

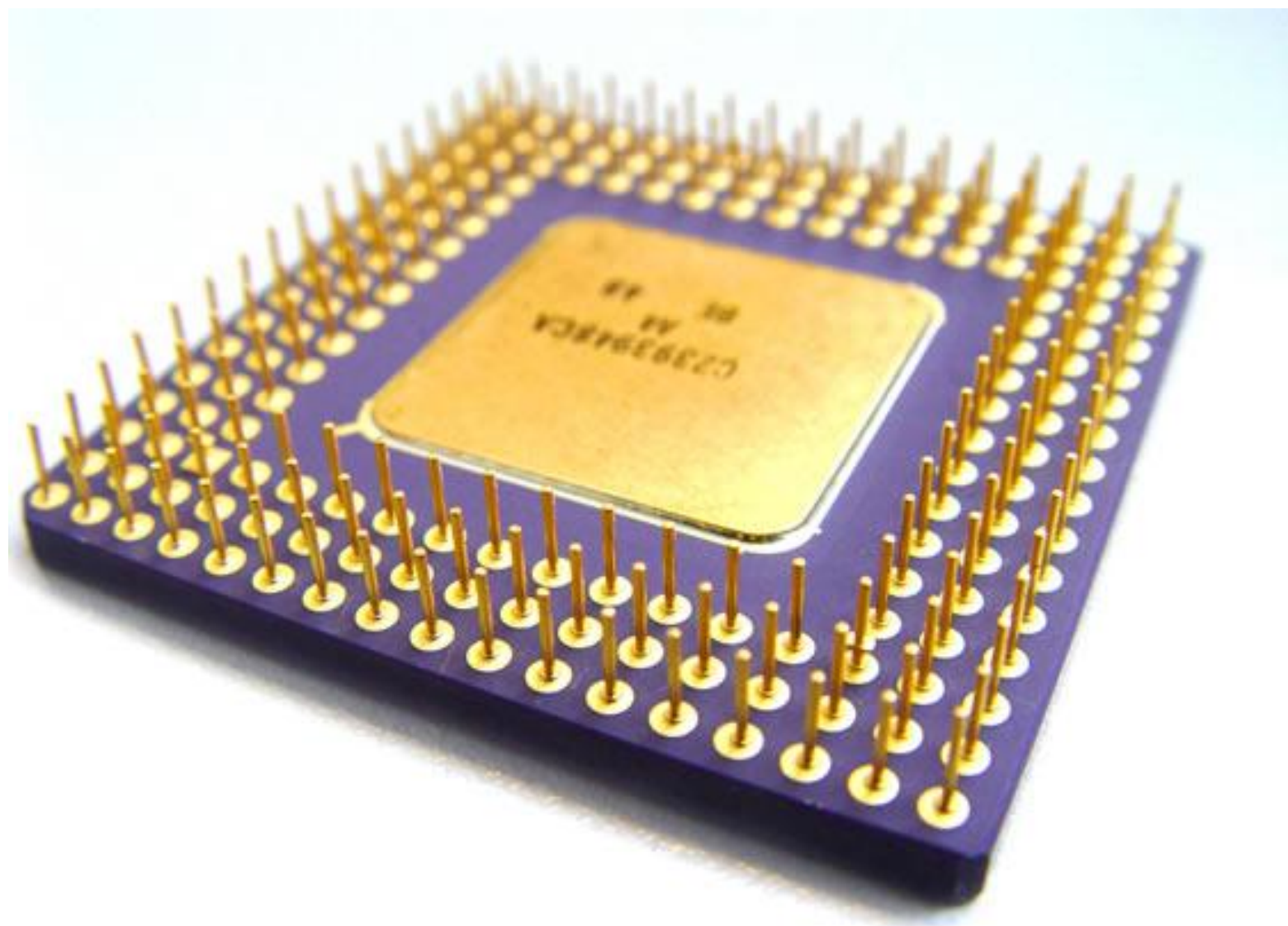
# CPU架构浅析

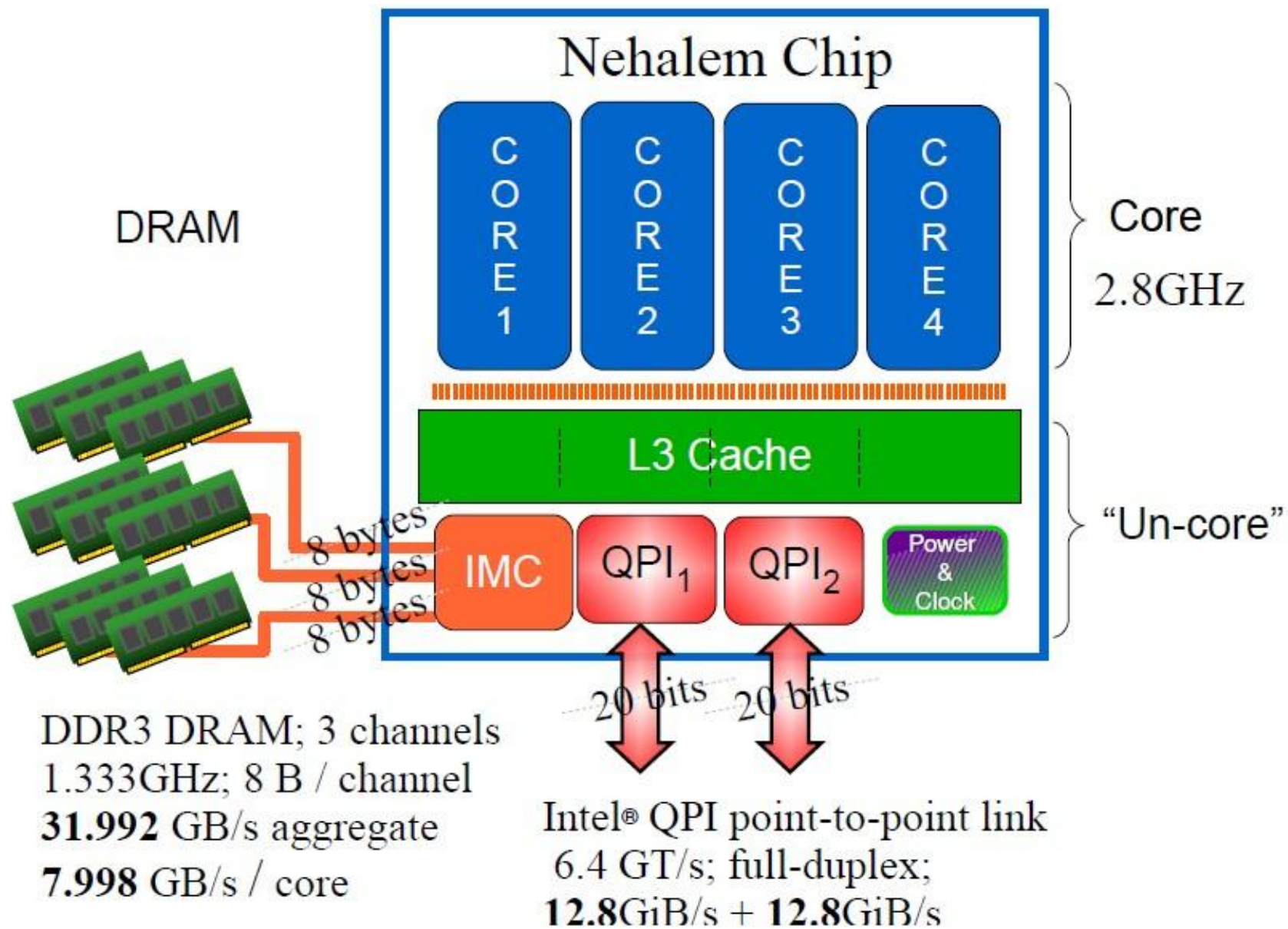
阿里数据库技术 —— 何登成

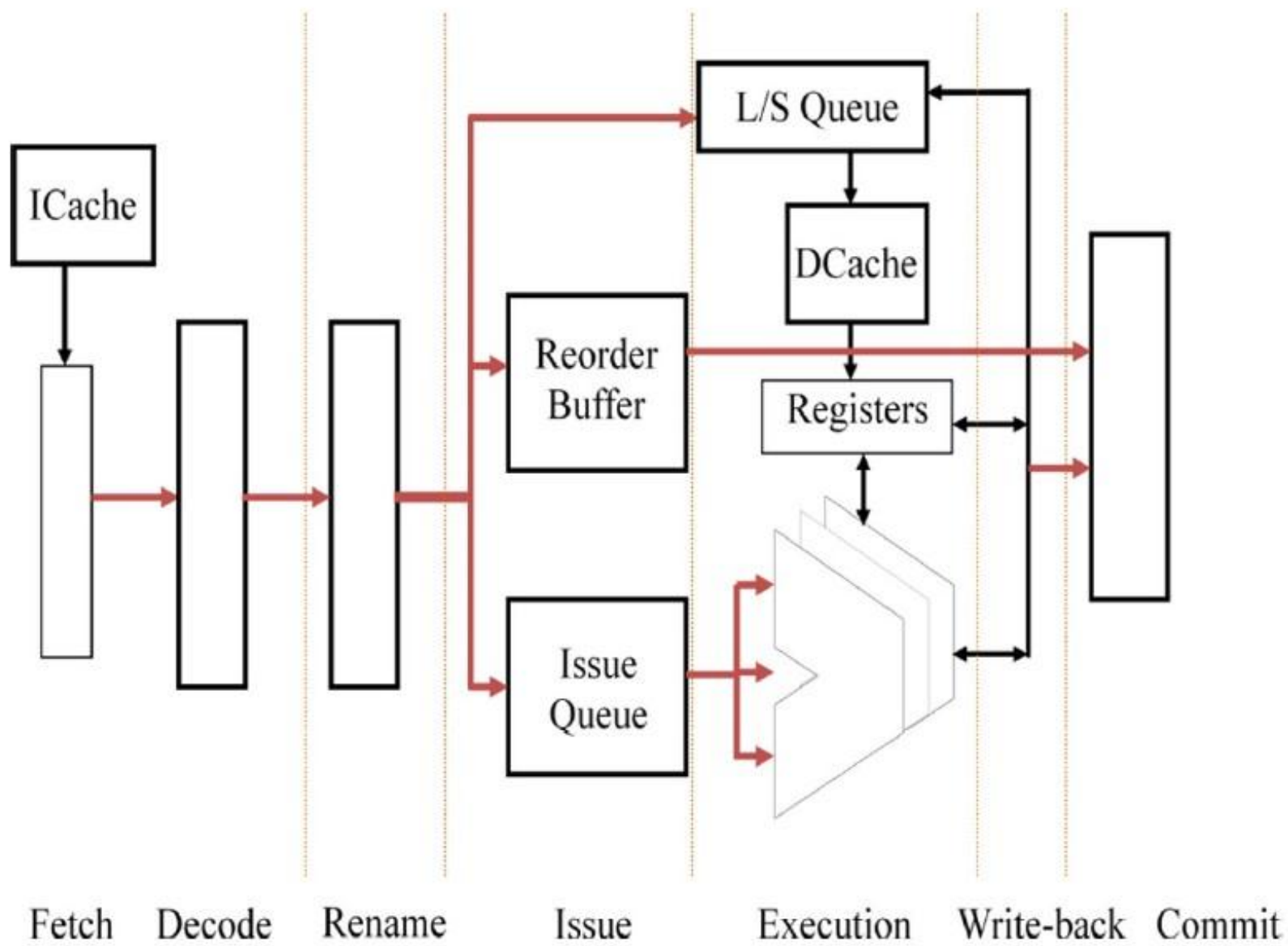
微博：何\_登成

# 分享大纲

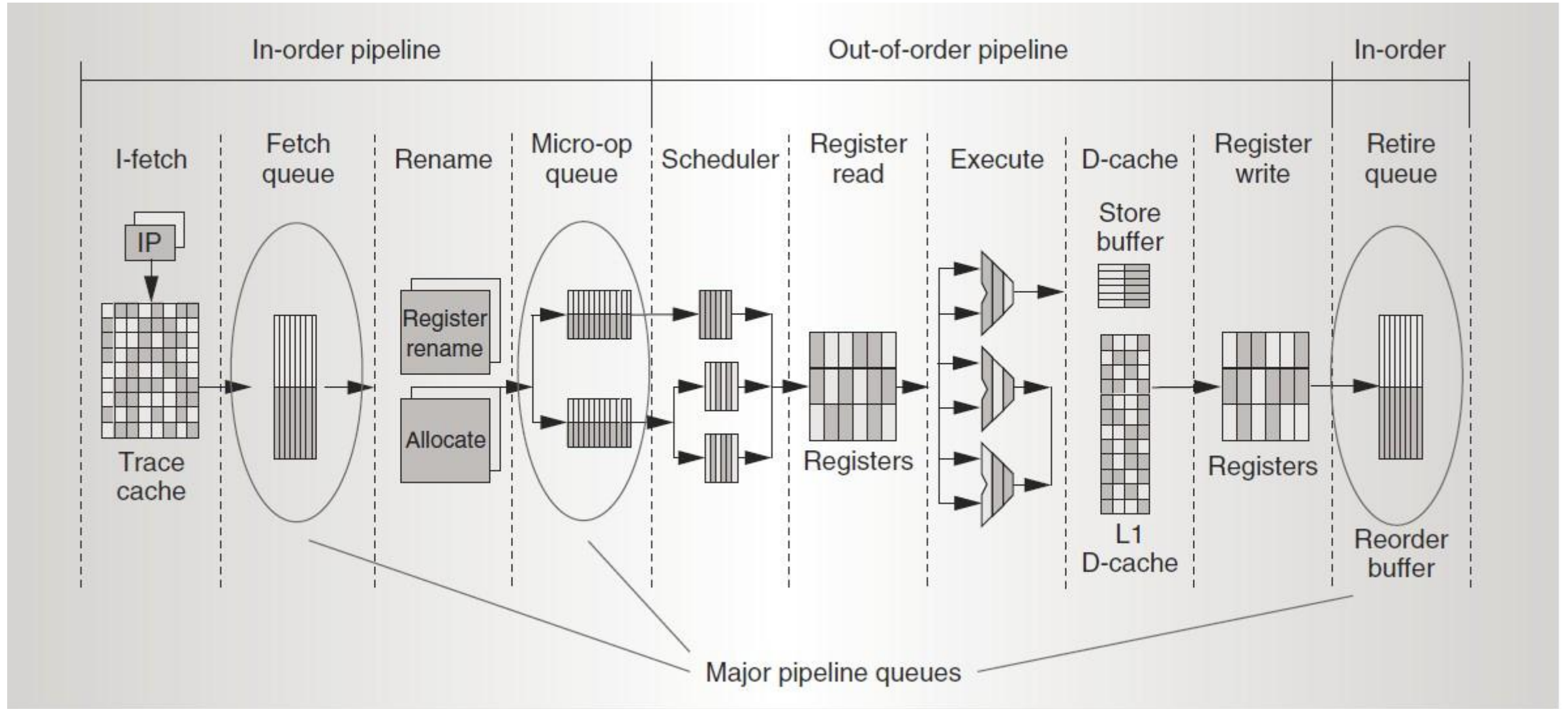
- 先介绍一个CPU的架构
  - Fetch to retire
- 通过CPU架构，理解一些技术点
  - Super pipelining
  - Superscalar
  - Simultaneous Multi-Threading (SMT)
  - Out-of-Order Execution
  - Performance vs Clock Cycle
- 再详细分析两个例子
  - Cache Coherence (缓存一致性)
  - Memory Consistency (Memory Model)
    - Intel, ARM



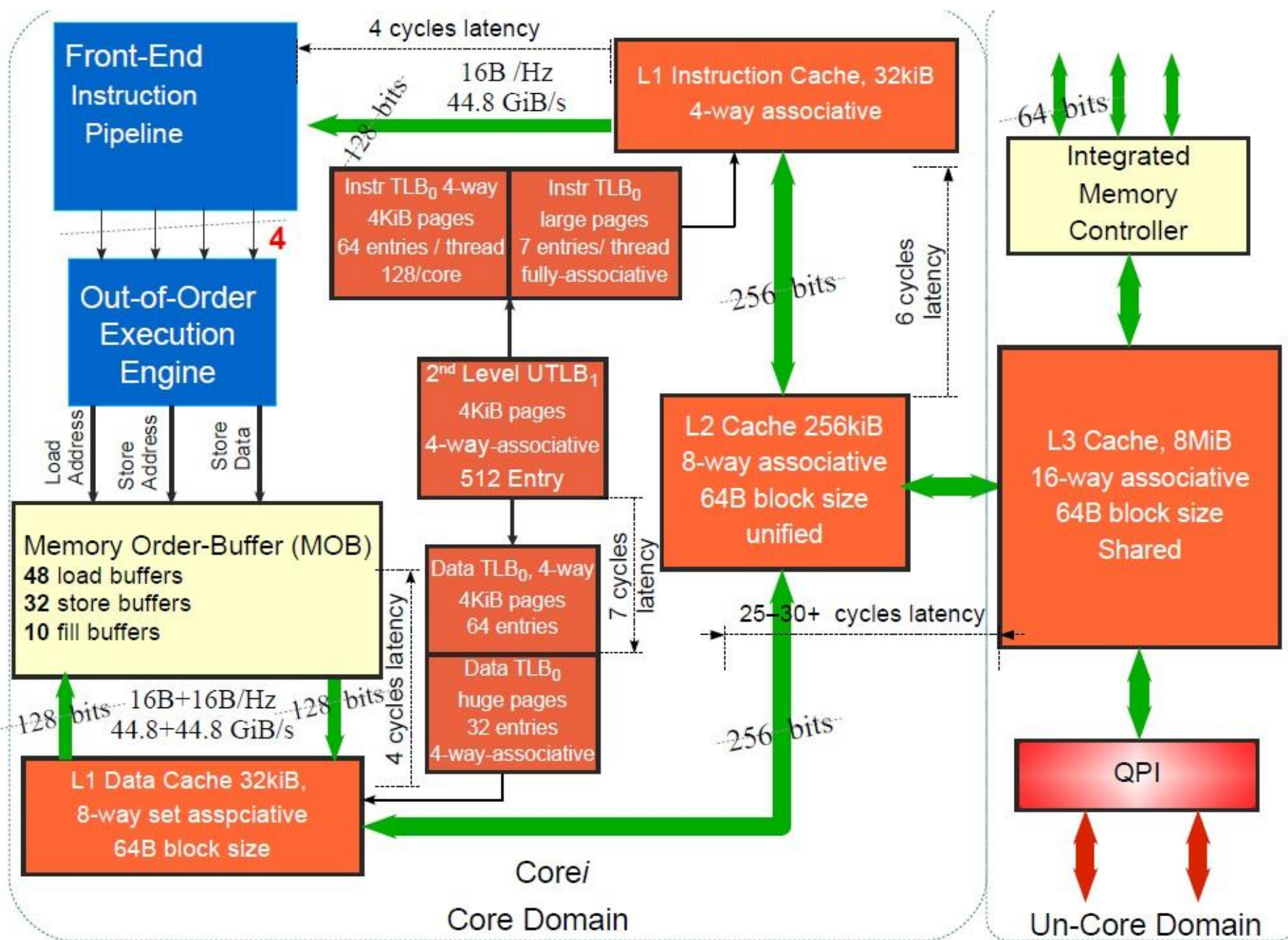




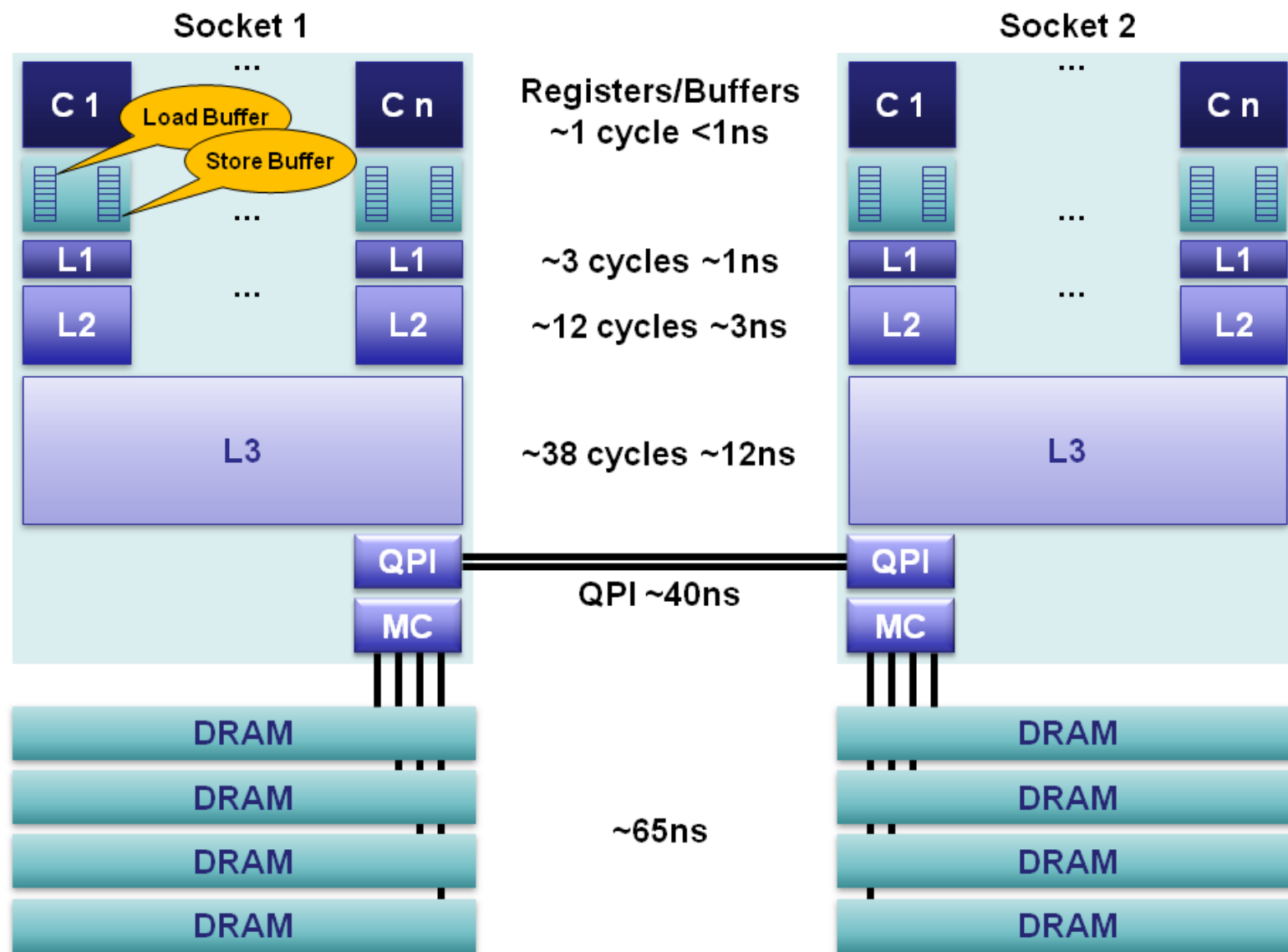
- The early stages of the processor **fetch-in** several **macro-instructions** at a time (say in a cache block) and
- **decode** them (break them down) into sequences of **micro-ops**.
- The micro-ops are **buffered** at various places where they can be picked up and scheduled to use the **FUs** in parallel if data dependencies are not violated. In Nehalem, micro-ops are issued to stations where they reserve their position for subsequent,
- **dispatching** as soon as **their input operands become available**.
- Finally, completed micro-ops **retire** and post their results to permanent storage.











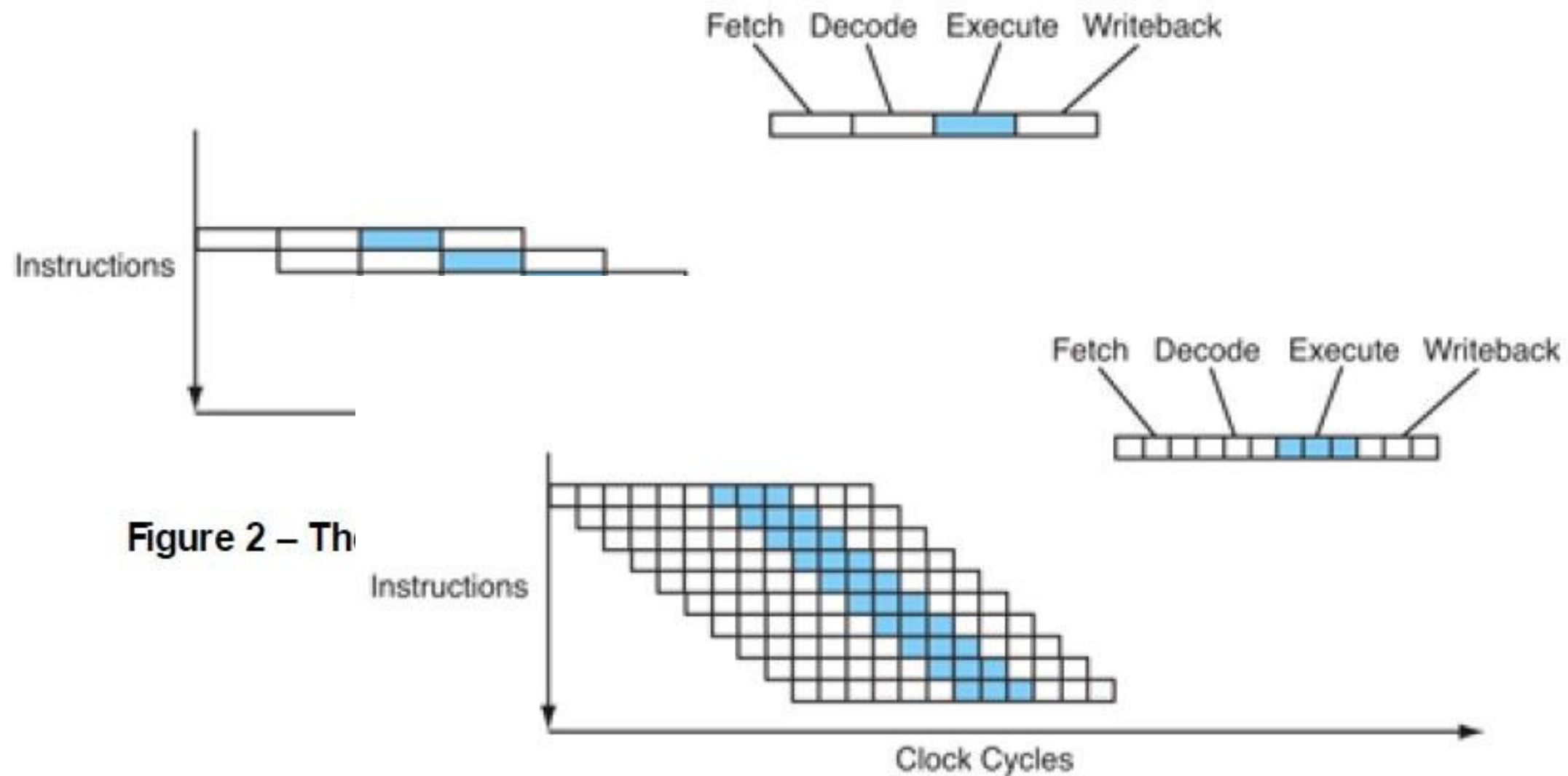


Figure 2 – Th

Figure 6 – The instruction flow of a superpipelined processor.

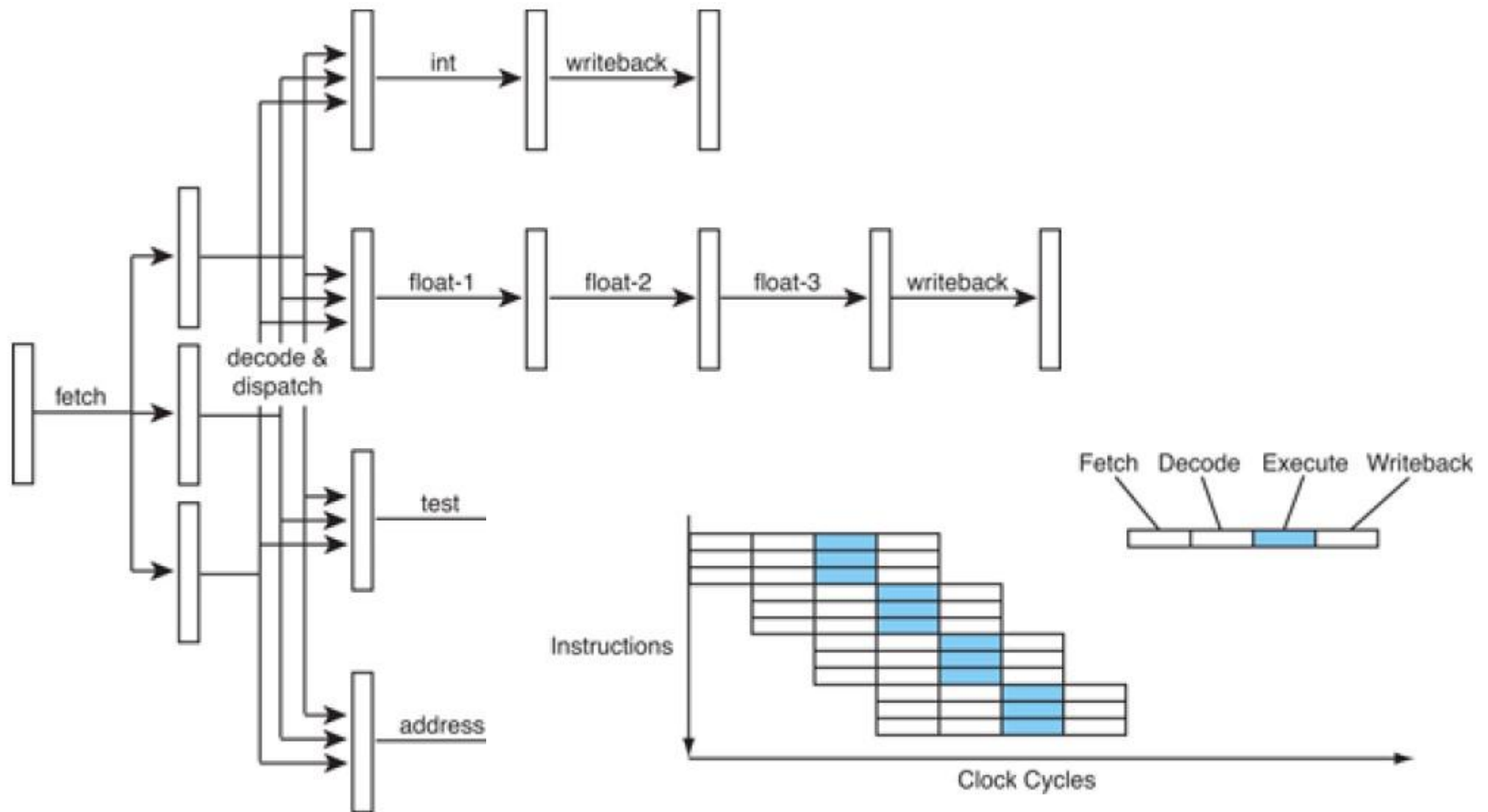


Figure 7 – A super

Figure 8 – The instruction flow of a superscalar processor.

- ILP
  - Instruction Level Parallelism
- CPI/IPC
  - Cycles per Instruction
  - Instructions per Cycle

```
# perf stat gzip file1

Performance counter stats for 'gzip file1':

    1920.159821 task-clock           #    0.991 CPUs utilized
           13 context-switches      #    0.007 K/sec
            0 CPU-migrations         #    0.000 K/sec
          258 page-faults            #    0.134 K/sec
 5,649,595,479 cycles                #    2.942 GHz                    [83.43%]
 1,808,339,931 stalled-cycles-frontend # 32.01% frontend cycles idle   [83.54%]
 1,171,884,577 stalled-cycles-backend # 20.74% backend cycles idle   [66.77%]
 8,625,207,199 instructions          #    1.53 insns per cycle
                                   #    0.21 stalled cycles per insn [83.51%]
 1,488,797,176 branches              # 775.351 M/sec                  [82.58%]
   53,395,139 branch-misses          #   3.59% of all branches       [83.78%]

    1.936842598 seconds time elapsed
```

```
$perf stat -p 14727
^C
Performance counter stats for process id '14727':

 37.366764 task-clock           #    0.006 CPUs utilized
       471 context-switches      #    0.013 M/sec
         4 CPU-migrations         #    0.000 M/sec
         0 page-faults           #    0.000 M/sec
57,734,472 cycles                #    1.545 GHz                    [78.03%]
44,125,966 stalled-cycles-frontend # 76.43% frontend cycles idle   [81.73%]
30,301,007 stalled-cycles-backend  # 52.48% backend cycles idle    [78.16%]
33,446,945 instructions          #    0.58 insns per cycle
                                   #    1.32 stalled cycles per insn [90.50%]
   6,291,083 branches            # 168.360 M/sec                  [90.70%]
    82,323 branch-misses         #   1.31% of all branches       [74.49%]

 6.380130669 seconds time elapsed
```

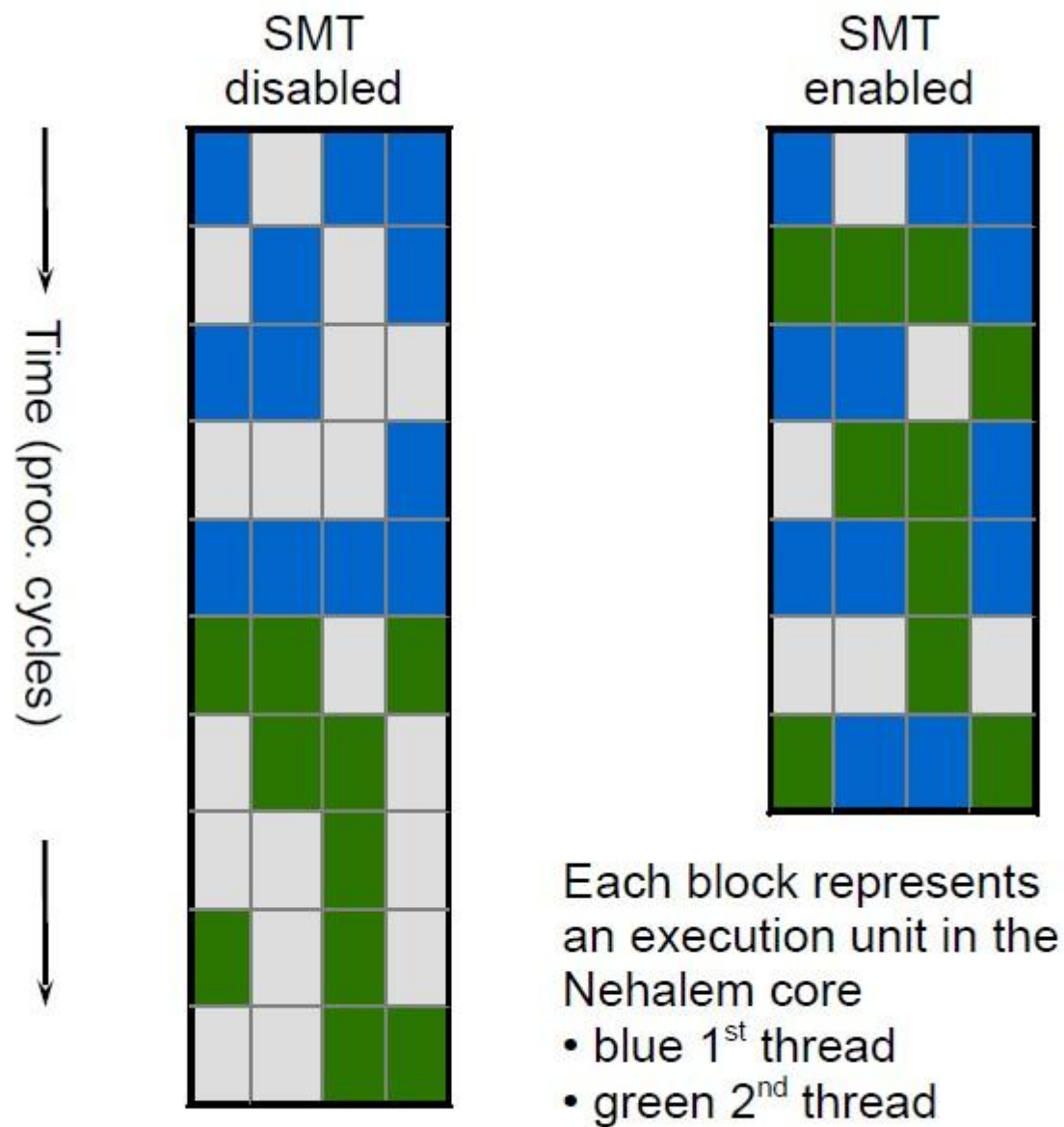
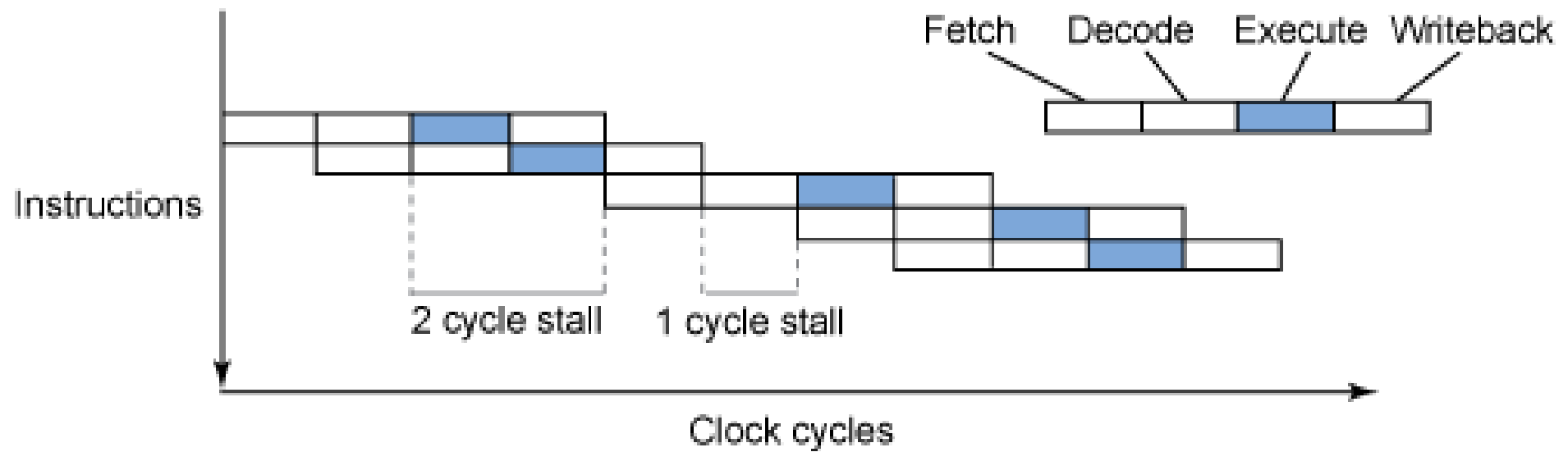


Figure 9: Simultaneous Multi-Treading (SMT) concept on Nehalem cores.



```
a = b * c;  
d = a + 1;
```



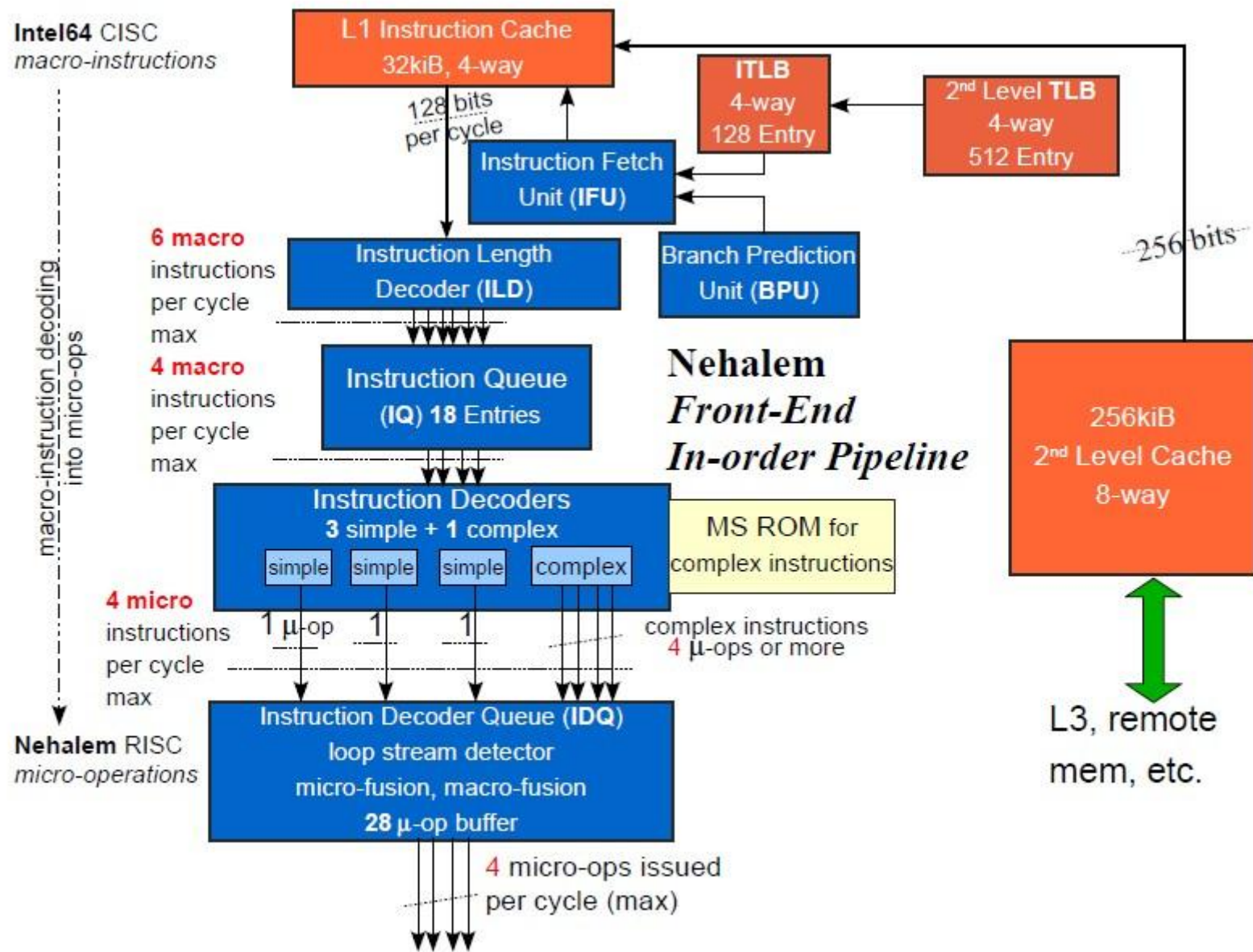


Figure 4: High-level diagram of the In-Order Front-End Nehalem Pipeline (FEP).

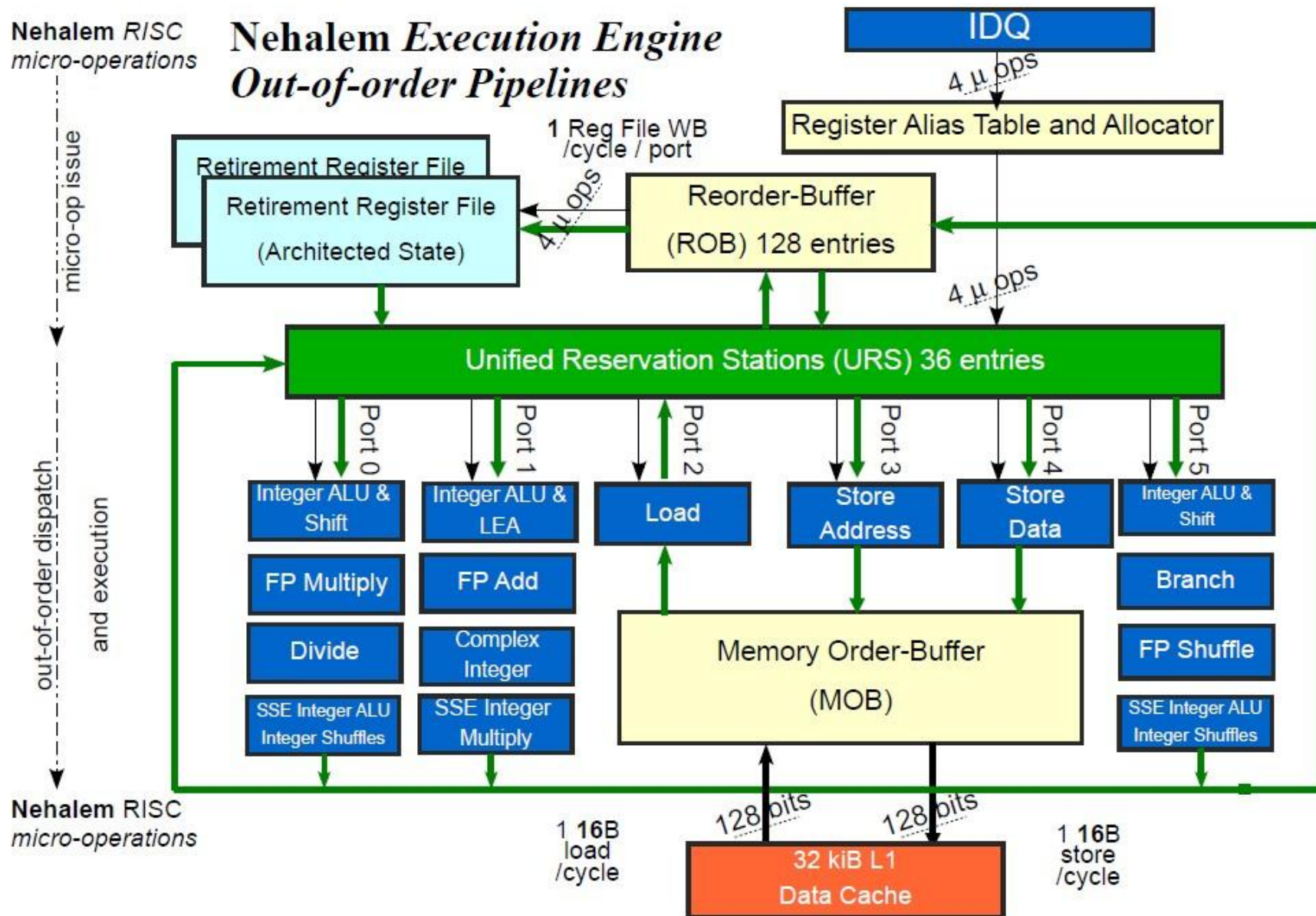
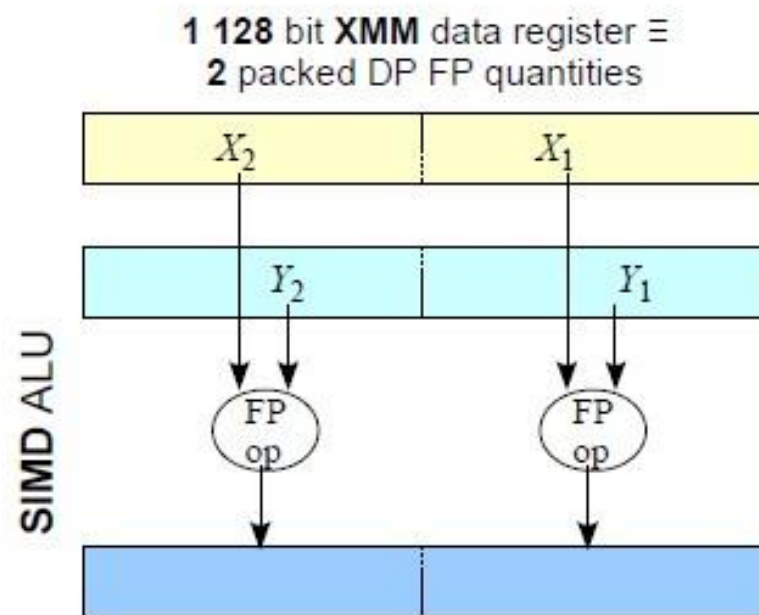
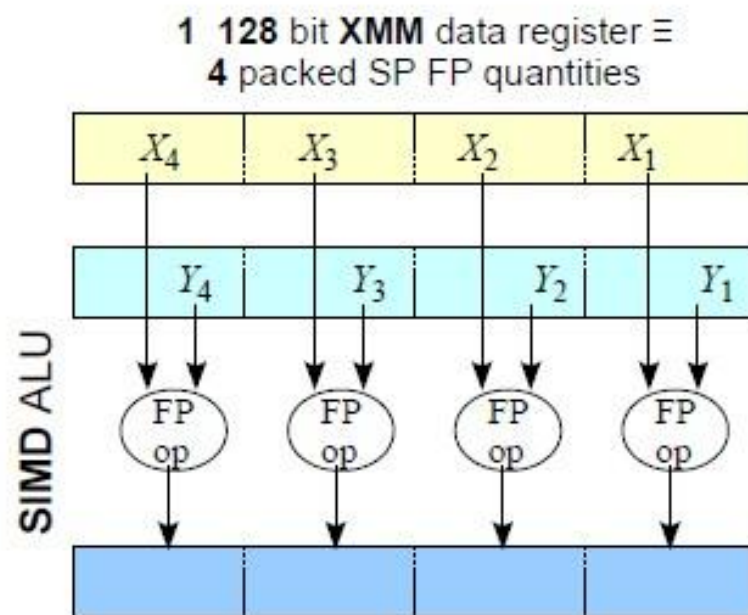


Figure 5: High-level diagram of a the out-of-order execution engine in the Nehalem core. All units are fully pipelined and can operate independently.



2 64-bit **Double**-Precision  
**SIMD** FP operations  
with **XMM** data registers

$\leftrightarrow$   
2 DP FP ops / cycle / port



4 32-bit **Single**-Precision  
**SIMD** FP operations  
with **XMM** data registers

$\leftrightarrow$   
4 SP FP ops / cycle / port

Figure 7: Floating-Point SIMD Operations in Nehalem.



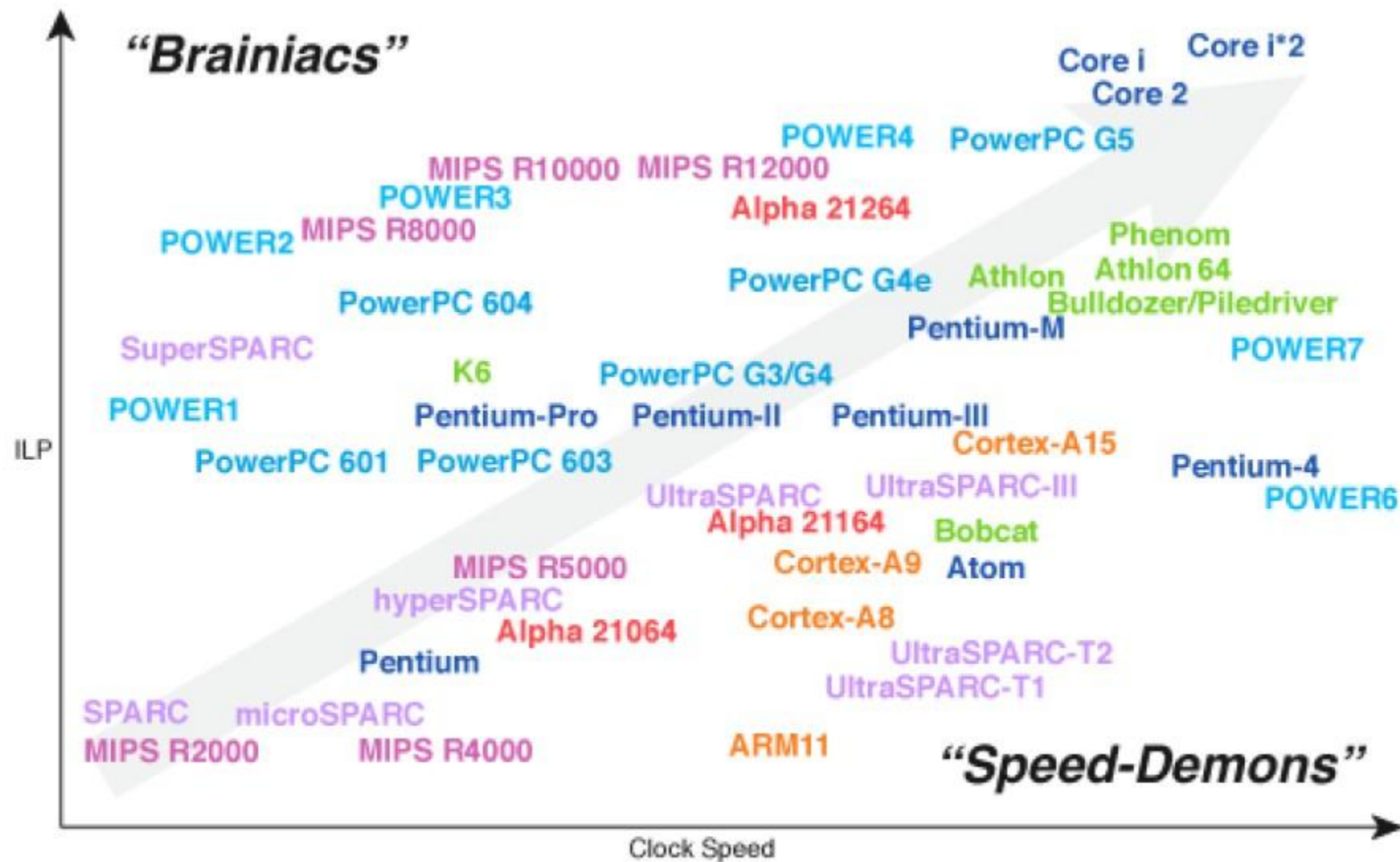


Figure 11 – Brainiacs vs speed-demons.

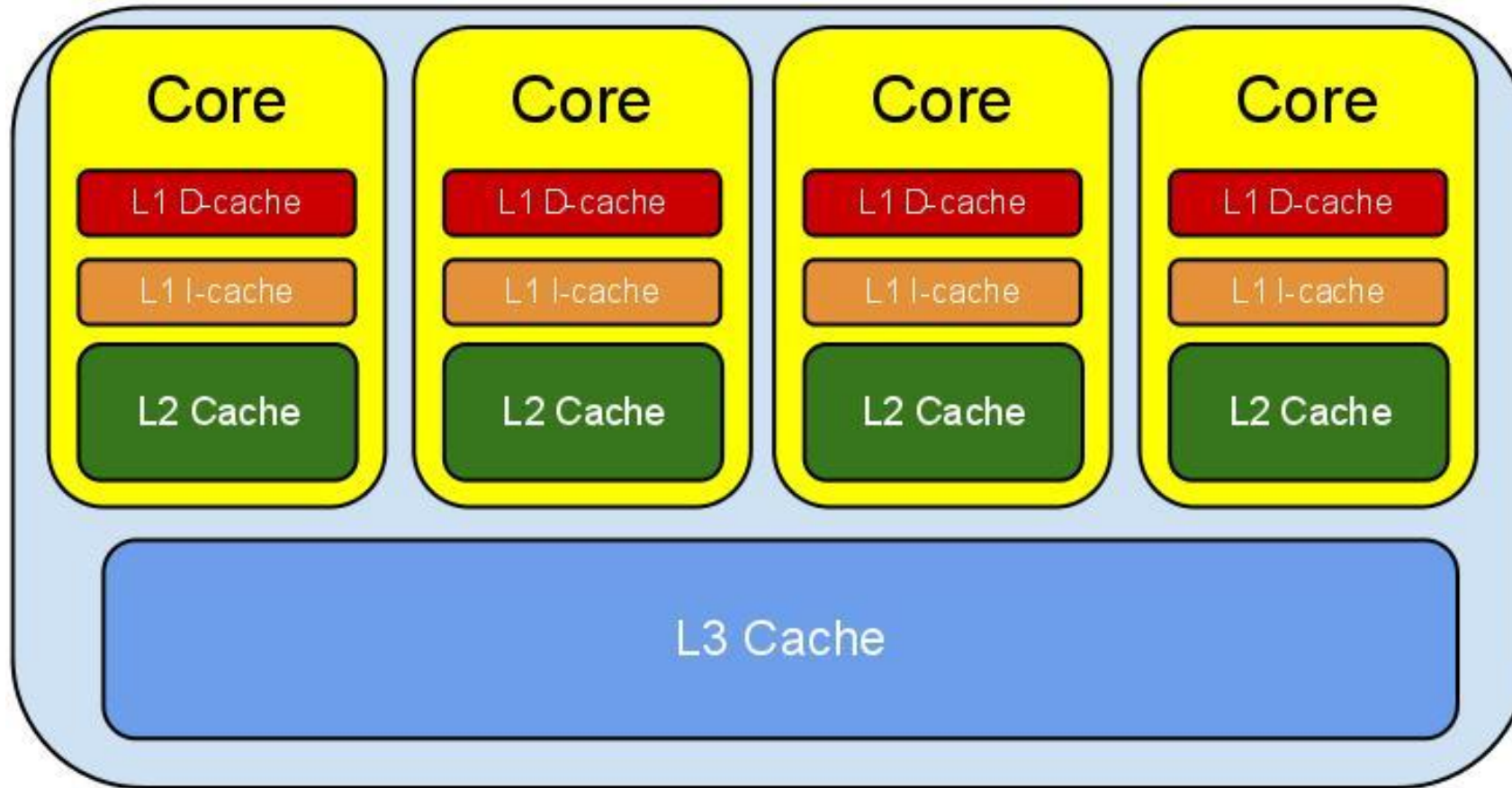


		<i>SPECint95</i>	<i>SPECfp95</i>
195 MHz	MIPS R10000	11.0	17.0
400 MHz	Alpha 21164	12.3	17.2
300 MHz	UltraSPARC	12.1	15.5
300 MHz	Pentium-II	11.6	8.8
300 MHz	PowerPC G3	14.8	11.4
135 MHz	POWER2	6.2	17.6

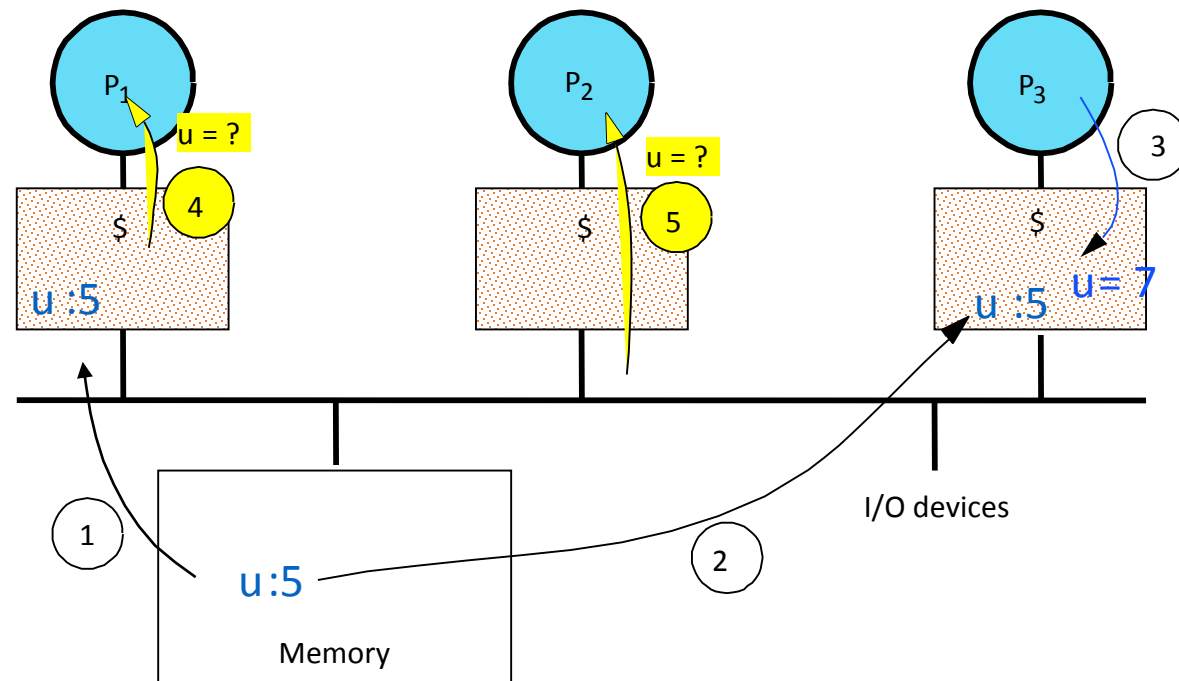
# Consistency vs Coherence

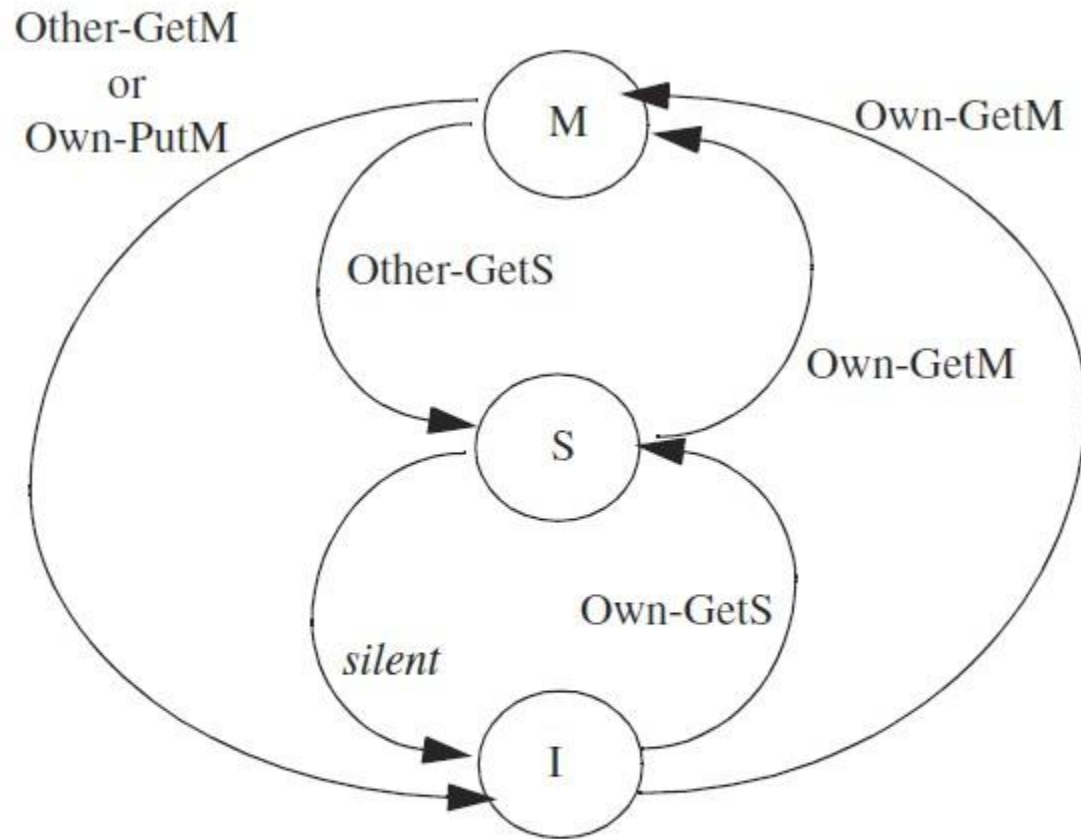
- **Consistency** (Memory Consistency, Memory Model: 内存模型)
  - Consistency definitions provide rules about **loads and stores** (or memory reads and writes) and how they act upon memory.
  - 可感知，对程序正确性有影响
- **Coherence** (Cache Coherence: 缓存一致性)
  - Coherence seeks to **make the caches of a shared-memory system as functionally invisible** as the caches in a single-core system. Correct coherence ensures that a programmer cannot determine whether and where a system has caches by analyzing the results of loads and stores.
  - 不可感知，CPU内部搞定

# CPU Cache

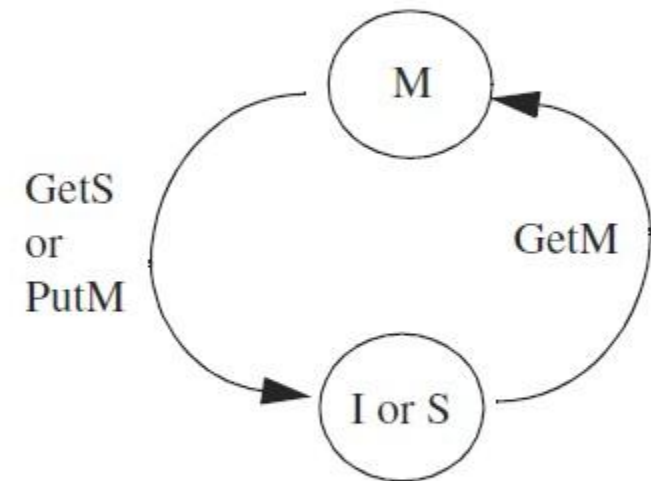


# Cache Coherence Problem





**FIGURE 7.1:** MSI: Transitions between stable states at cache controller.



**FIGURE 7.2:** MSI: Transitions between stable states at memory controller.



**TABLE 7.5:** Simple Snooping (Atomic Requests, Atomic Transactions): Cache Controller

States	Processor Core Events			Bus Events						
				Own Transaction				Transactions For Other Cores		
	Load	Store	Replacement	Own-GetS	Own-GetM	Own-PutM	Data	Other-GetS	Other-GetM	Other-PutM
I	issue GetS /IS <sup>D</sup>	issue GetM /IM <sup>D</sup>								
IS <sup>D</sup>	stall Load	stall Store	stall Evict				copy data into cache, load hit /S	(A)	(A)	(A)
IM <sup>D</sup>	stall Load	stall Store	stall Evict				copy data into cache, store hit /M	(A)	(A)	(A)
S	load hit	issue GetM /SM <sup>D</sup>	-/I						-/I	
SM <sup>D</sup>	load hit	stall Store	stall Evict				copy data into cache, store hit /M	(A)	(A)	(A)
M	load hit	store hit	issue PutM, send Data to memory /I					send Data to req and memory /S	send Data to req /I	

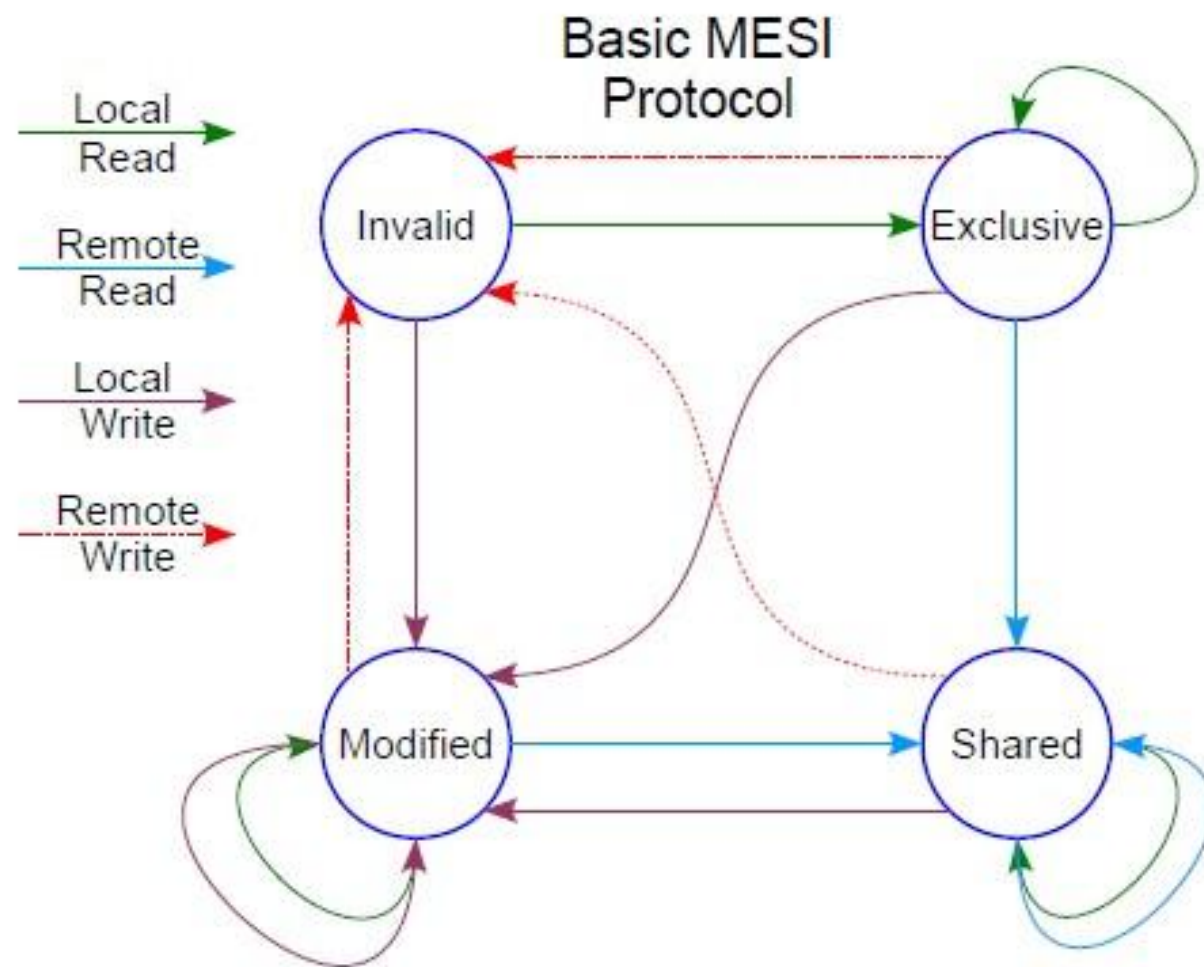


Figure 12: The State Transitions of Cache Blocks in the Basic MESI Cache-Coherence Protocol.

# TestCase

- [ccBench](#)

Supported events:

- 0 - STORE\_ON\_MODIFIED
- 1 - STORE\_ON\_MODIFIED\_NO\_SYNC
- 2 - STORE\_ON\_EXCLUSIVE
- 3 - STORE\_ON\_SHARED
- 4 - STORE\_ON\_OWNED\_MINE
- 5 - STORE\_ON\_OWNED
- 6 - STORE\_ON\_INVALID
- 7 - LOAD\_FROM\_MODIFIED
- 8 - LOAD\_FROM\_EXCLUSIVE
- 9 - LOAD\_FROM\_SHARED
- 10 - LOAD\_FROM\_OWNED
- 11 - LOAD\_FROM\_INVALID
- 12 - CAS
- 13 - FAI
- 14 - TAS
- 15 - SWAP
- 16 - CAS\_ON\_MODIFIED
- 17 - FAI\_ON\_MODIFIED
- 18 - TAS\_ON\_MODIFIED
- 19 - SWAP\_ON\_MODIFIED
- 20 - CAS\_ON\_SHARED
- 21 - FAI\_ON\_SHARED
- 22 - TAS\_ON\_SHARED
- 23 - SWAP\_ON\_SHARED
- 24 - CAS\_CONCURRENT
- 25 - FAI\_ON\_INVALID
- 26 - LOAD\_FROM\_L1
- 27 - LFENCE
- 28 - SFENCE
- 29 - MFENCE
- 30 - PROFILER
- 31 - PAUSE
- 32 - NOP

```
./ccbench -e 8 -x 0 -y 8 -t 0
test: STORE_ON_MODIFIED / #cores: 2 / #repetitions: 10000 / stride: 2048 (128 kiB) / fence: load/full
core1: 0 / core2: 8
[00] *** Core 0 ****
---- statistics:
[00] avg : 36.3 abs dev : 5.8 std dev : 9.5 num : 10000
[00] min : 24.0 (element: 10) max : 256.0 (element: 5618)
[00] 0-10% : 3517 ( 35.2% | avg: 36.0 | abs dev: 0.0 | std dev: 0.0 = 0.0%)
[00] 10-25% : 5471 ( 54.7% | avg: 33.0 | abs dev: 5.0 | std dev: 5.4 = 16.5%)
[00] 25-50% : 257 ( 2.6% | avg: 28.8 | abs dev: 7.7 | std dev: 9.8 = 34.2%)
[00] 50-75% : 159 ( 1.6% | avg: 59.8 | abs dev: 0.3 | std dev: 0.8 = 1.4%)
[00] 75-100% : 596 ( 6.0% | avg: 64.9 | abs dev: 1.7 | std dev: 9.8 = 15.2%)

[01] *** Core 1 ****
---- statistics:
[01] avg : 38.1 abs dev : 6.4 std dev : 9.7 num : 10000

./ccbench -e 8 -x 0 -y 1 -t 0
test: STORE_ON_MODIFIED / #cores: 2 / #repetitions: 10000 / stride: 2048 (128 kiB) / fence: load/full 4.8%)
[00] *** Core 0 **** 16.4%)
---- statistics:
[00] avg : 224.5 abs dev : 15.3 std dev : 39.3 num : 10000
[00] min : 0.0 (element: 4787) max : 1004.0 (element: 8309) 18.5%)
[00] 0-10% : 9466 ( 94.7% | avg: 232.1 | abs dev: 2.0 | std dev: 3.4 = 1.5%) 1.3%)
[00] 10-25% : 131 ( 1.3% | avg: 236.0 | abs dev: 22.5 | std dev: 25.4 = 10.8%) 15.5%)
[00] 25-50% : 0 ( 0.0% | avg: -nan | abs dev: -nan | std dev: -nan = -nan%)
[00] 50-75% : 4 ( 0.0% | avg: 139.0 | abs dev: 110.5 | std dev: 127.6 = 91.8%)
[00] 75-100% : 399 ( 4.0% | avg: 40.2 | abs dev: 6.7 | std dev: 51.6 = 128.3%)

[01] *** Core 1 ****
---- statistics:
[01] avg : 222.8 abs dev : 15.9 std dev : 39.3 num : 10000
[01] min : 0.0 (element: 5903) max : 392.0 (element: 5346)
[01] 0-10% : 9434 ( 94.3% | avg: 230.9 | abs dev: 2.9 | std dev: 3.7 = 1.6%)
[01] 10-25% : 143 ( 1.4% | avg: 232.9 | abs dev: 26.3 | std dev: 28.1 = 12.1%)
[01] 25-50% : 0 ( 0.0% | avg: -nan | abs dev: -nan | std dev: -nan = -nan%)
[01] 50-75% : 2 ( 0.0% | avg: 210.0 | abs dev: 146.0 | std dev: 146.0 = 69.5%)
[01] 75-100% : 421 ( 4.2% | avg: 37.8 | abs dev: 2.8 | std dev: 17.5 = 46.2%)
```

# Sequential Consistency

- 大家做一



- SC的缺点？（为什么会有乱序出现）

- 性能，性能，性能
- 乱序执行 → 消除Stall

- 有哪些常见的乱序情况？

- 两类内存操作：Load，Store
- 4种组合：
  - LoadLoad
  - LoadStore
  - StoreLoad
  - StoreStore



		Loads Reordered After Loads?	Loads Reordered After Stores?	Stores Reordered After Stores?	Stores Reordered After Loads?	Atomic Instructions Reordered With Loads?	Atomic Instructions Reordered With Stores?	Dependent Loads Reordered?	Incoherent Instruction Cache/Pipeline?
Loads reordered		Y	Y	Y	Y	Y	Y	Y	Y
Loads reordered		Y	Y	Y	Y	Y	Y	Y	Y
Stores reordered		Y	Y	Y	Y	Y	Y	Y	Y
Stores reordered	Alpha	Y	Y	Y	Y	Y	Y	Y	Y
Atomic reordered	AMD64	Y			Y				Y
Atomic reordered	IA64	Y	Y	Y	Y	Y	Y		Y
Dependent loads reordered	(PA-RISC)	Y	Y	Y	Y				
Incoherent instruction cache/pipeline	PA-RISC CPUs								Y
	POWER	Y	Y	Y	Y	Y	Y		Y
	SPARC RMO	Y	Y	Y	Y	Y	Y		Y
	(SPARC PSO)			Y	Y		Y		Y
	SPARC TSO				Y				Y
	x86	Y	Y		Y				Y
	(x86 OOSTore)	Y	Y	Y	Y				Y
	zSeries				Y				Y

memory ordering in  
ARMv7 PA-RISC POWER

	Loads Reordered After Loads?	Loads Reordered After Stores?	Stores Reordered After Stores?	Stores Reordered After Loads?	Atomic Instructions Reordered With Loads?	Atomic Instructions Reordered With Stores?	Dependent Loads Reordered?	Incoherent Instruction Cache/Pipeline?
Alpha	Y	Y	Y	Y	Y	Y	Y	Y
AMD64				Y				
ARMv7-A/R	Y	Y	Y	Y	Y	Y		Y
IA64	Y	Y	Y	Y	Y	Y		Y
(PA-RISC)	Y	Y	Y	Y				
PA-RISC CPUs								
POWER™	Y	Y	Y	Y	Y	Y		Y
(SPARC RMO)	Y	Y	Y	Y	Y	Y		Y
(SPARC PSO)			Y	Y		Y		Y
SPARC TSO				Y				Y
x86				Y				Y
(x86 OOSTore)	Y	Y	Y	Y				Y
zSeries®				Y				Y

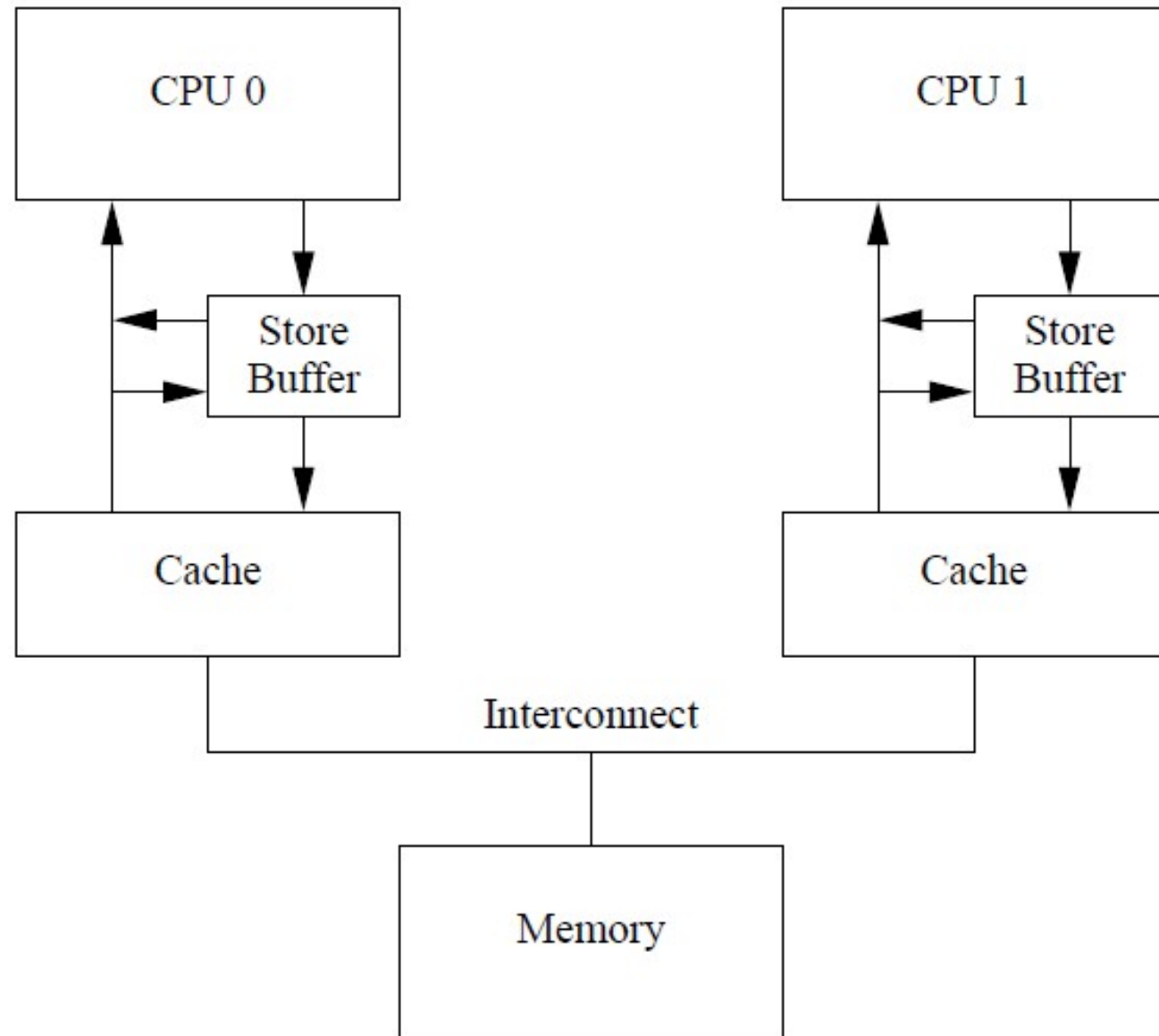
IA-64

Y
Y
Y
Y
Y
Y
Y

Table 5: Summary of Memory Ordering

@何登成

weibo.com/u/2216172320



## WEAK

## STRONG

Really weak



Weak with  
data dependency  
ordering



*Usually strong*  
(implicit acquire/  
release & TSO, usually)



Sequentially  
consistent

DEC Alpha



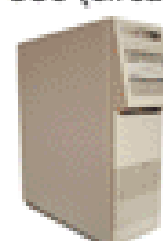
ARM



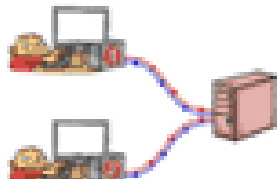
x86/64



dual 386 (circa 1989)



C/C++11  
low-level atomics



Source control  
analogy

PowerPC



SPARC TSO



Java volatile

C/C++11  
default atomics

Or, run on  
a single core  
without optimization

# Test Case

- [ordering.zip](#)
  - LoadLoad
  - LoadStore
  - StoreLoad
  - StoreStore

```
$. /ordering 1
```

```
StoreLoad Reorder Tests.
```

```
312 storeload reorders detected after 10000 iterations
1203 storeload reorders detected after 20000 iterations
2049 storeload reorders detected after 30000 iterations
2822 storeload reorders detected after 40000 iterations
3787 storeload reorders detected after 50000 iterations
4491 storeload reorders detected after 60000 iterations
4492 storeload reorders detected after 70000 iterations
4495 storeload reorders detected after 80000 iterations
4496 storeload reorders detected after 90000 iterations
4497 storeload reorders detected after 90000 iterations
```

```
$. /ordering 3
```

```
0 loadstore reorders detected after 10000 iterations
0 loadstore reorders detected after 20000 iterations
0 loadstore reorders detected after 30000 iterations
0 loadstore reorders detected after 40000 iterations
0 loadstore reorders detected after 50000 iterations
0 loadstore reorders detected after 60000 iterations
0 loadstore reorders detected after 70000 iterations
0 loadstore reorders detected after 80000 iterations
0 loadstore reorders detected after 90000 iterations
```

```
$. /ordering 2
```

```
0 loadload & storestore reorders detected after 10000 iterations
0 loadload & storestore reorders detected after 20000 iterations
0 loadload & storestore reorders detected after 30000 iterations
0 loadload & storestore reorders detected after 40000 iterations
0 loadload & storestore reorders detected after 50000 iterations
0 loadload & storestore reorders detected after 60000 iterations
0 loadload & storestore reorders detected after 70000 iterations
0 loadload & storestore reorders detected after 80000 iterations
0 loadload & storestore reorders detected after 90000 iterations
```

# 参考资料

- **Gallery of Processor Cache Effects**
- HyperThreading Technology in the NetBurst MicroArchitecture
- Intel's Sandy Bridge Microarchitecture
- **Memory Barriers a Hardware View for Software Hackers**
- **Modern Microprocessors - A 90 Minute Guide!**
- SMP Primer for Android
- **The Architecture of the Nehalem Processor and Nehalem-EP SMP Platforms**
- **Processor Microarchitecture an Implementation Perspective**
- **A Primer on Memory Consistency and Cache Coherence**
- **Everything You Always Wanted to Know About Synchronization but Were Afraid to Ask**
- [Jeff Preshing's Blog](#)
- [Martin Thompson's Blog](#)
- **CPU Cache and Memory Ordering**
- ... ..

Questions?