# Intel Core i7 Memory Hierarchy
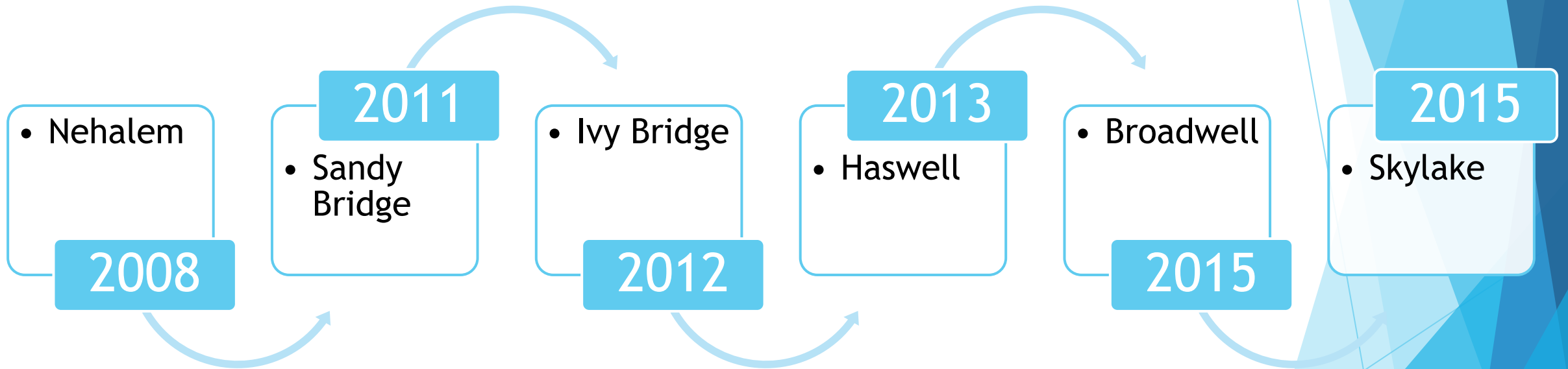
Amanda Adkins, Brett Ammeson, James Anouna,
Tony Garside, Lukas Hunker, Sam Mailand
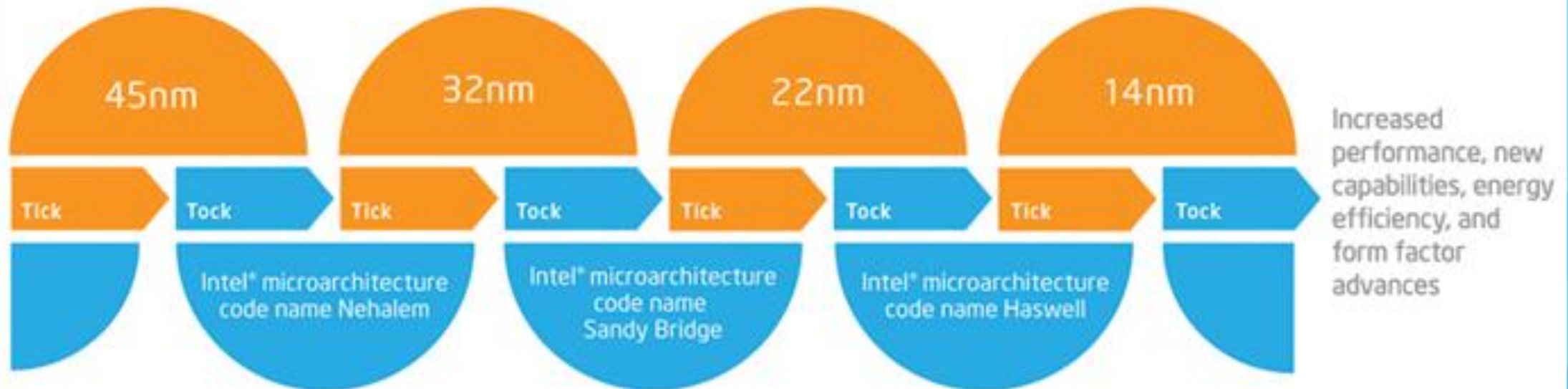
# Intel i7 Timeline

**2008**
- Nehalem

**2011**
- Sandy Bridge

**2012**
- Ivy Bridge

**2013**
- Haswell

**2015**
- Broadwell

**2015**
- Skylake

# The Tick-Tock model through the years

Manufacturing process technology    Microarchitectures

45nm    32nm    22nm    14nm

Tick    Tock    Tick    Tock    Tick    Tock    Tick    Tock

Intel® microarchitecture code name Nehalem

Intel® microarchitecture code name Sandy Bridge

Intel® microarchitecture code name Haswell

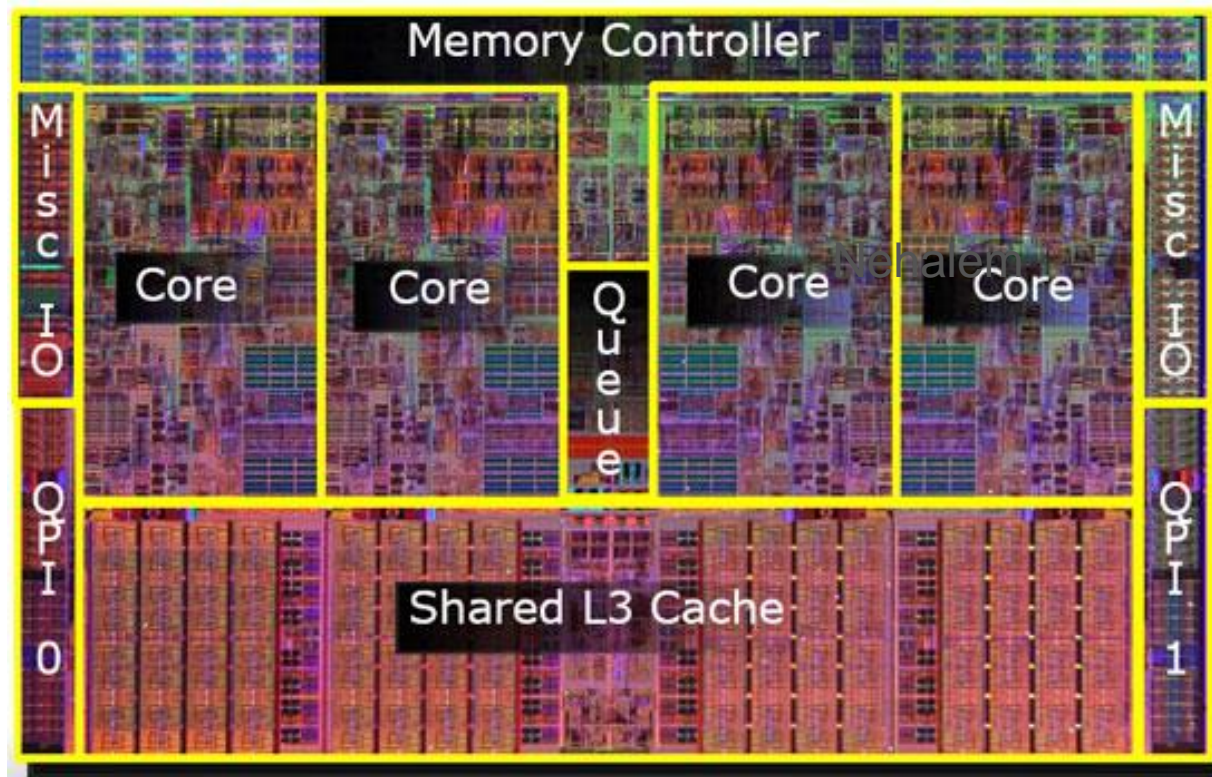Increased performance, new capabilities, energy efficiency, and form factor advances
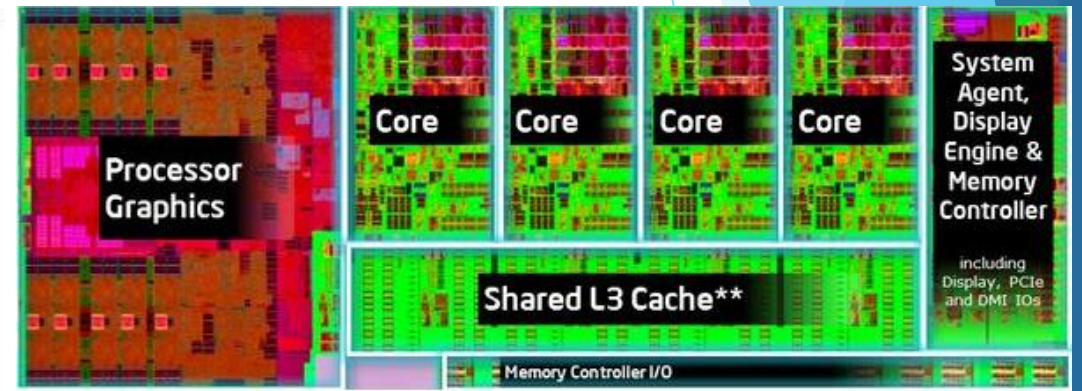
# Core i7 Basic Structure

- 4 cores
- Hyper threaded – 8 threads
- Pipelined with 16 stages
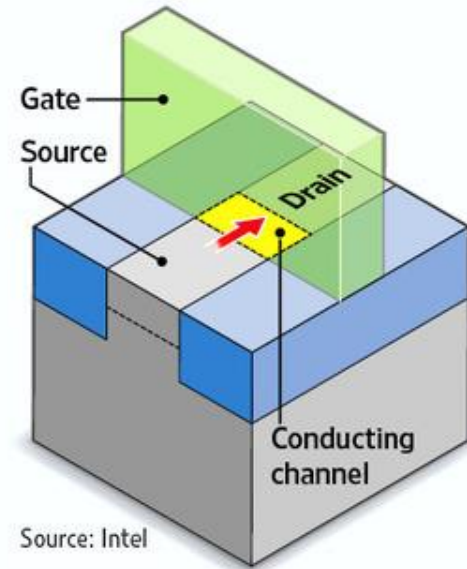
# Footprint



Nehalem (First Gen)



Haswell (Fourth Gen)

# Major Developments



**Intel's Move Into 3-D**

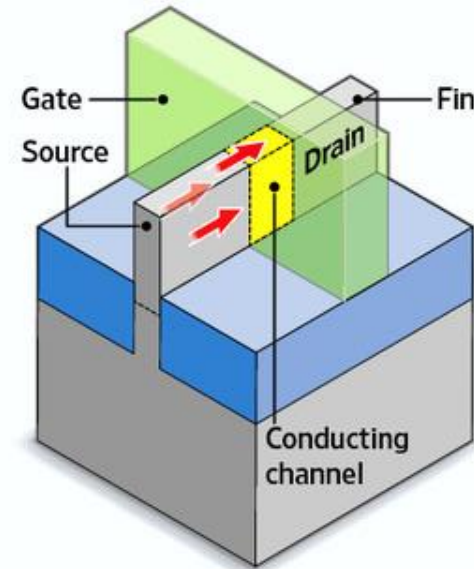The chip maker breaks from conventional approaches to make transistors.
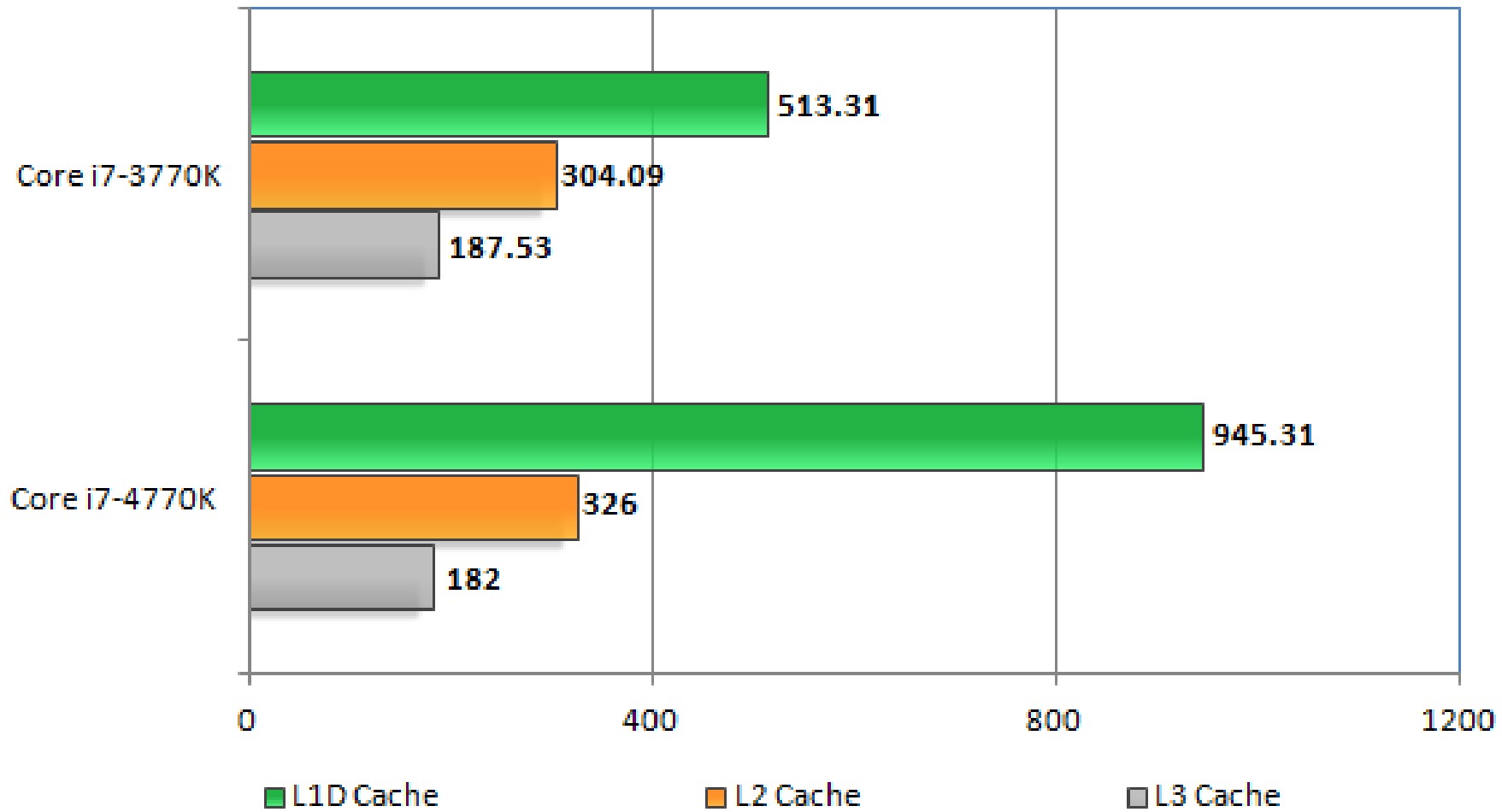
**Conventional transistor:** Electrons flow between components called a **source** and a **drain**, forming a two-dimensional **conducting channel**. A component called a **gate** starts and stops the flow, switching a transistor on or off.

**Intel's new transistor:** A fin-like **structure** rises above the surface of the transistor with the **gate** wrapped around it, forming **conducting channels** on three sides. The design takes less space on a chip, and improves speed and reduces power consumption.
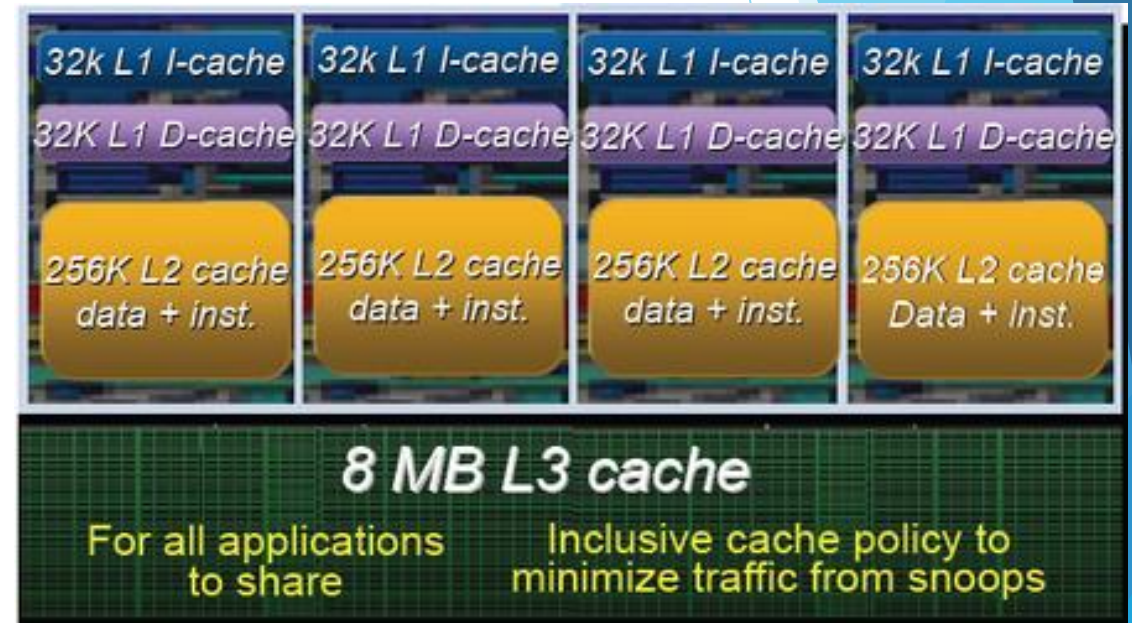
Gate — Source — Drain — Conducting channel

Gate — Source — Drain — Fin — Conducting channel

Source: Intel

Increased Cache Bandwidth

# Intel Core i7 Caching Basics

- Intel core i7 processors feature three levels of caching.
    - Separate L1 and L2 cache for each core.
    - L1 cache broken up into to halves, instruction/data.
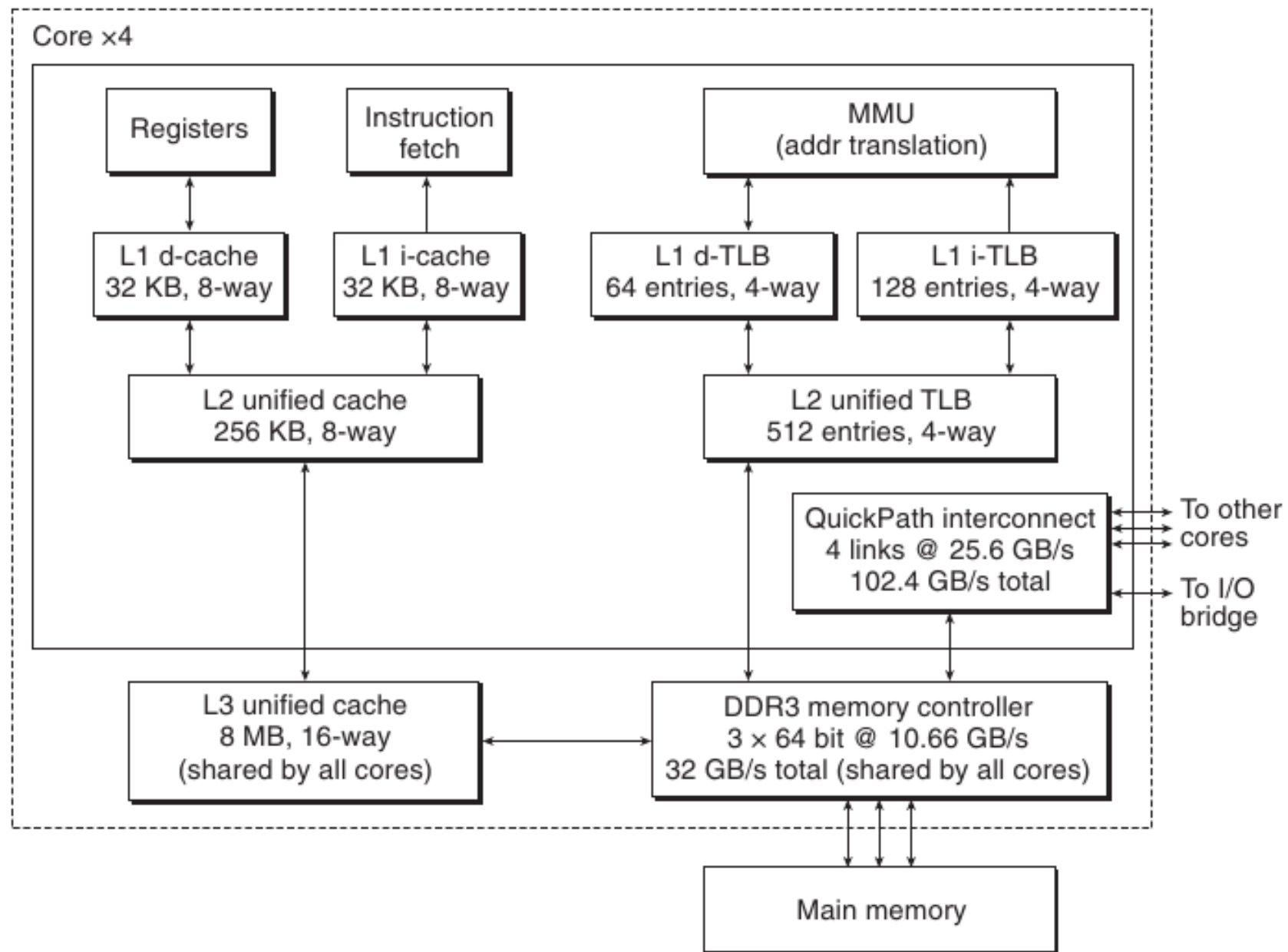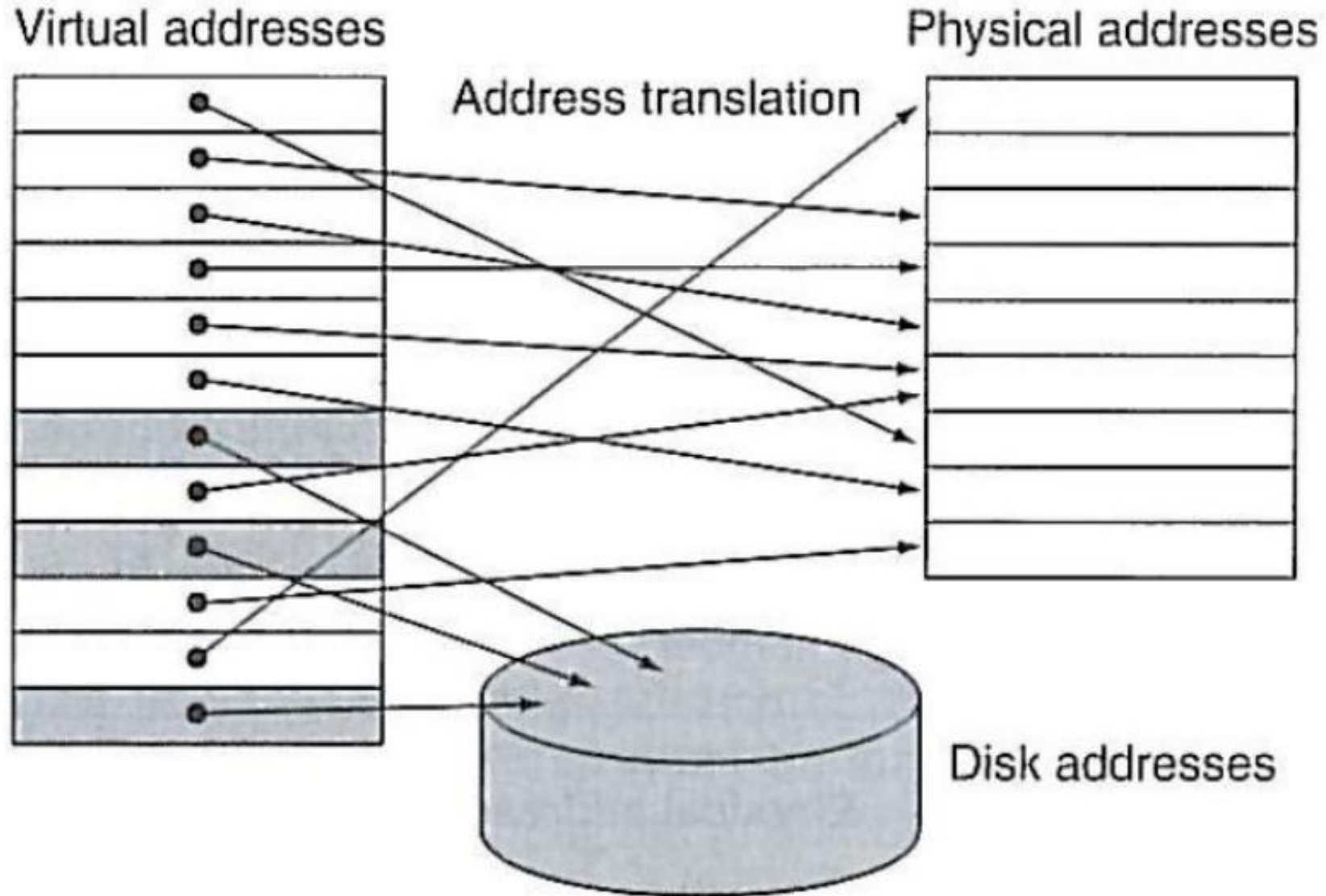    - L3 cache shared among all cores and is inclusive.

Figure 9.21    The Core i7 memory system.

# Virtual Addressing
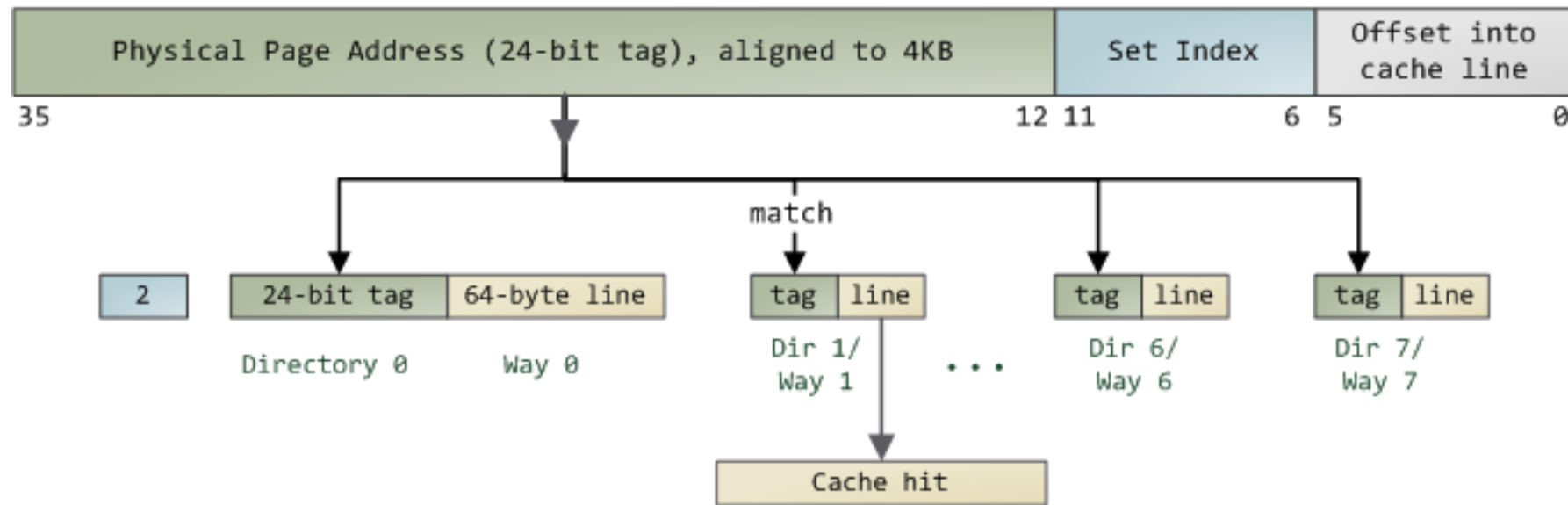


Virtual addresses      Address translation      Physical addresses

Disk addresses

# Physical Addressing



2. Search for matching tag in the set

36-bit memory location as interpreted by the L1 cache:

| Physical Page Address (24-bit tag), aligned to 4KB | Set Index | Offset into cache line |
|---|---|---|

35                                                  12 11      6 5        0

match

| 2 | | 24-bit tag | 64-byte line | | tag | line | | tag | line | | tag | line |

Directory 0      Way 0         Dir 1/    . . .    Dir 6/         Dir 7/
                               Way 1              Way 6          Way 7

Cache hit

# N-way set associativity (Review)

- Multiple entries per index
- Narrows search area needed to find unused slot
- i7 4790
  - L1 4x32 KB 8-way
  - L2 4/256 KB 8-way
  - L3 shared 8 MB 16-way

# Intel's core i7 TLB design

- Memory cache that stores recent translations of virtual memory to physical addresses for faster retrieval.

- Uses a 2 level cache system

- L1 TLB
  - Divided into 2 parts
  - Data TLB: 64 4KB entries
  - Instruction TLB: 128 4KB entries

- L2 TLB (Services misses in L1 DTLB)
  - Can hold translations for 4KB and 2 MB pages (vs. only 4KB)
  - 1024 entries (vs. 512)
  - 8-way associative (vs. 4-way)

# TLB Comparisons between generations

## Nehalem

| Cache | | Page Size | |
|---|---|---|---|
| Name | Level | 4 KB | 2 MB |
| DTLB | 1st | 64 | 32 |
| ITLB | 1st | 128 | 7 / logical core |
| STLB | 2nd | 512 | none |

## Sandy Bridge and Ivy Bridge

| Cache | | Page Size | | |
|---|---|---|---|---|
| Name | Level | 4 KB | 2 MB | 1 GB |
| DTLB | 1st | 64 | 32 | 4 |
| ITLB | 1st | 128 | 8 / logical core | none |
| STLB | 2nd | 512 | none | none |

## Haswell

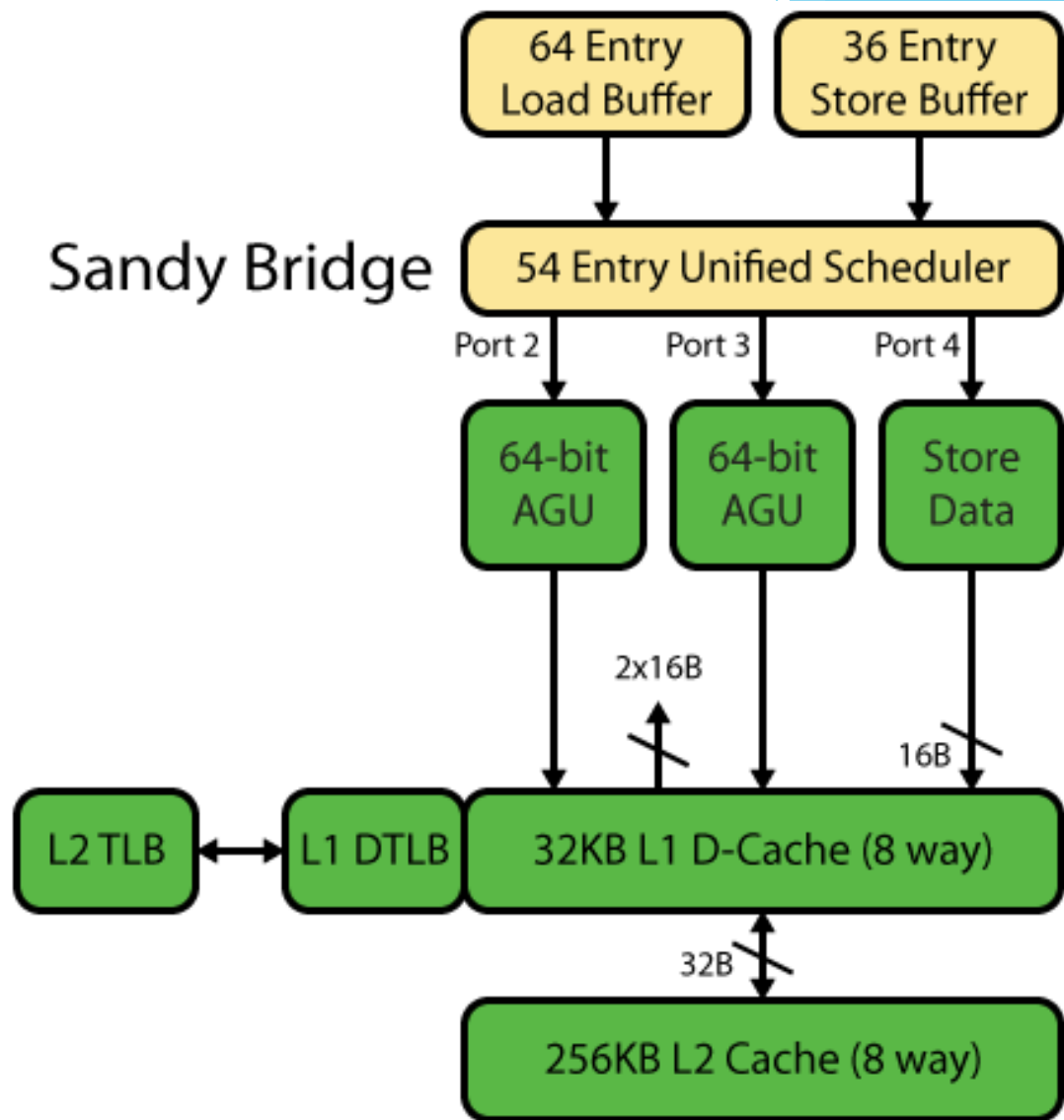| Cache | | Page Size | | |
|---|---|---|---|---|
| Name | Level | 4 KB | 2 MB | 1 GB |
| DTLB | 1st | 64 | 32 | 4 |
| ITLB | 1st | 128 | 8 / logical core | none |
| STLB | 2nd | 1024 | | none |

# Pseudo-LRU (Intel's core i7 caching algorithm)

- One bit per cache line
- Resets after all lines' bit is set
- Lowest line index with a '0' replaced

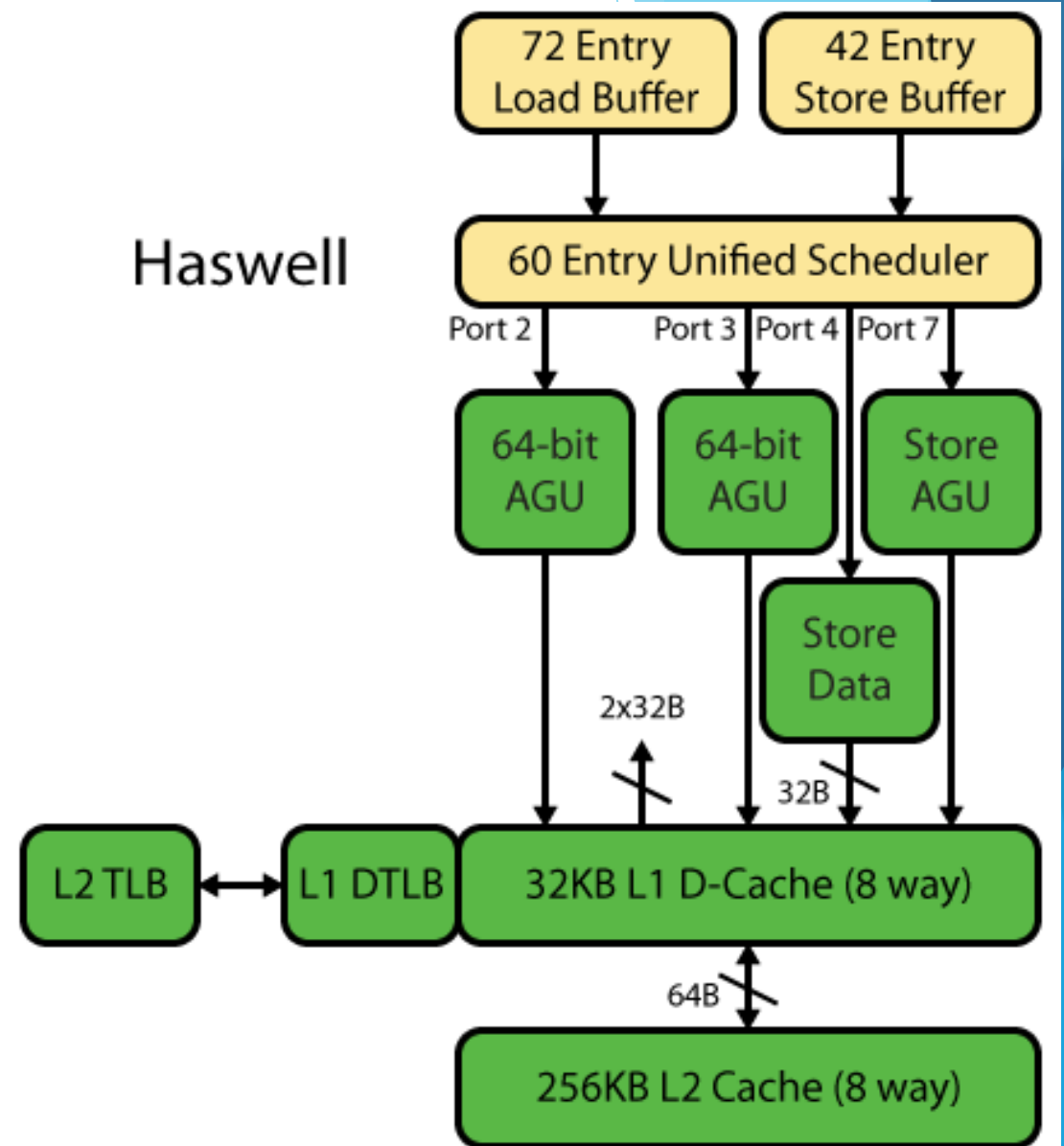- Port 2 and 3 are the Address Generation Units

- Port 4 for writing data from the core to the L1 Cache

- Additional port added to Haswell

- Haswell can sustain 2 loads and 1 store per cycle "under nearly any circumstances"

- Forwarding latency for AVX loads decreased from 2 to 1 cycle

- AVX: Set of instructions for doing SIMD operations on Intel CPUs

- 4 Split line buffers to resolve unaligned loads (vs 2 in Sandy-bridge)

- Decrease impact of unaligned access



Haswell

72 Entry Load Buffer
42 Entry Store Buffer
60 Entry Unified Scheduler
Port 2  Port 3  Port 4  Port 7
64-bit AGU
64-bit AGU
Store AGU
Store Data
2x32B
32B
L2 TLB
L1 DTLB
32KB L1 D-Cache (8 way)
64B
256KB L2 Cache (8 way)

# Haswell L1 Cache

- 32 kb
- 8 way associative
- Writeback
- TLB access & cache tag can occur in parallel
- Does not suffer from bank conflicts (unlike Sandy Bridge)
- Minimum latency: 4 cycles (same as Sandy-Bridge)
- Minimum lock latency of haswell is 12 cycles (sandy-bridge was 16)

# Haswell L2 Cache

- Bandwidth doubled

- Can deliver 64 bit line to data or instruction cache every cycle

- 11 cycle latency

- 256 KB for each cache

# Haswell L3 Cache

- Shared between all cores
- Size varies between models and generations between 6MB and 15MB
- Most Haswell models have an 8MB cache
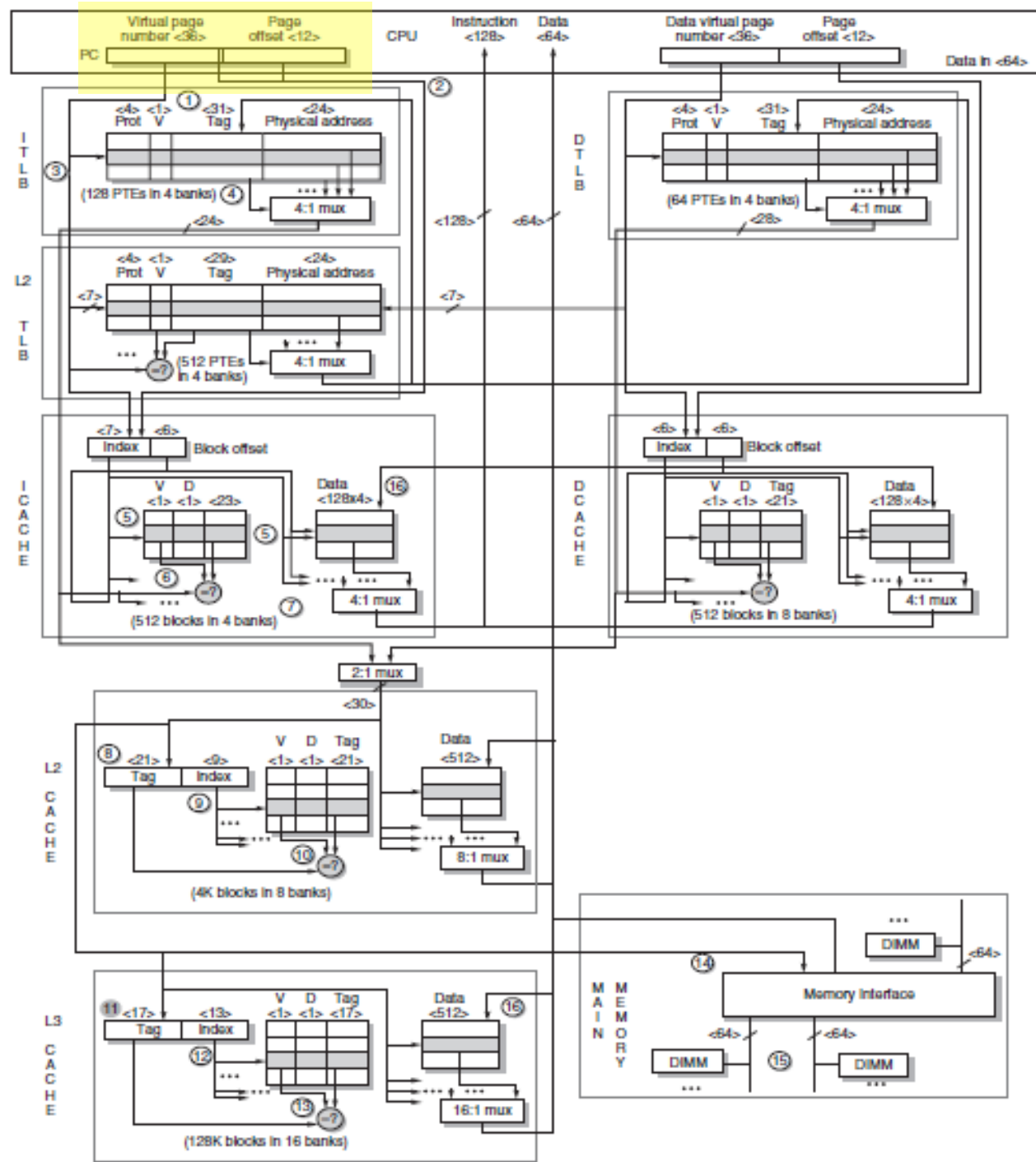- Size reduced for power efficiency

# Shared Data

- Transactional Synchronization Extensions
  - Transactional memory
- Hardware Lock Elision
  - Backwards Compatible, Windows only
  - Uses instruction prefixes to lock and release
- Restricted Transactional Memory
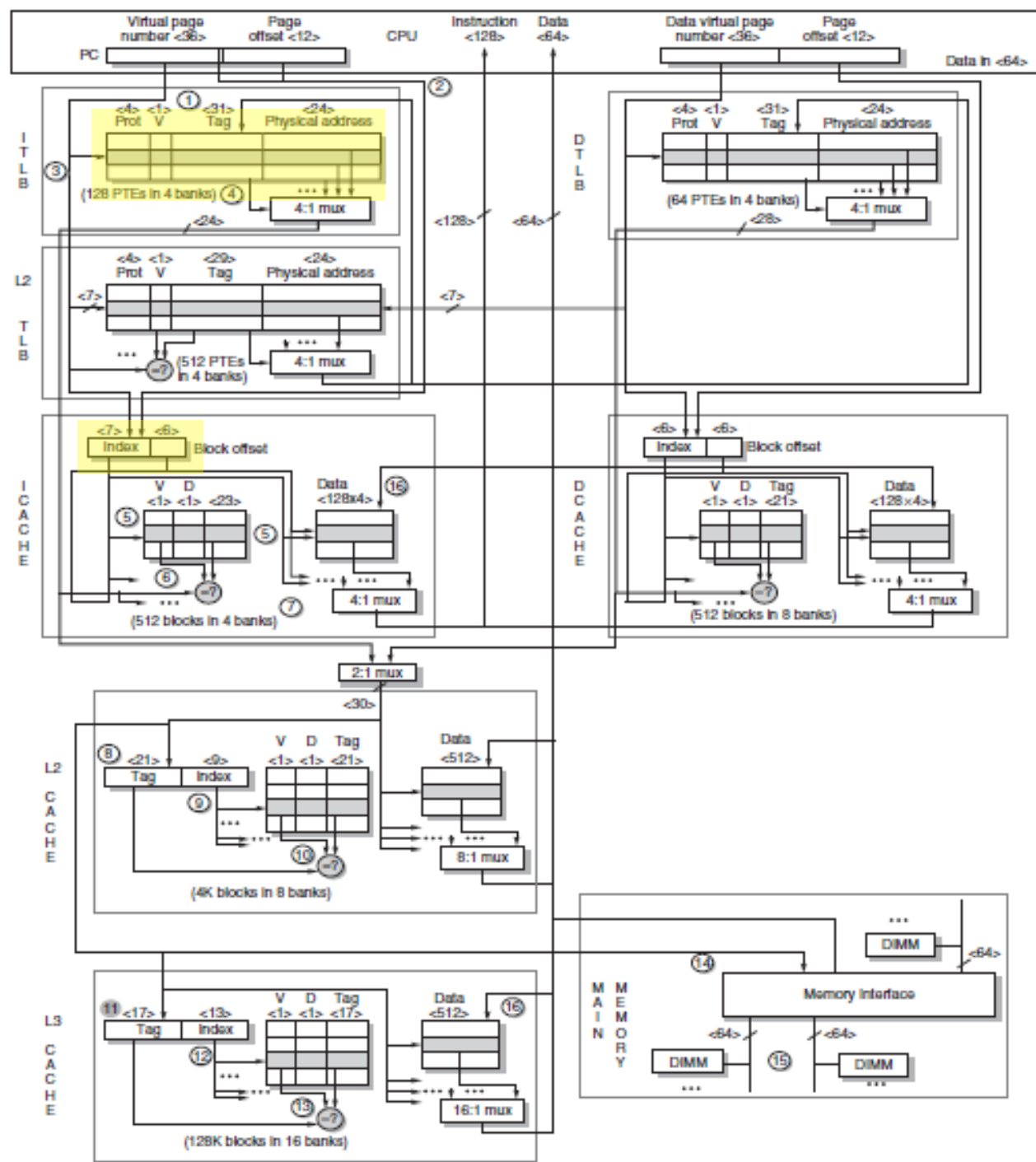  - Newer, more flexible
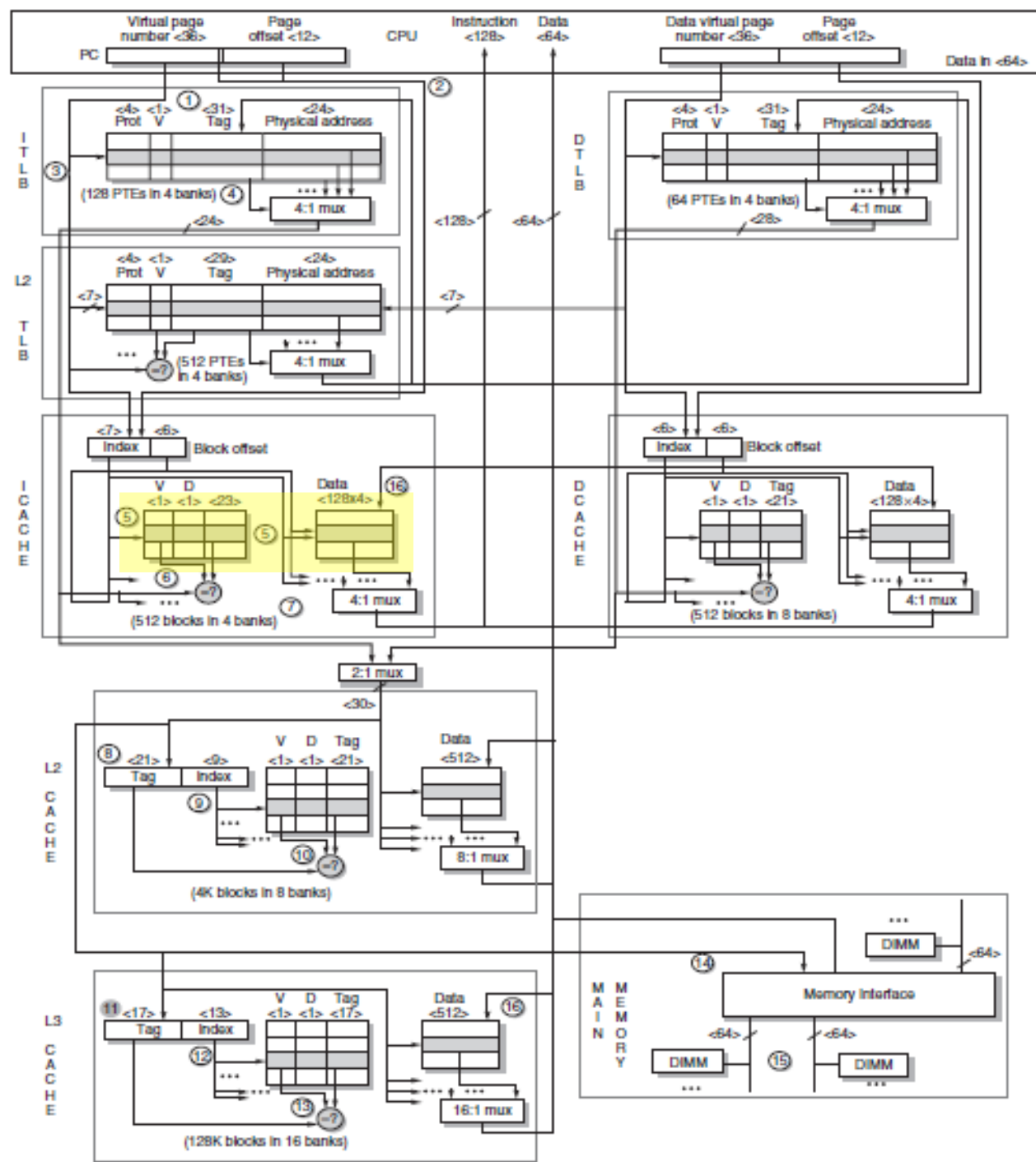  - Fallback code in case of failure

# Pre-fetching

- Fetch Instructions/Data before needed
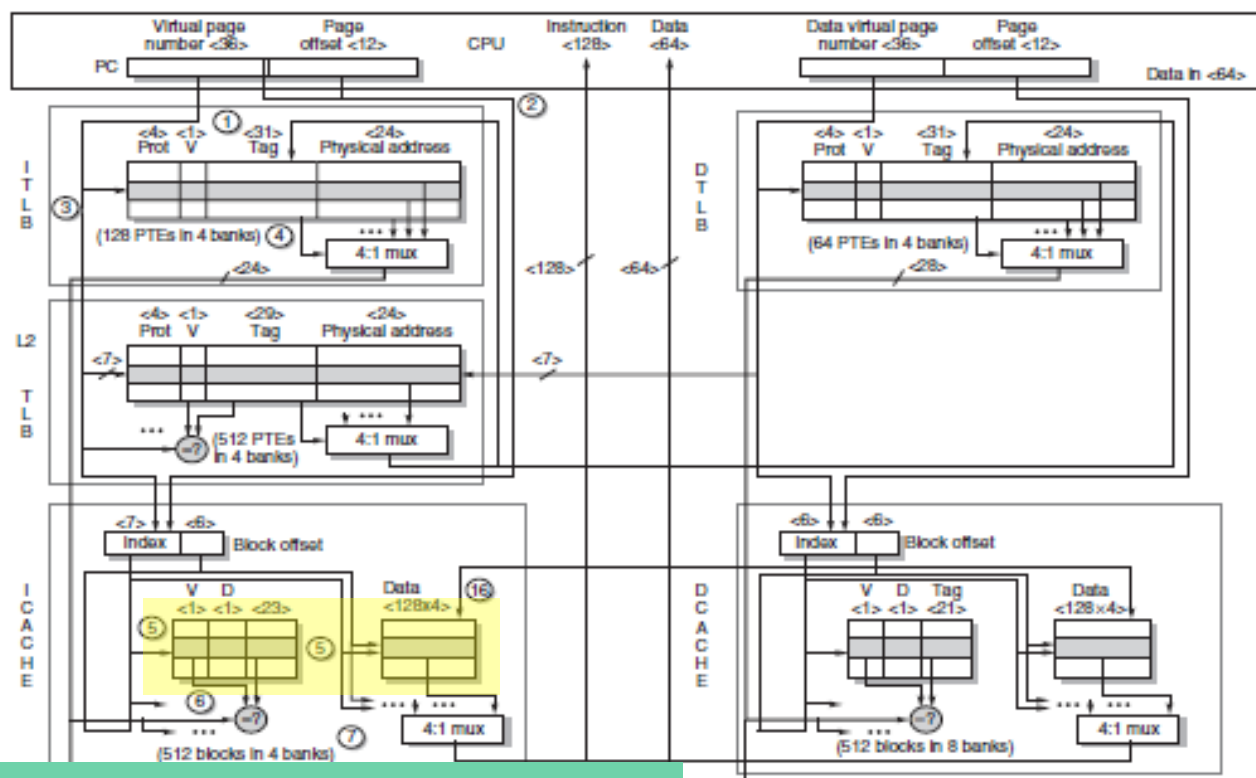  - On a miss 2 blocks are fetched
- If successful, miss will grab from buffer, and pre-fetch next block
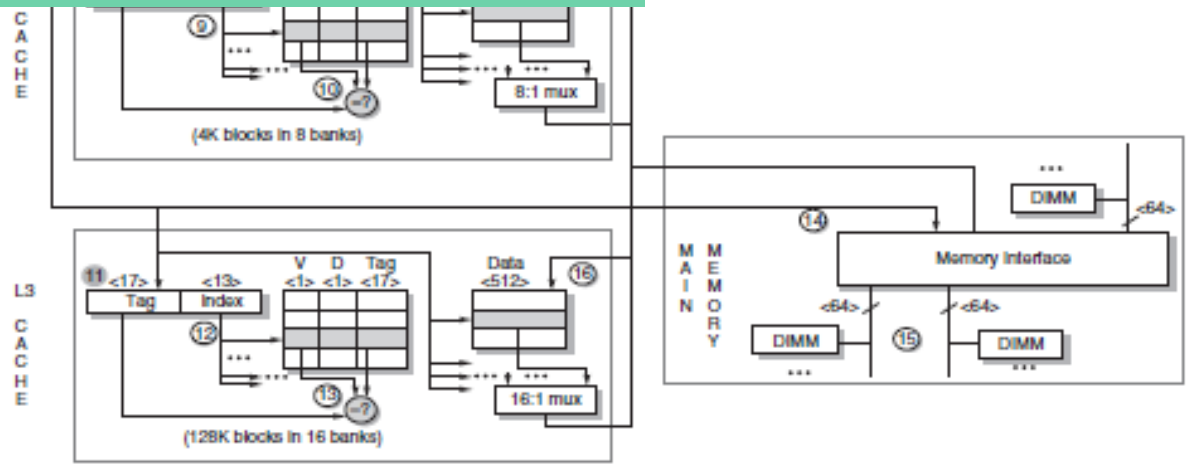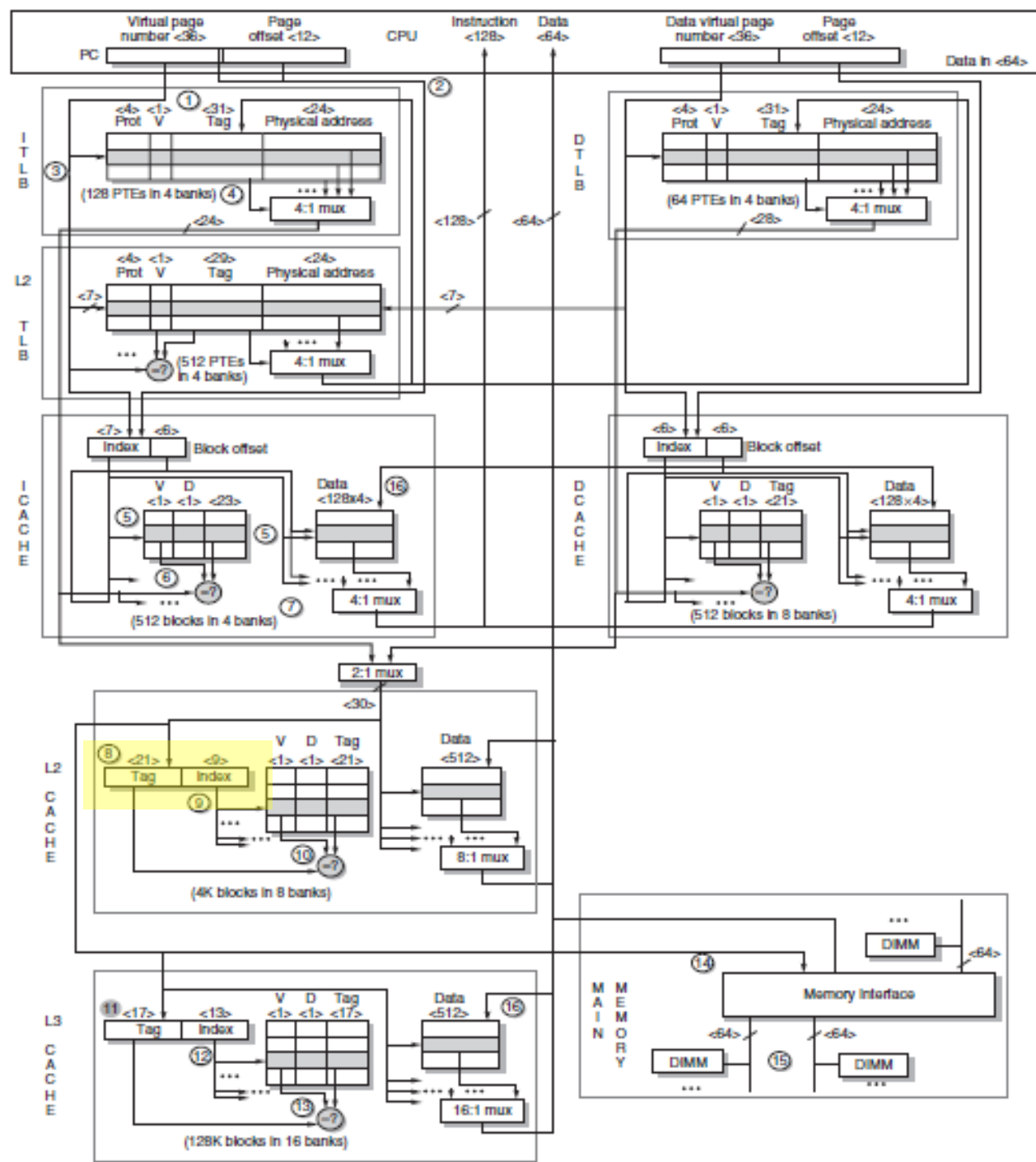
# Memory Hierarchy Access Steps

Virtual page number <36>   Page offset <12>   CPU   Instruction <128>   Data <64>   Data virtual page number <36>   Page offset <12>

PC

Data in <64>

**ITLB**

① ② ③ ④

<4> <1>   <31>   <24>
Prot  V   Tag   Physical address

(128 PTEs in 4 banks) ④

4:1 mux

<24>

**DTLB**

<4> <1>   <31>   <24>
Prot  V   Tag   Physical address

(64 PTEs in 4 banks)

4:1 mux

<28>

**L2 TLB**

<4> <1>   <29>   <24>
Prot  V   Tag   Physical address

<7>

<7>

=? (512 PTEs in 4 banks)

4:1 mux

<128>   <64>

**ICACHE**

<7>   <6>
Index   Block offset

V   D        Data
<1> <1> <23>   <128x4>   ⑯

⑤   ⑤

⑥   =?   ⑦

4:1 mux

(512 blocks in 4 banks)

**DCACHE**

<6>   <6>
Index   Block offset

V   D   Tag   Data
<1> <1> <21>   <128x4>

=?

4:1 mux

(512 blocks in 8 banks)

2:1 mux

<30>

**L2 CACHE**

⑧ <21>   <9>   V   D   Tag   Data
Tag   Index   <1> <1> <21>   <512>

⑨

=? ⑩

8:1 mux

(4K blocks in 8 banks)

DIMM   <64>

**MAIN MEMORY**

⑭

Memory Interface

**L3 CACHE**

⑪ <17>   <13>   V   D   Tag   Data
Tag   Index   <1> <1> <17>   <512>   ⑯

⑫

=? ⑬

16:1 mux

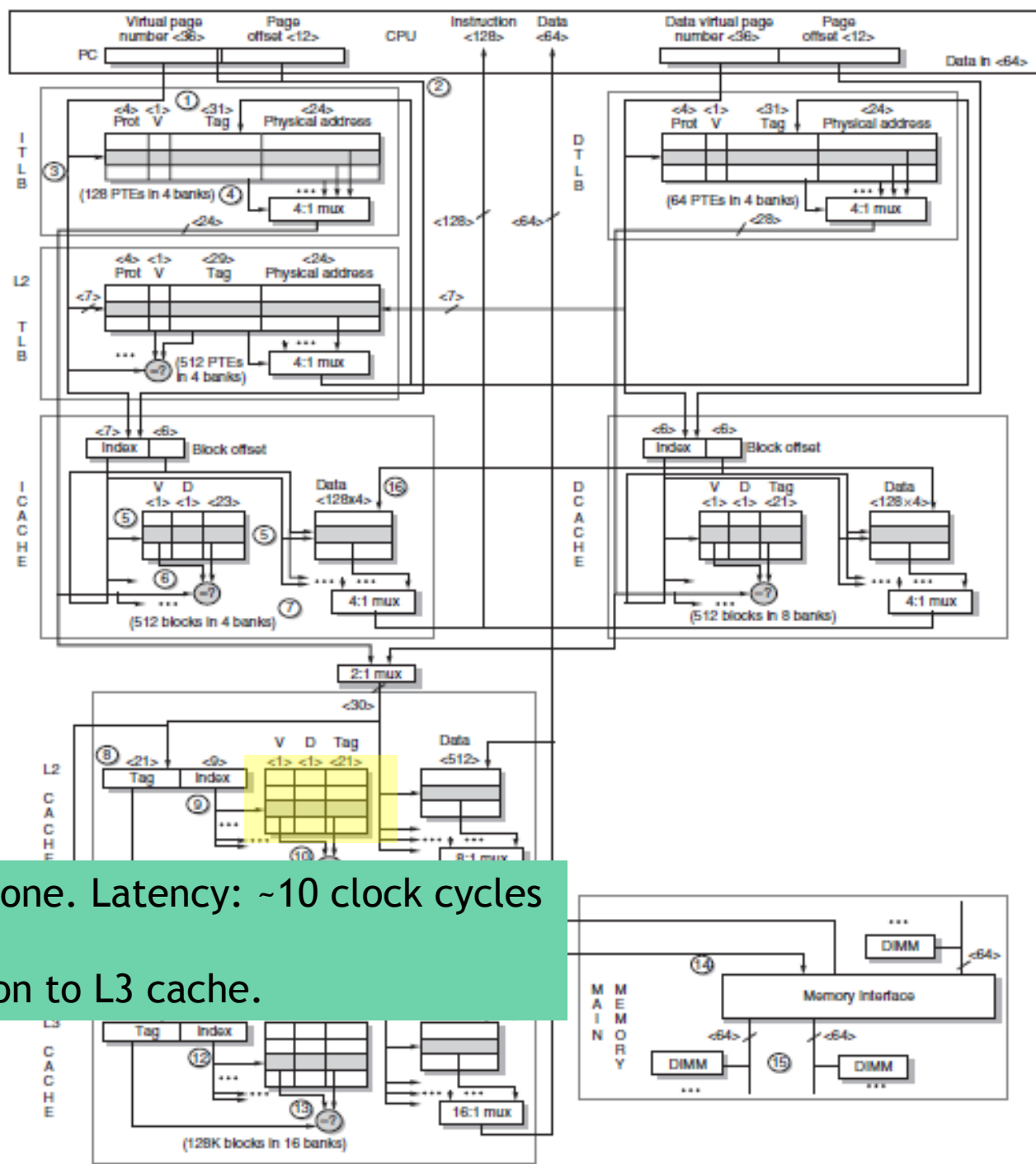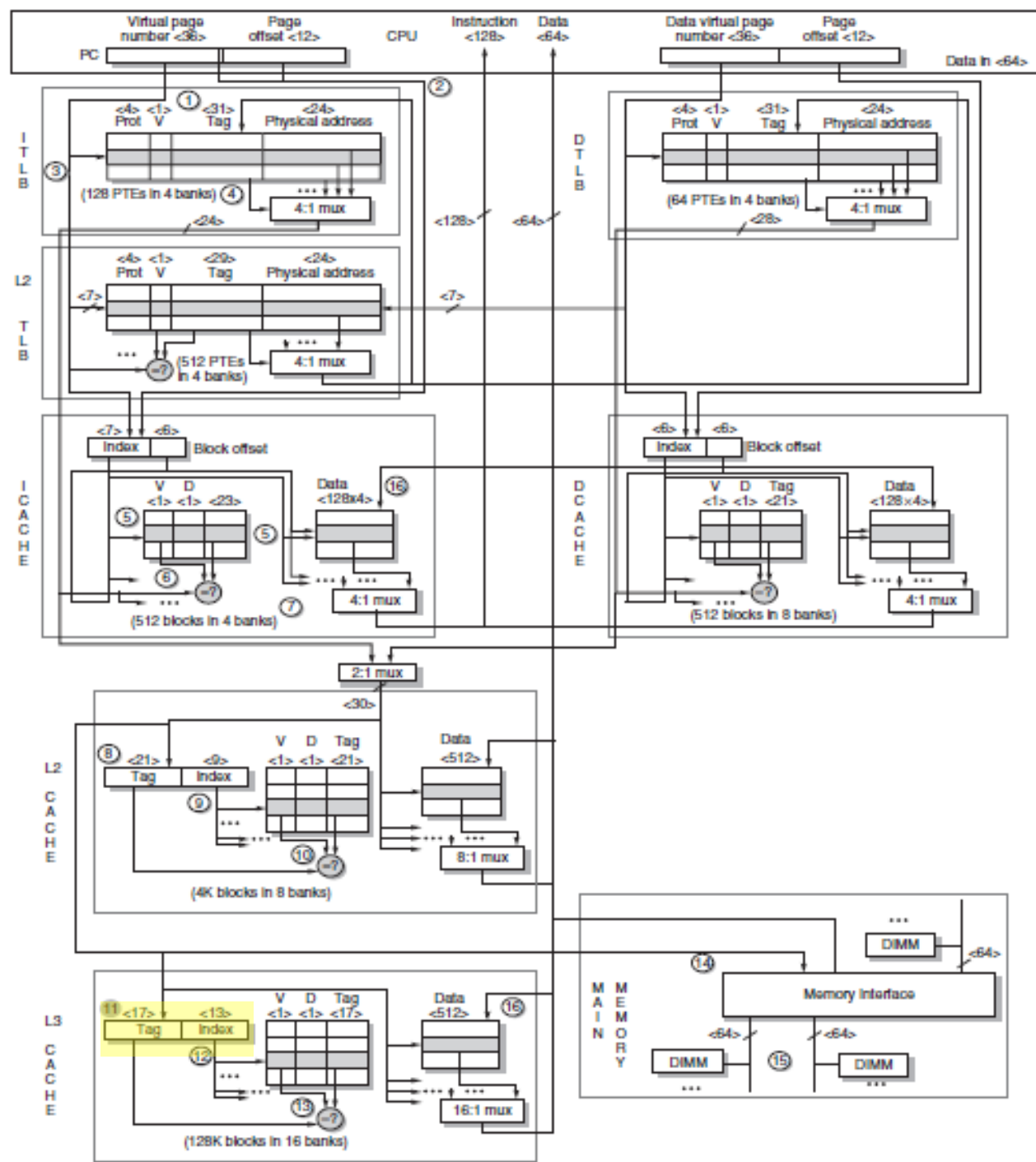(128K blocks in 16 banks)

<64>   <64>
DIMM   ⑮   DIMM

Cache hit! We're done. Latency: ~4 clock cycles
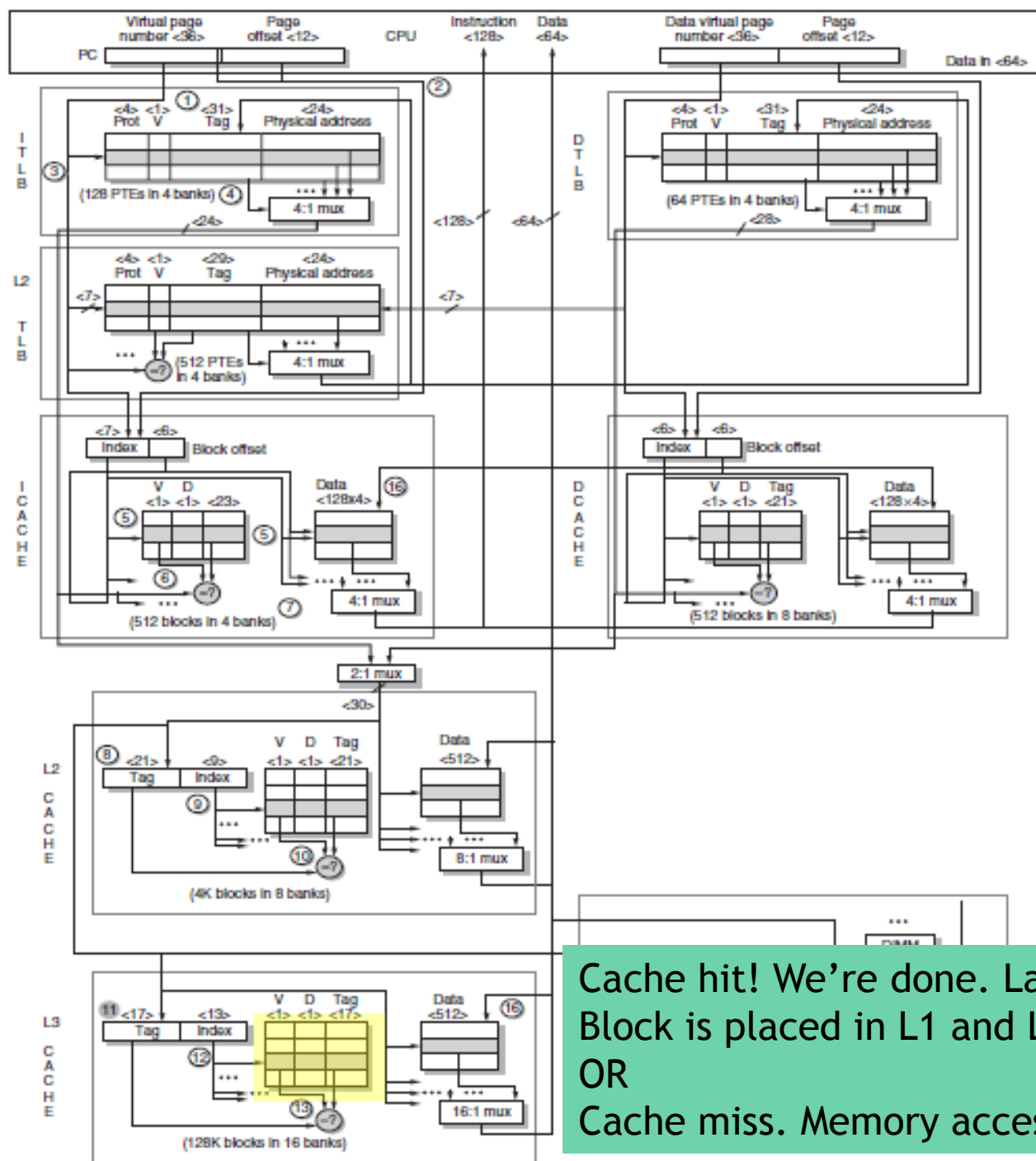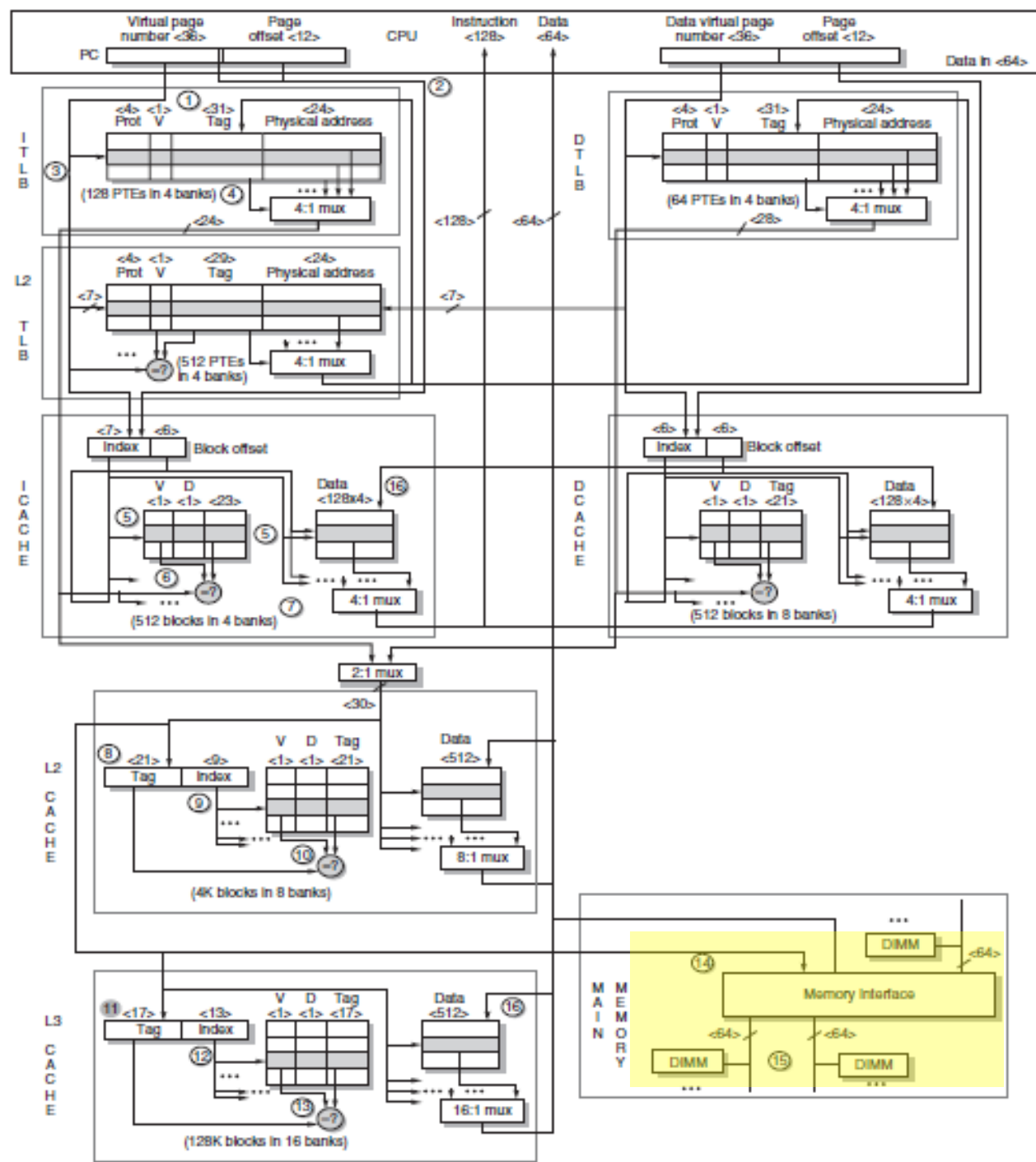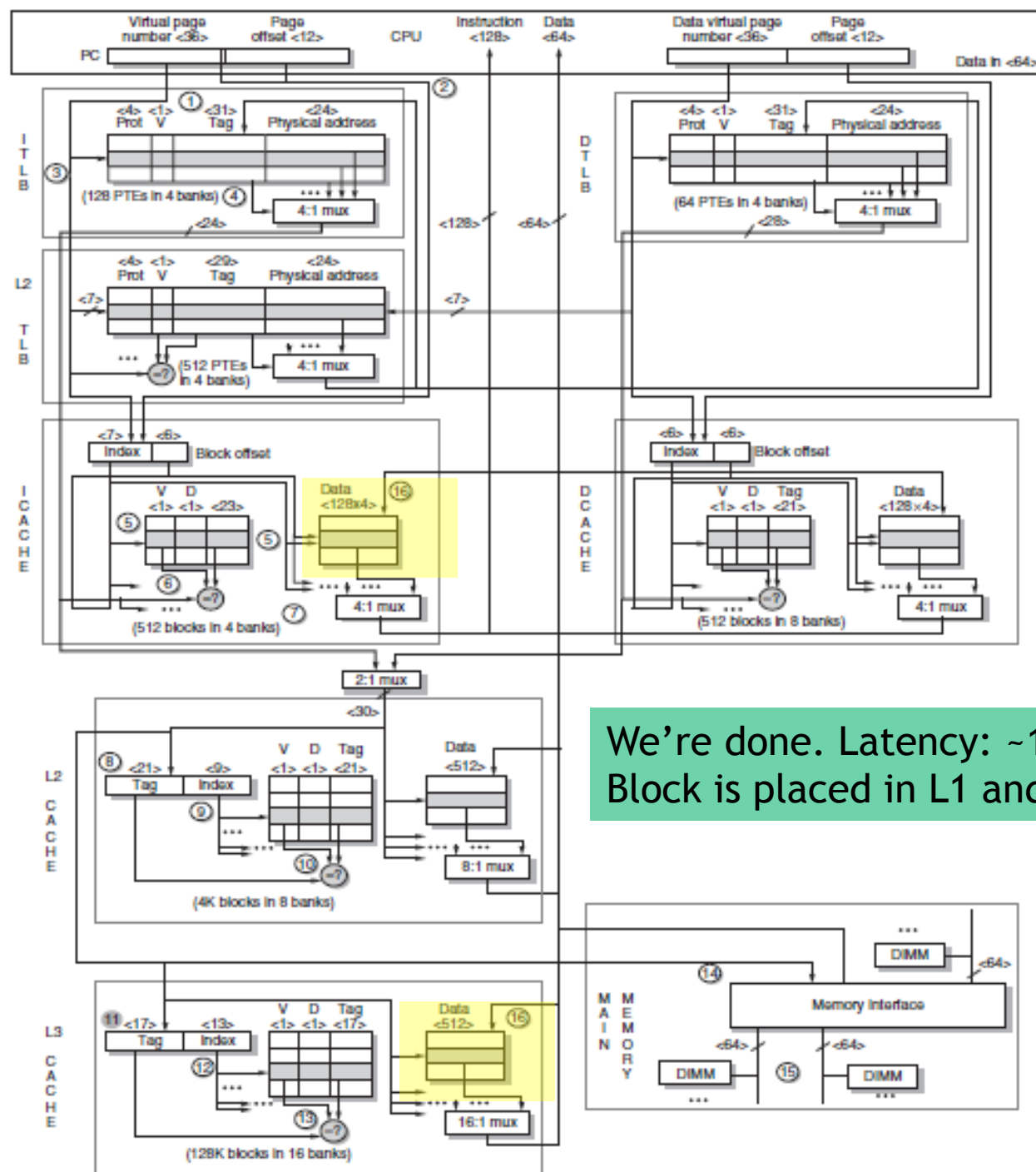OR
Cache miss. Move on to L2 cache.

Cache hit! We're done. Latency: ~10 clock cycles
OR
Cache miss. Move on to L3 cache.

Cache hit! We're done. Latency: ~35 clock cycles
Block is placed in L1 and L3 cache
OR
Cache miss. Memory access is initiated.

We're done. Latency: ~135 clock cycles
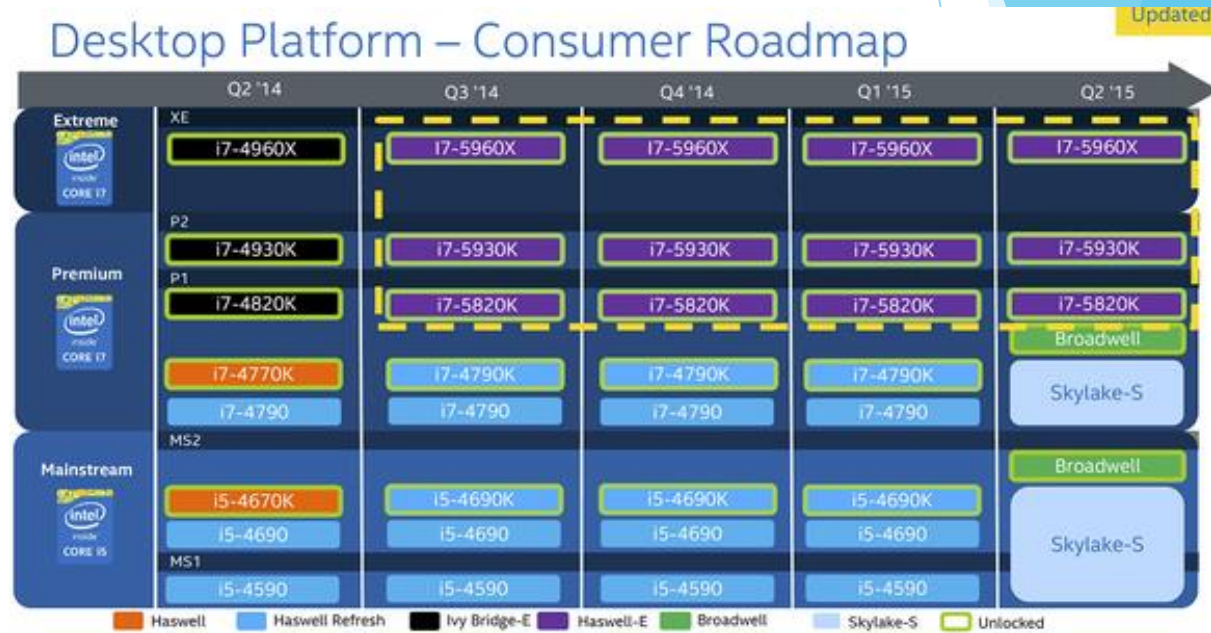Block is placed in L1 and L3 cache.

# Generation 5 (Broadwell)

- Currently mobile only (Lower power systems)
  - Two cores
- Shrunk to 14 nm
- Power Consumption down to 15 w
- No low-end desktop processors
- Extended instruction set

# Future Releases

- Broadwell Desktop
  - Many manufacturers plan to skip
  - Possibly due to lack of low-end offerings
- Skylake
  - Second half of 2015



Desktop Platform – Consumer Roadmap

| Metric | Nehalem | Sandy Bridge | Haswell |
|---|---|---|---|
| L1 Instruction Cache | 32K, 4-way | 32K, 8-way | 32K, 8-way |
| L1 Data Cache | 32K, 8-way | 32K, 8-way | 32K, 8-way |
| Fastest Load-to-use | 4 cycles | 4 cycles | 4 cycles |
| Load bandwidth | 16 Bytes/cycle | 32 Bytes/cycle (banked) | 64 Bytes/cycle |
| Store bandwidth | 16 Bytes/cycle | 16 Bytes/cycle | 32 Bytes/cycle |
| L2 Unified Cache | 256K, 8-way | 256K, 8-way | 256K, 8-way |
| Fastest load-to-use | 10 cycles | 11 cycles | 11 cycles |
| Bandwidth to L1 | 32 Bytes/cycle | 32 Bytes/cycle | 64 Bytes/cycle |
| L1 Instruction TLB | 4K: 128, 4-way 2M/4M: 7/thread | 4K: 128, 4-way 2M/4M: 8/thread | 4K: 128, 4-way 2M/4M: 8/thread |
| L1 Data TLB | 4K: 64, 4-way 2M/4M: 32, 4-way 1G: fractured | 4K: 64, 4-way 2M/4M: 32, 4-way 1G: 4, 4-way | 4K: 64, 4-way 2M/4M: 32, 4-way 1G: 4, 4-way |
| L2 Unified TLB | 4K: 512, 4-way | 4K: 512, 4-way | 4K+2M shared: 1024, 8-way |

All caches use 64-byte lines

# Conclusion

- Why is it faster?
    - Increased Bandwidth
    - Doubled the associativity in L2 TLB
    - Tri Gate Transistors
- Smaller chip size
- Lower power requirements
    - Decreased L3 Cache Size

# Questions?