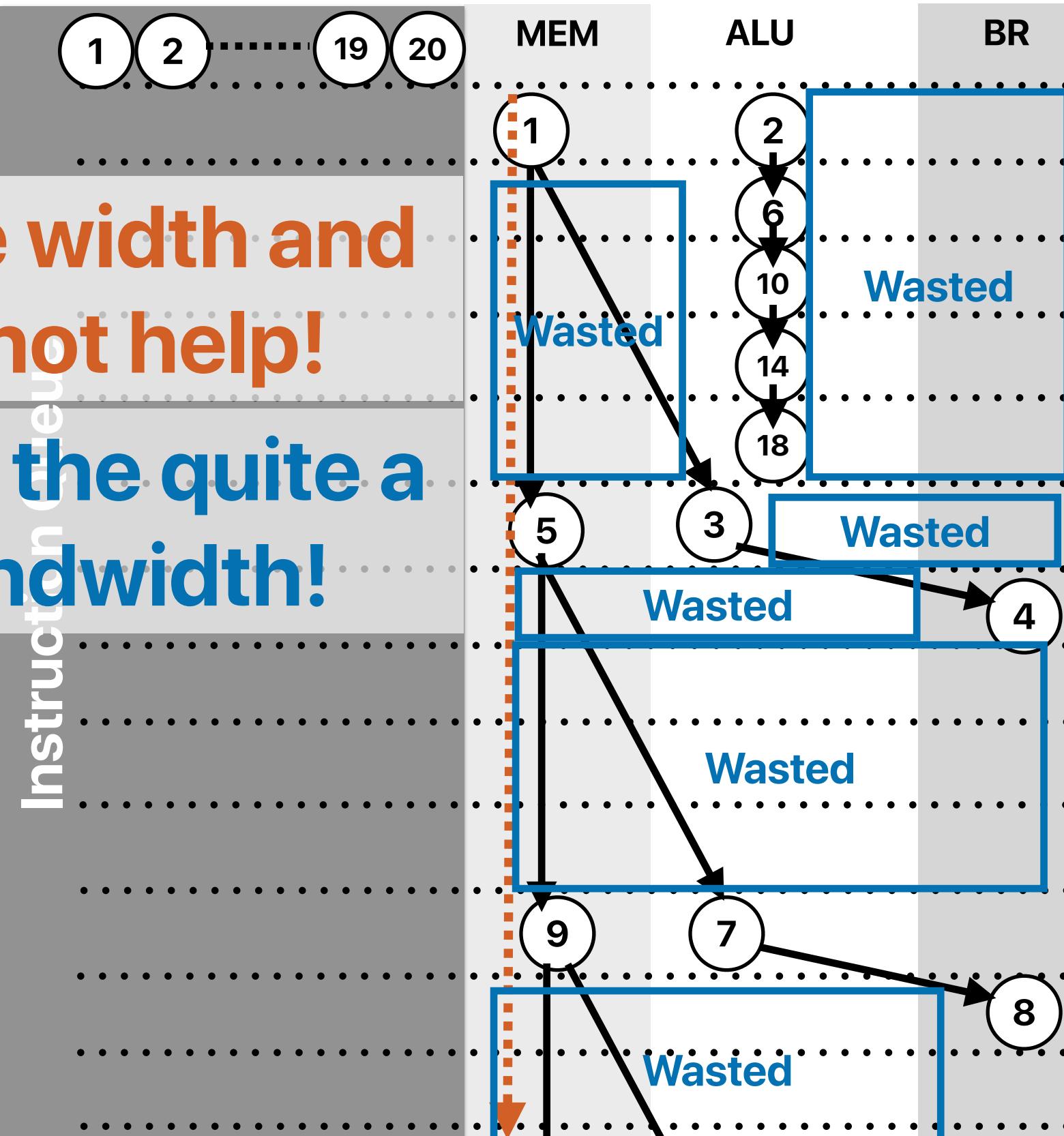


# What if we have “unlimited” fetch/issue width — “linked list”

Even unlimited issue width and  
a perfect cache cannot help!

```
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



# 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

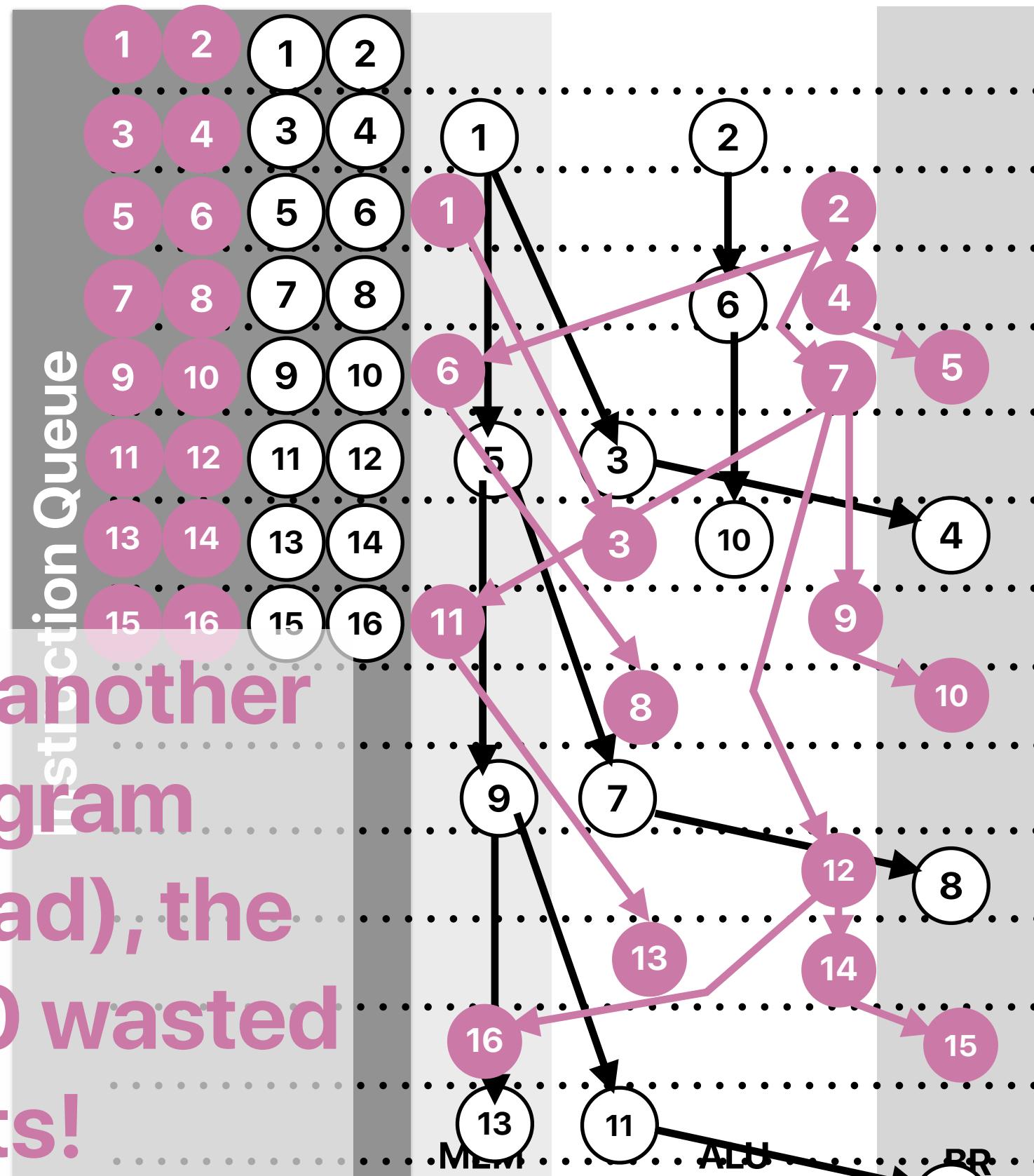
# **Parallel architectures**

# **Simultaneous multithreading**

# 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

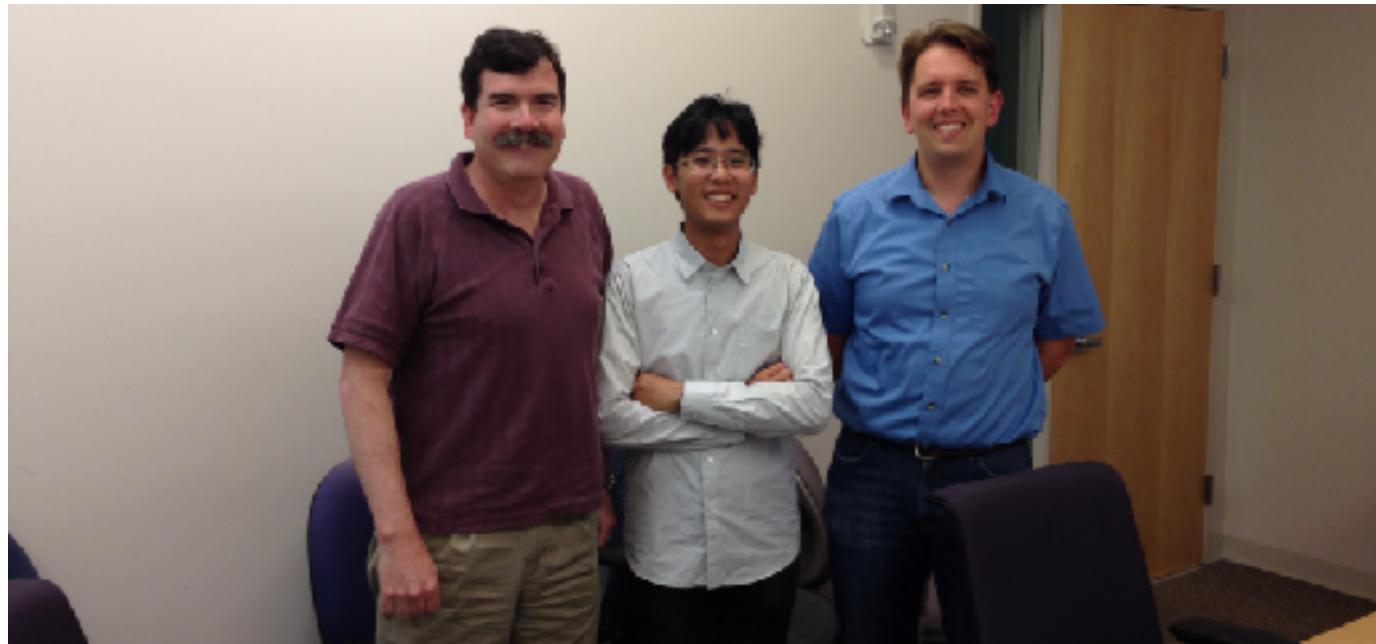
By scheduling another running program instance (thread), the processor has 0 wasted issue slots!



- 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

# 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
- Invented by Dean Tullsen (Now a professor at **UCSD CSE**)



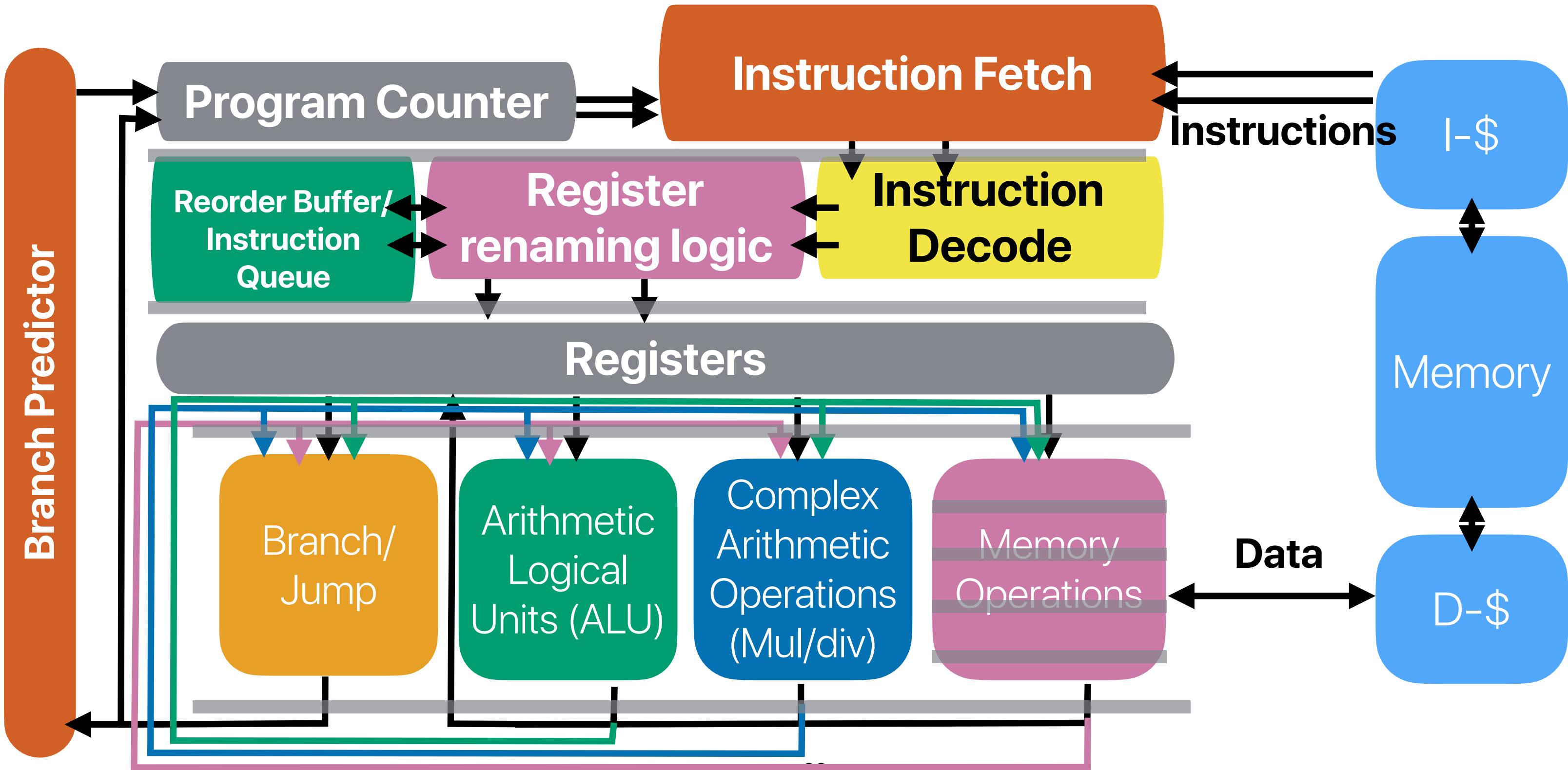
# SMT from the user/OS' perspective



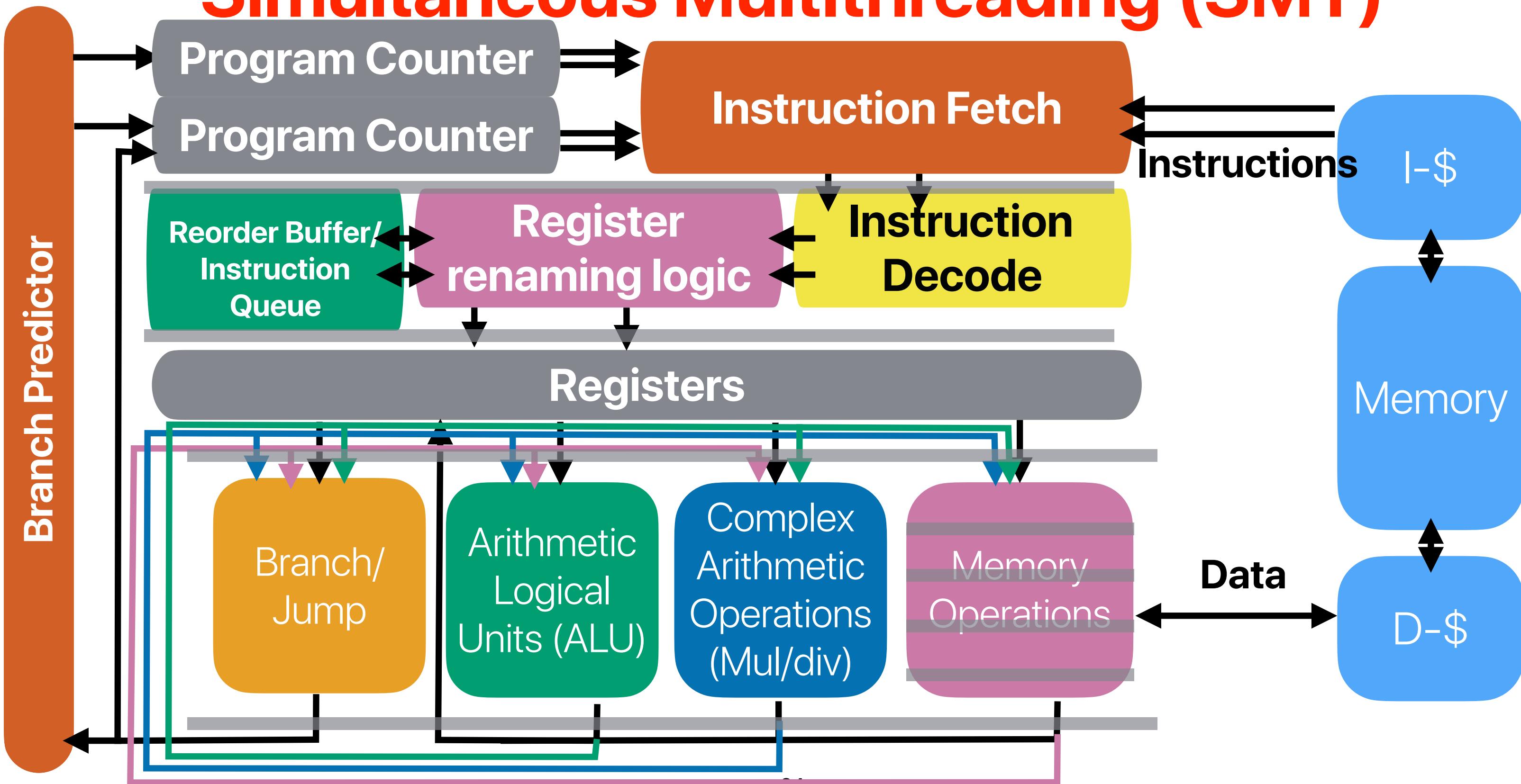
# How do we support two running programs in one pipeline?

- We need two program counters
- We need two sets of architectural to physical register mappings
- We do not need
  - Duplicated cache — virtually indexed, physically tagged cache already addressed that
  - Duplicated pipeline functional units — isn't sharing the whole purpose?
  - Duplicated reorder buffer — you simply need to tag which process the instruction belongs to

# Recap: Register renaming



# Simultaneous Multithreading (SMT)



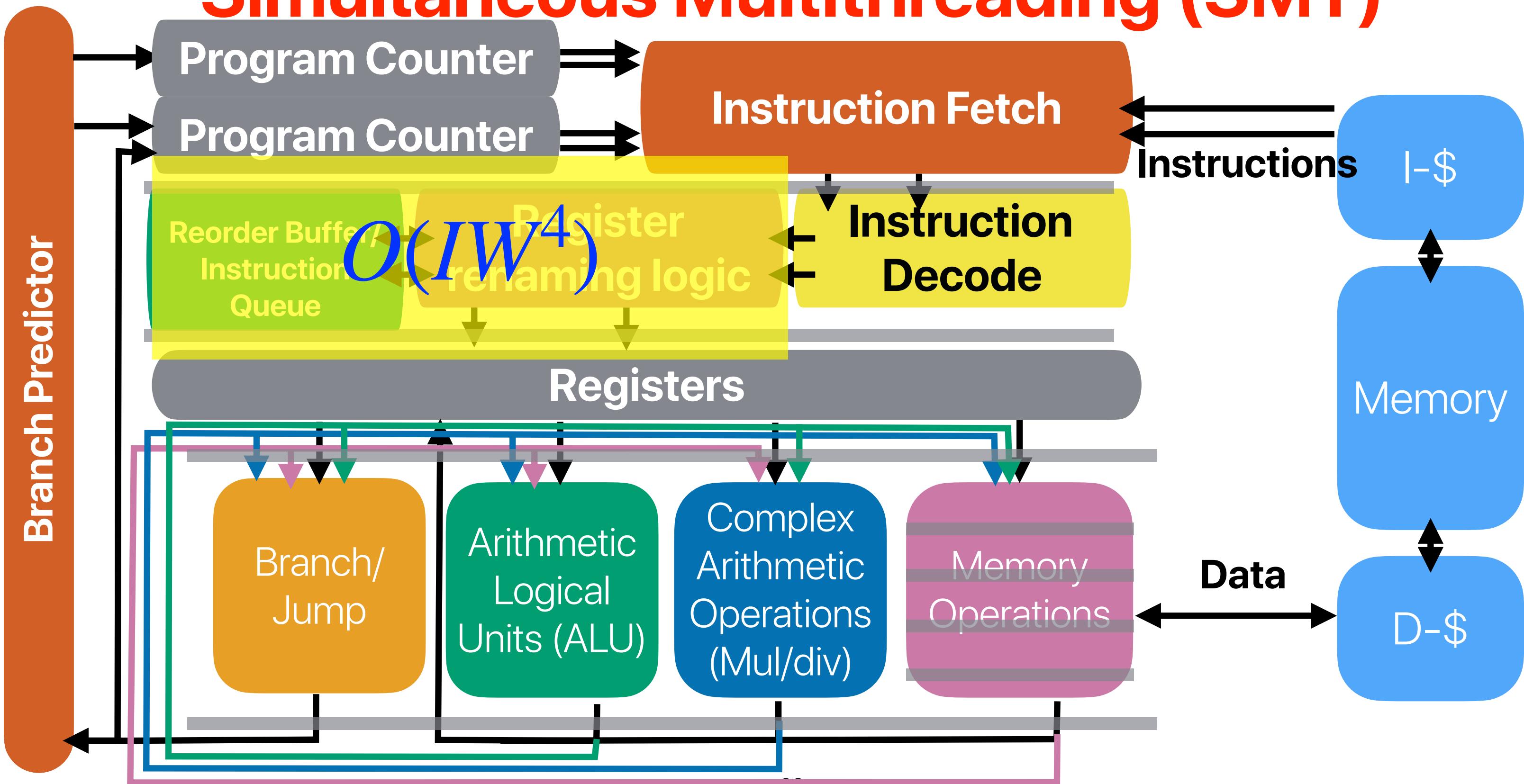
Architecture:	x86_64
CPU op-mode(s):	32-bit, 64-bit
Byte Order:	Little Endian
Address sizes:	39 bits physical, 48 bits virtual
CPU(s):	8
On-line CPU(s) list:	0-7
Thread(s) per core:	2
Core(s) per socket:	4
Socket(s):	1
NUMA node(s):	1
Vendor ID:	GenuineIntel
CPU family:	6
Model:	151
Model name:	12th 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

# SMT in real practice

- Intel HyperThreading (supports up to two threads per core)
  - Intel Pentium 4, Intel Atom, All Intel Core i7s and Core i9s
- AMD RyZen (Zen microarchitecture)
- If you see a processor with “threads” more than “cores”, that must 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

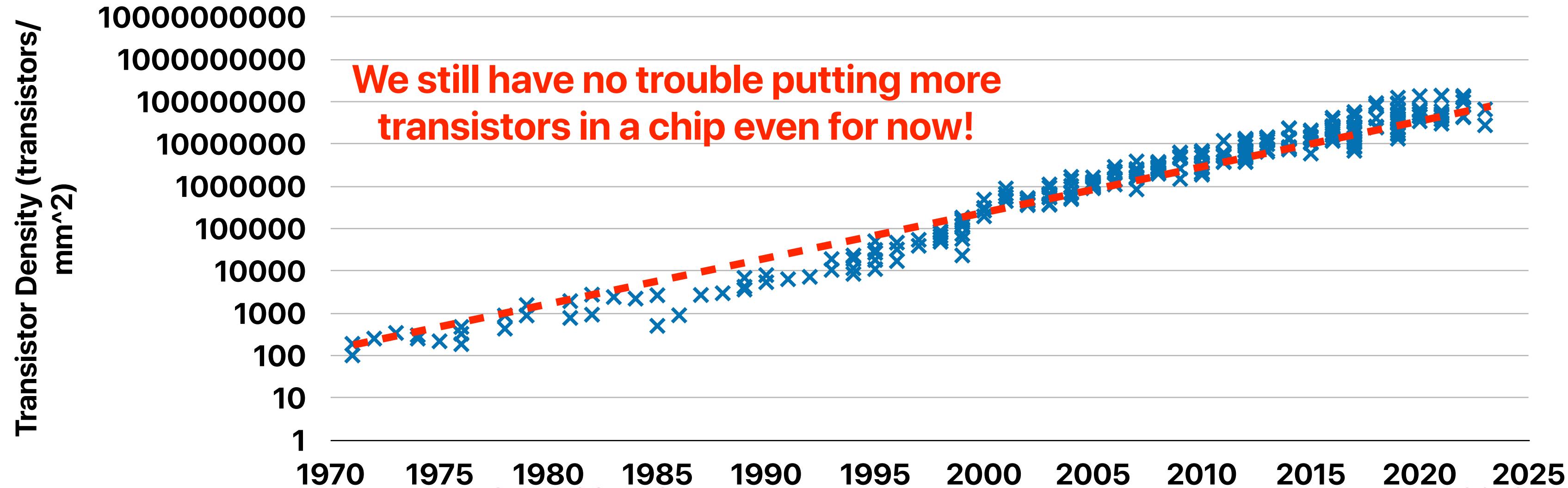
# Simultaneous Multithreading (SMT)



# **Chip-Multiprocessors (CMP) or Multi-core processors**

# Recap: Moore's Law<sup>(1)</sup>

- 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



(1) Moore, G. E. (1965), 'Cramming more components onto integrated circuits', Electronics 38 (8).

# 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. 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 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 12700K	425.8 million

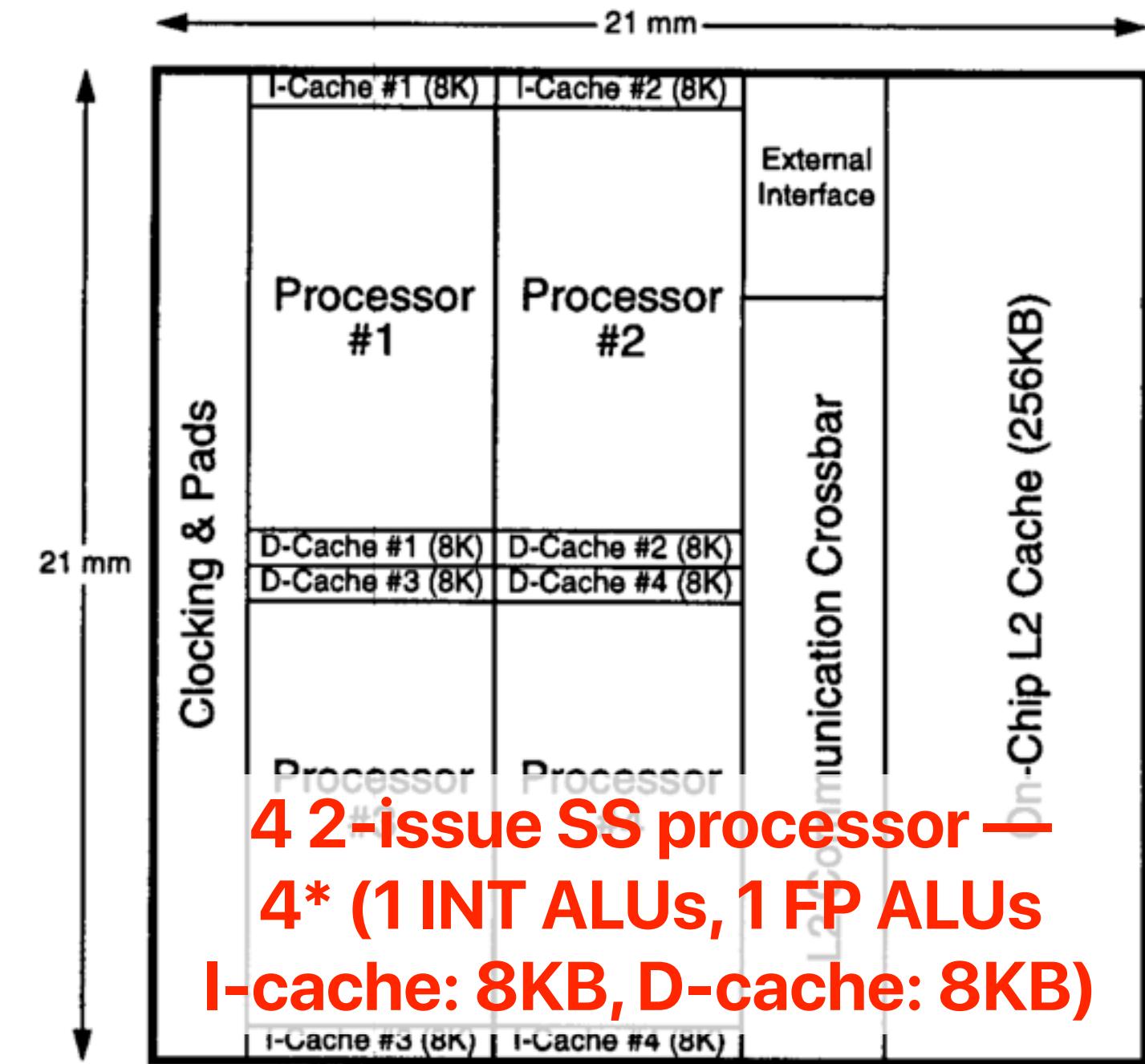
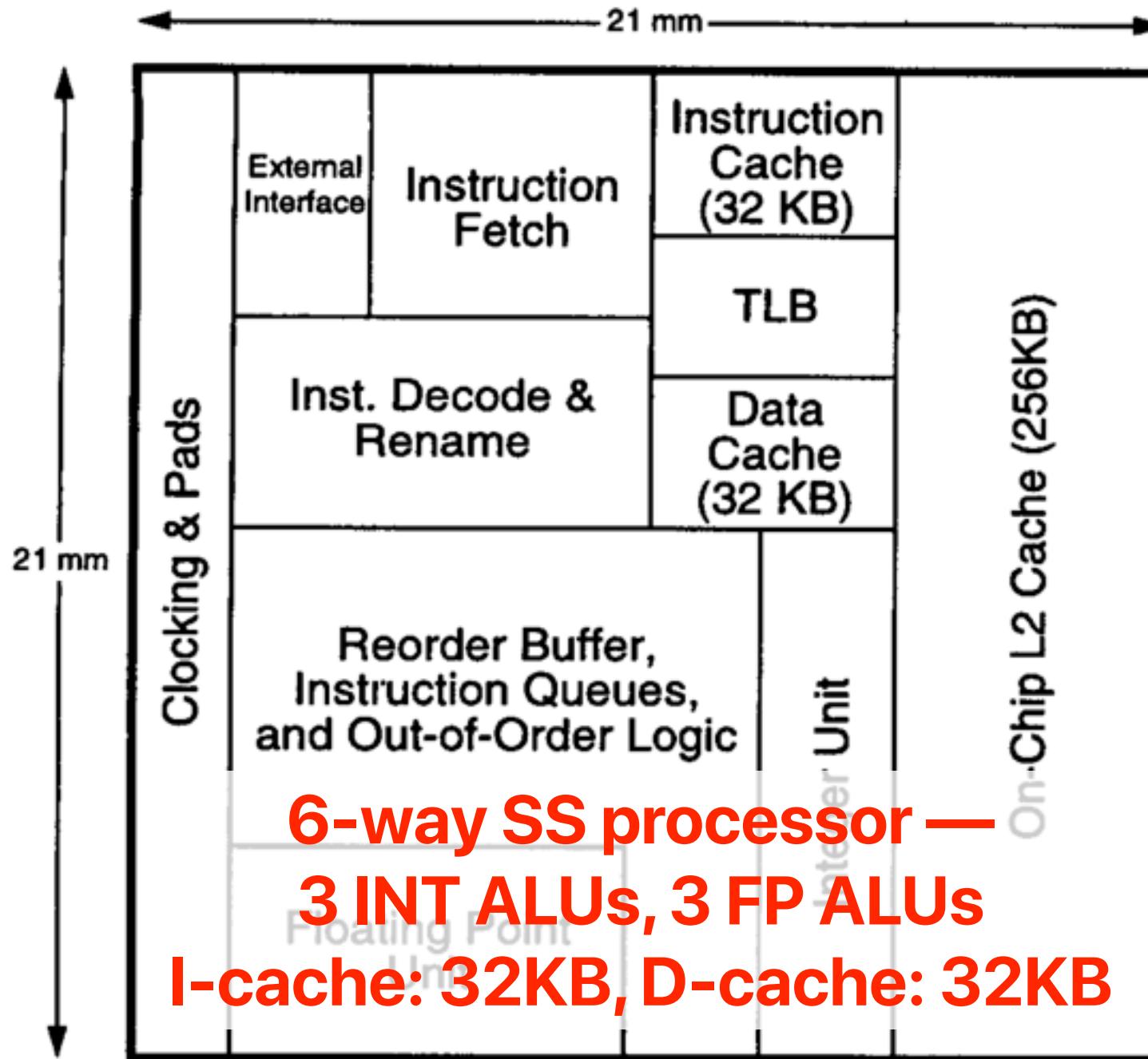
## 2x 3-issue ALUs Nehalem

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

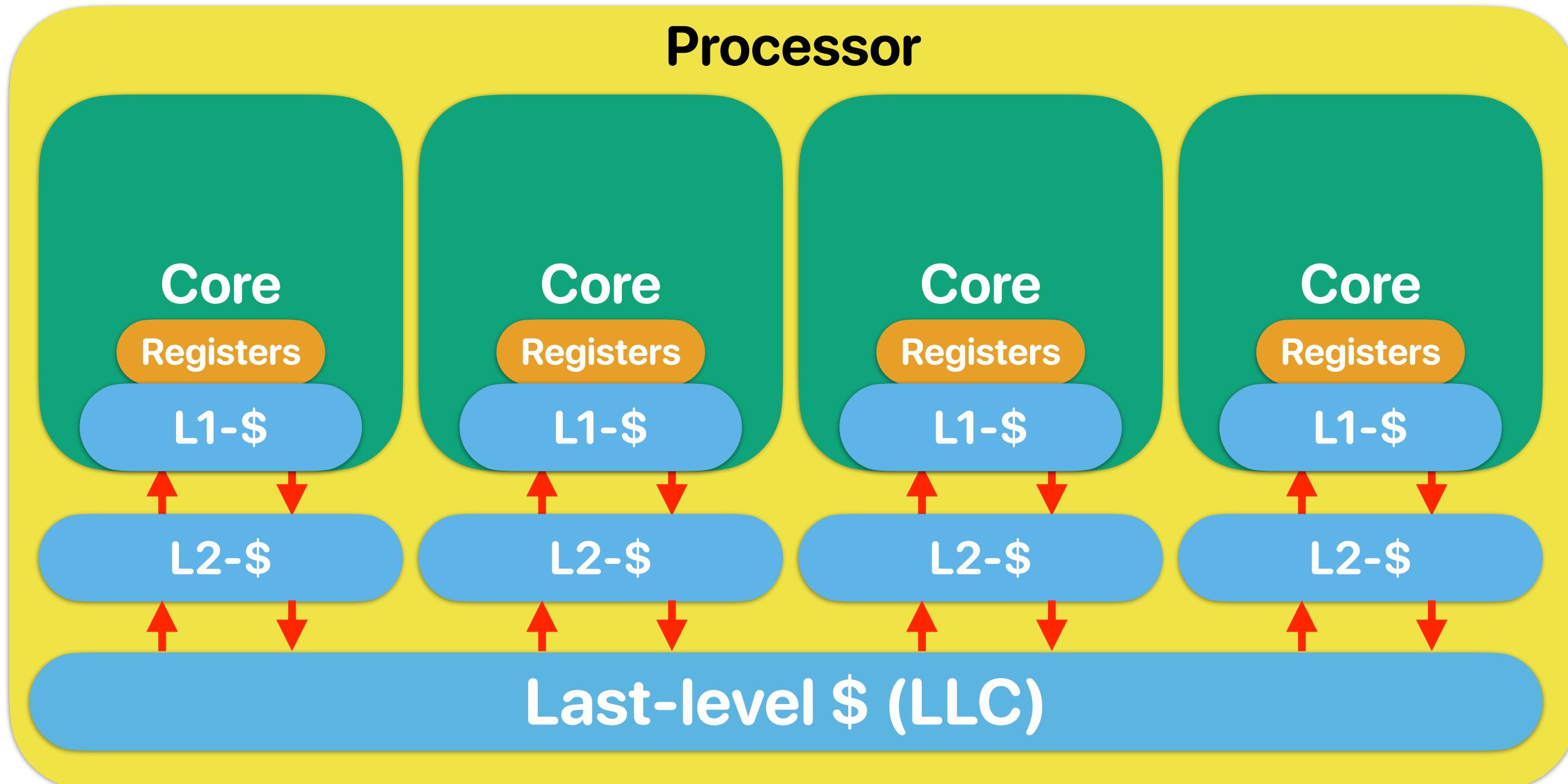
## 1x 5-issue ALUs Alder Lake

Based on [https://en.wikipedia.org/wiki/Transistor\\_count](https://en.wikipedia.org/wiki/Transistor_count)

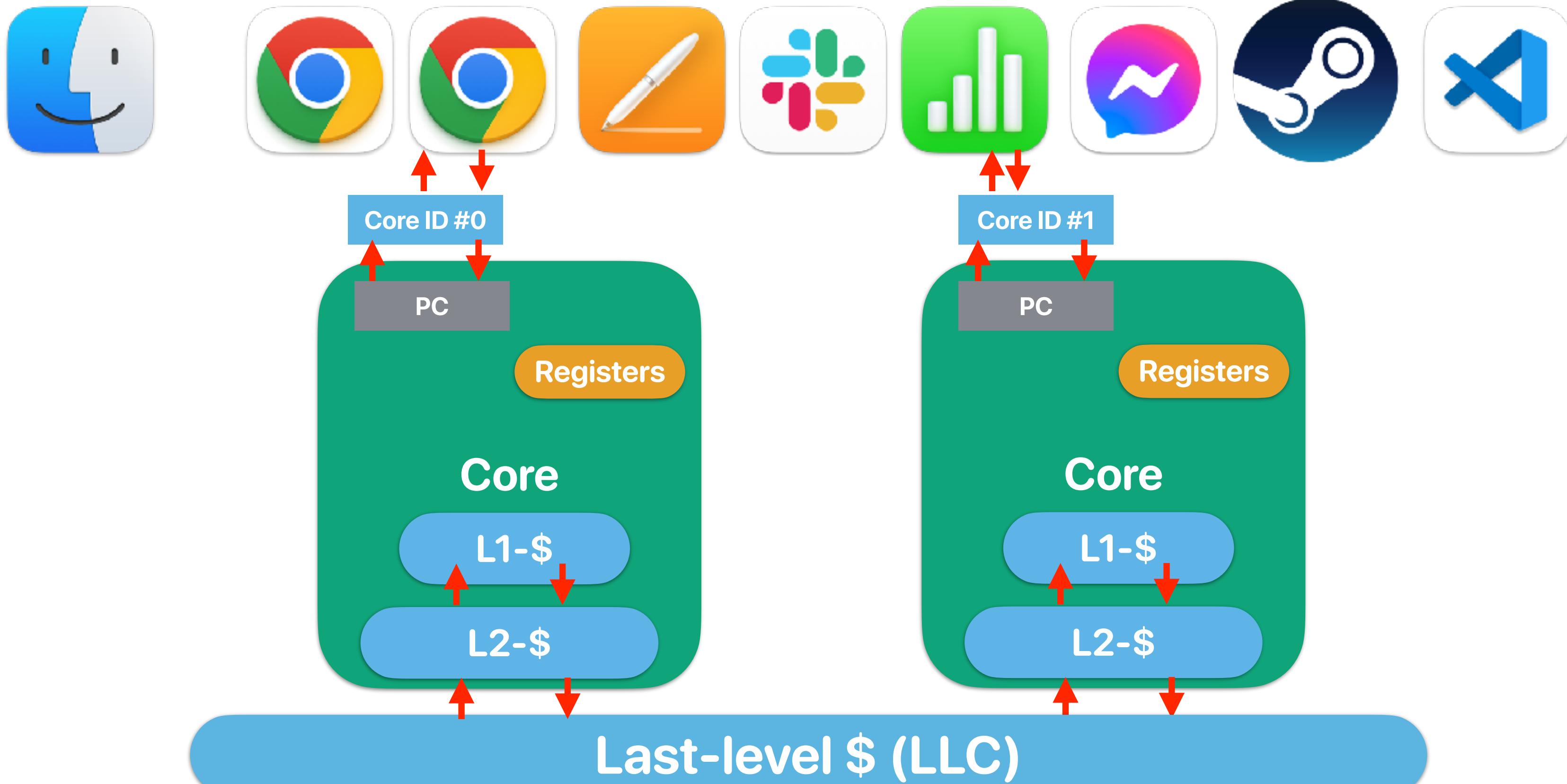
# Wide-issue SS processor v.s. multiple narrower-issue SS processors



# Concept of CMP



# CMP from the user/OS' perspective

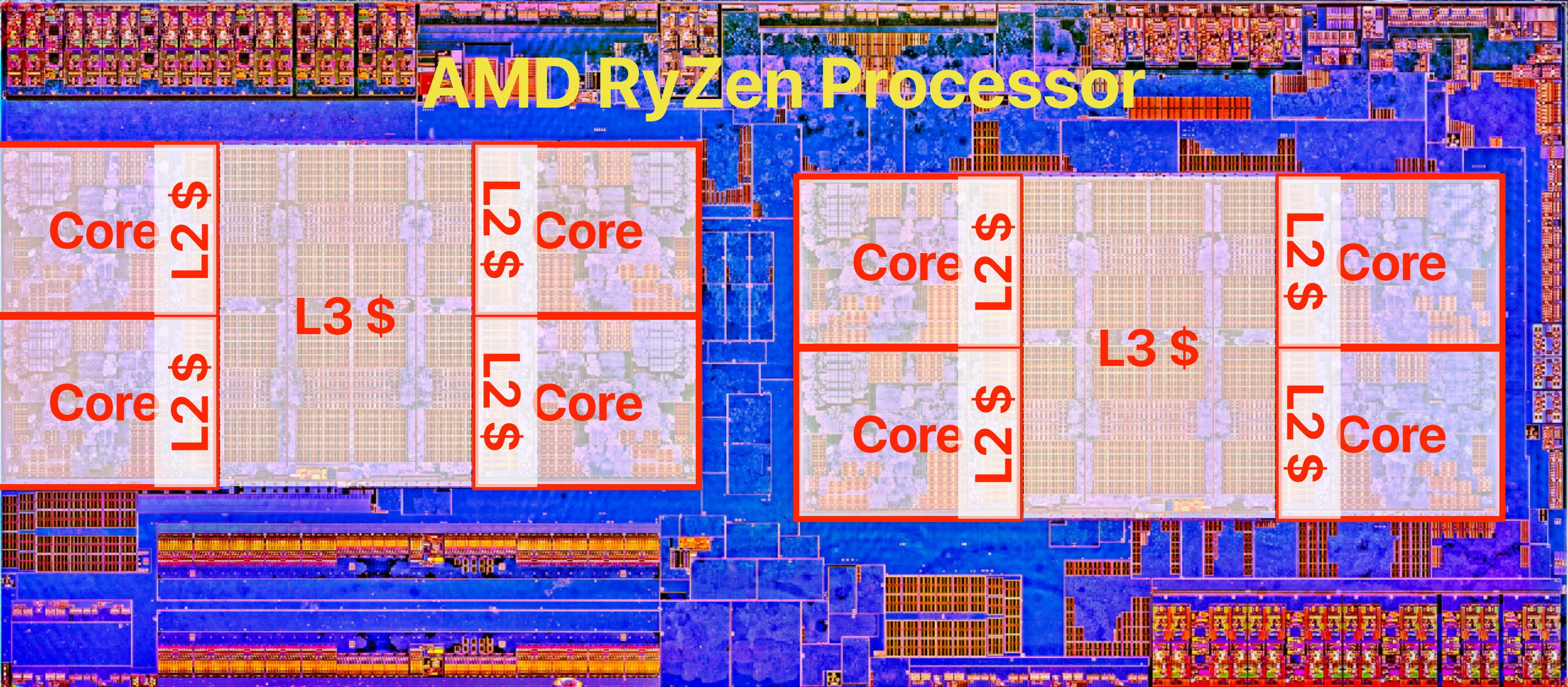


Architecture:	x86_64
CPU op-mode(s):	32-bit, 64-bit
Byte Order:	Little Endian
Address sizes:	39 bits physical, 48 bits virtual
CPU(s):	8
On-line CPU(s) list:	0-7
Thread(s) per core:	2
Core(s) per socket:	4
Socket(s):	1
NUMA node(s):	1
Vendor ID:	GenuineIntel
CPU family:	6
Model:	151
Model name:	12th 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

# Modern processors have both CMP/SMT



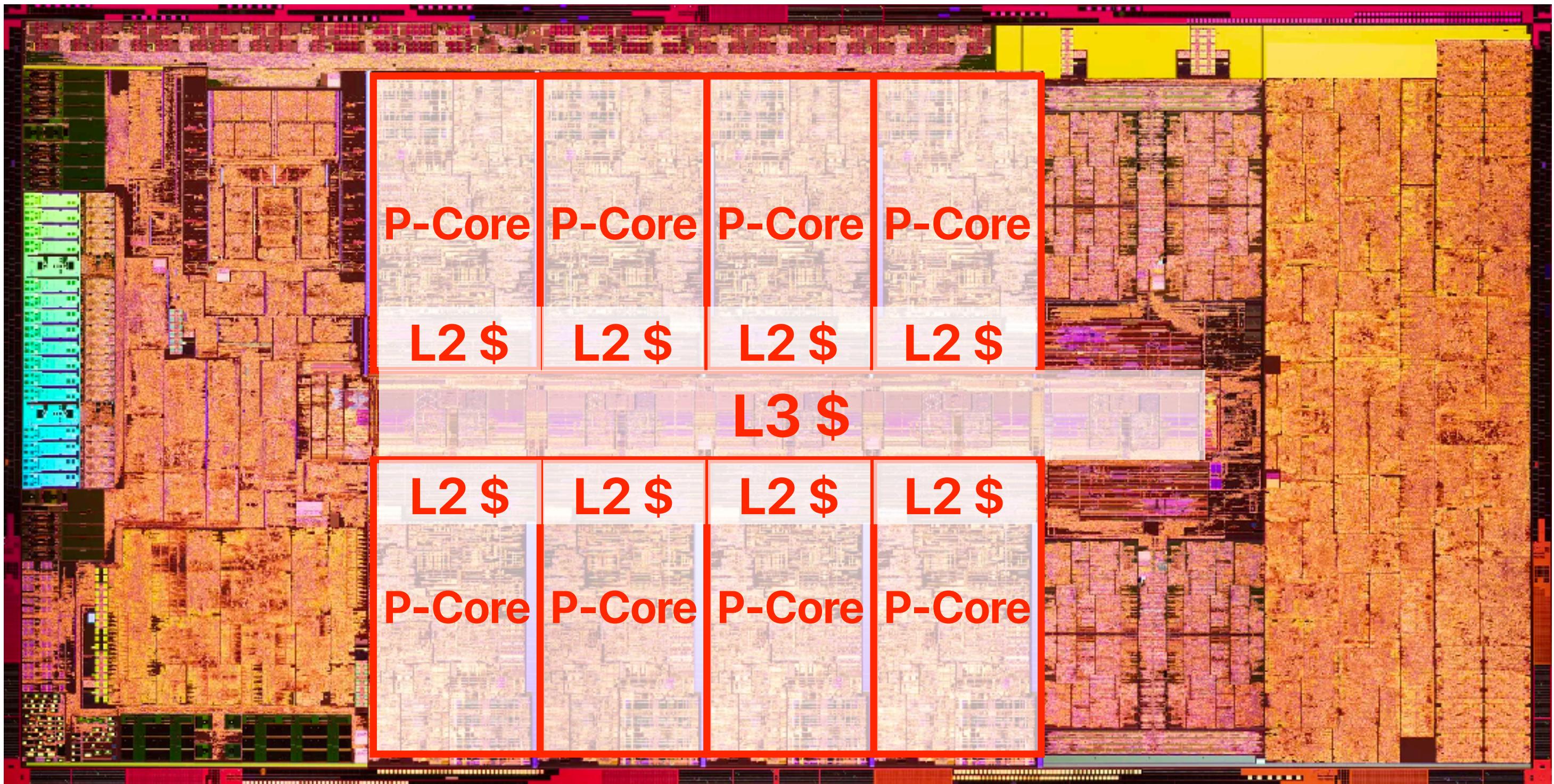
# AMD Ryzen Processor



AMD

RYZEN

# Intel Alder Lake



# Announcements

- Assignment #4 — due this Saturday
- Last reading quiz — due next Tuesday

# Computer Science & Engineering

142

つづく

