# CPSC 313: Computer Hardware and Operating Systems

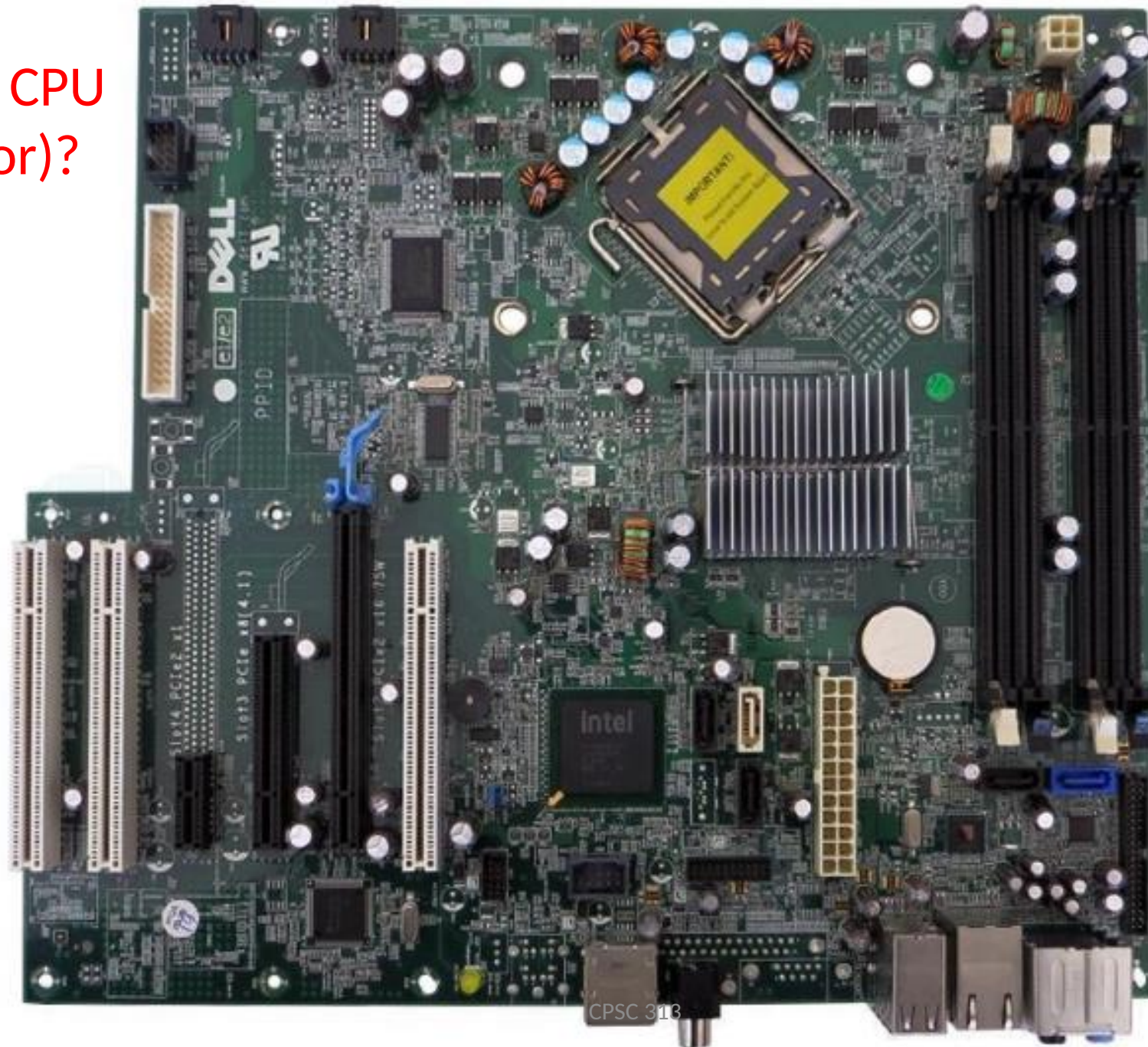## Unit 3: The Memory Hierarchy

# Administration

- Quiz 2: Don't miss your time!

  - Quiz 2 information and practice quiz were released on Friday.

- Lab 5:

  - Due Sunday October 20th (in two weeks).

  - No class Friday because of Quiz 2

  - but we'll be here to answer questions.
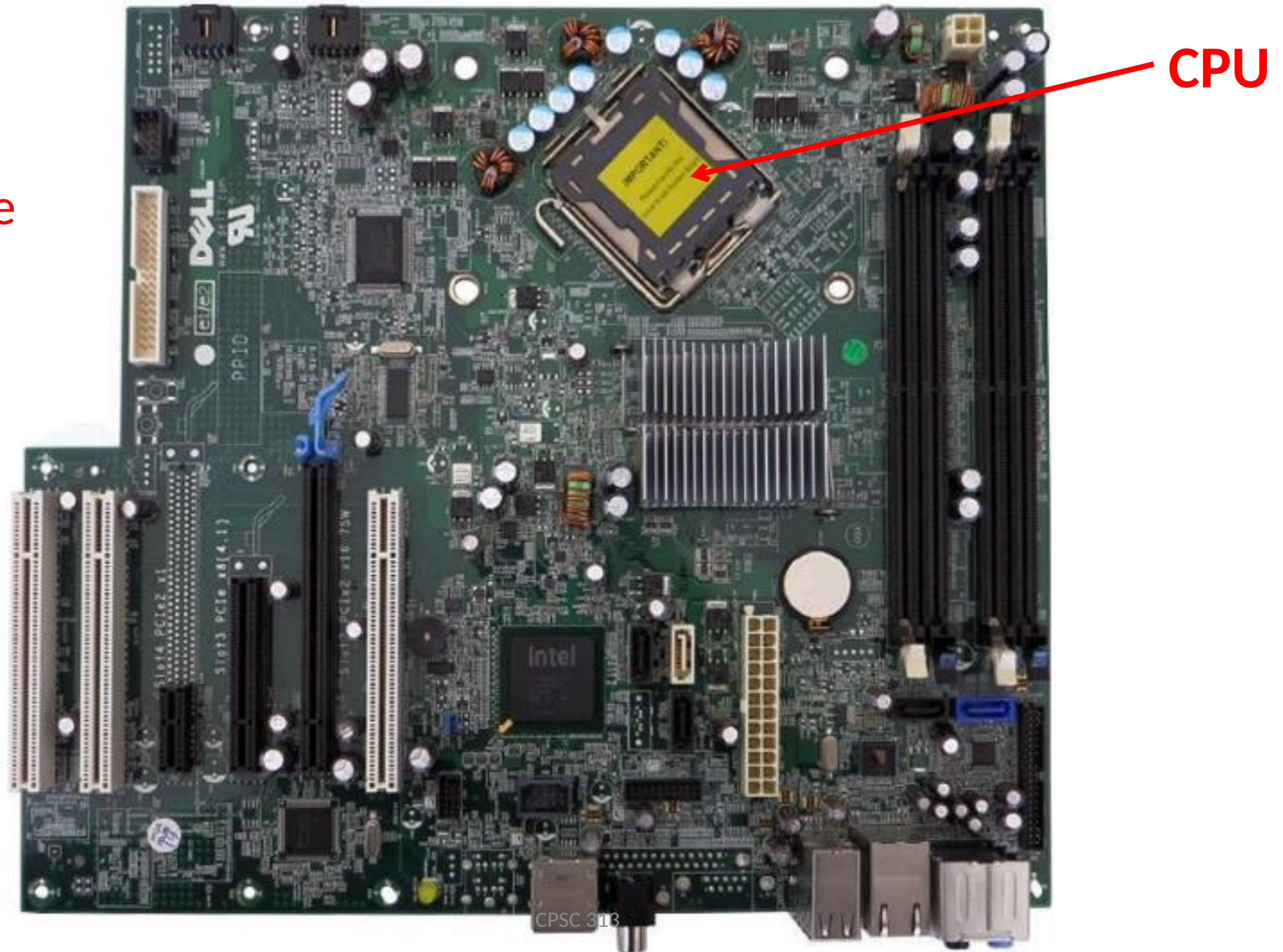
- Tutorial 4 this week!

# Today

- Learning Outcomes
  - Define memory hierarchy.
  - Evaluate the performance differences found at the different layers of the hardware memory hierarchy.
  - Explain the different kinds of caching that processors and hardware systems perform to mitigate the performance differences between the levels of the memory hierarchy.
- Reading
  - 6.2, 6.3

# Where is the CPU (aka processor)?

**CPU**

Where is the memory?

CPU

**Factor of 100x in performance!**

Memory Slots

# The Memory Hierarchy

Registers

L1 Cache

L2 Cache

L3 Cache

L4 Cache          (rare)

Main Memory

Flash Drive

Disk Drive

Bigger

Faster

# The Memory Hierarchy



Bigger

Faster

Registers

L1 Cache

L2 Cache

L3 Cache

L4 Cache    (rare)

Main Memory

Flash Drive

Disk Drive

THE UNIVERSITY OF BRITISH COLUMBIA

Computer Science
Faculty of Science

CPSC 313

# Intel Golden Cove (late 2021)

- Cycle time: 3.2 to 5.1 GHz => .2 to .3 ns/cycle (registers)

- L1: 1.25 ns (5 clock cycles; 80K/core = 640K total) PER CORE

- L2: 3.7 ns (15 clock cycles; 1280K/core = 10 MB total) PER CORE

- L3: 17 ns (67 clock cycles; 30 MB) SHARED

- Memory: 80 ns (i.e., Main memory, RAM)

THE UNIVERSITY OF BRITISH COLUMBIA
Computer Science
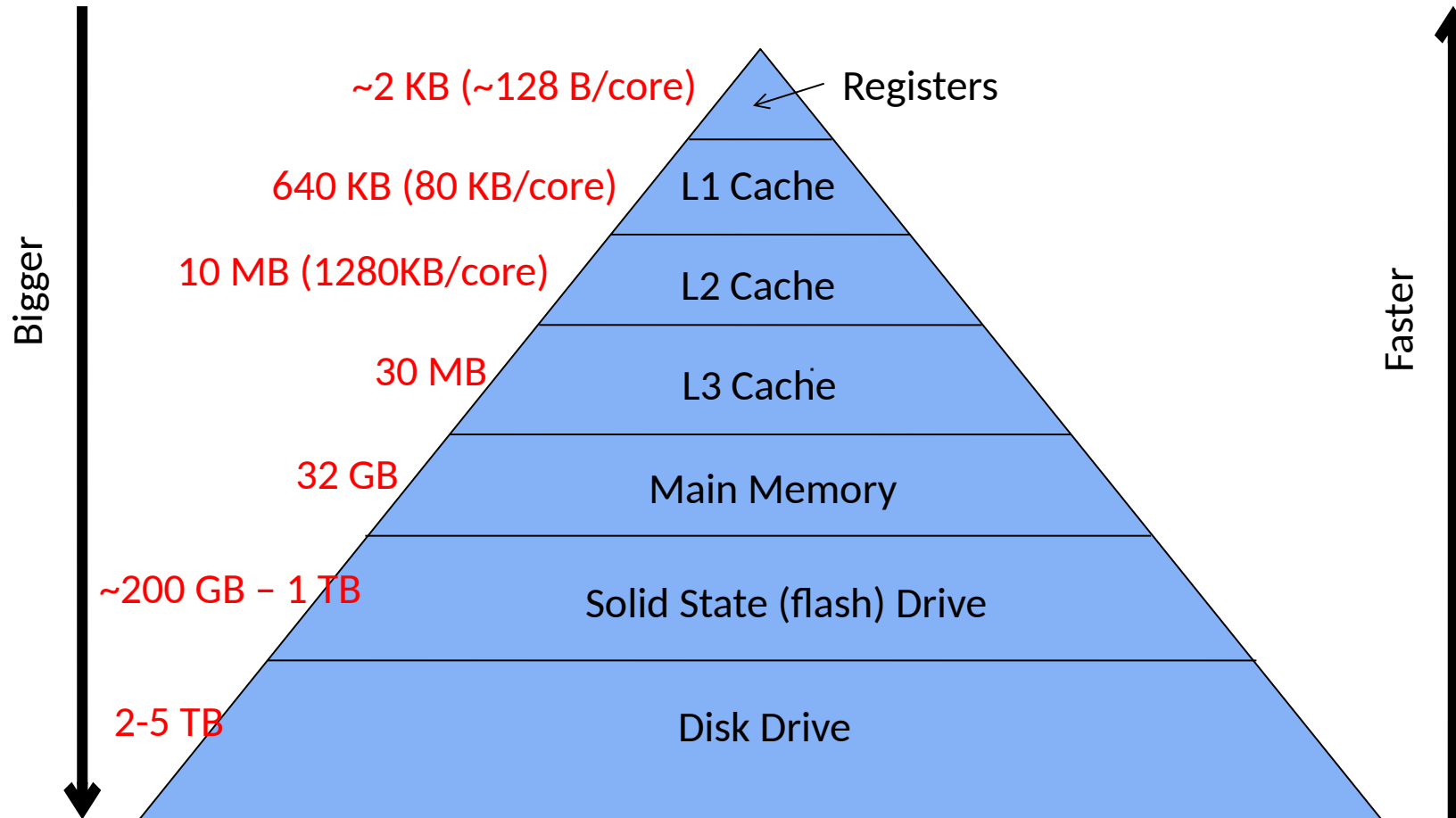Faculty of Science

# Intel Golden Cove (late 2021)

- Cycle time: 3.2 to 5.1 GHz => .2 to .3 ns/cycle (registers)

- L1: 1.25 ns (5 clock cycles; 80K/core = 640K total) PER CORE

- L2: 3.7 ns (15 clock cycles; 1280K/core = 10 MB total) PER CORE

- L3: 17 ns (67 clock cycles; 30 MB) SHARED

- Memory: 80 ns (i.e., Main memory, RAM)

THE UNIVERSITY OF BRITISH COLUMBIA
Computer Science
Faculty of Science

# The Memory Hierarchy -- Size

**Bigger** ↓

**Faster** ↑

~2 KB (~128 B/core) — Registers

640 KB (80 KB/core) — L1 Cache

10 MB (1280KB/core) — L2 Cache

30 MB — L3 Cache

32 GB — Main Memory

~200 GB – 1 TB — Solid State (flash) Drive

2-5 TB — Disk Drive

# The Memory Hierarchy -- Latency



Bigger

~2 KB (~128 B/core) — Registers — 0.3 ns

640 KB (80 KB/core) — L1 Cache — 1.25 ns

10 MB (1280KB/core) — L2 Cache — 3.7 ns

30 MB — L3 Cache — 17 ns

32 GB — Main Memory — 80 ns

~200 GB – 1 TB — Solid State (flash) Drive — ~.1 ms

2-5 TB — Disk Drive — ~3 ms

Faster

THE UNIVERSITY OF BRITISH COLUMBIA
Computer Science
Faculty of Science

# The Memory Hierarchy -- Price



Registers

~2 KB (~128 B/core)                                    0.3 ns

640 KB (80 KB/core)          L1 Cache                  1.25 ns

10 MB (1280KB/core)          L2 Cache                  3.7 ns

30 MB                        L3 Cache                  17 ns

32 GB            $3.15/**GB**   Main Memory            80 ns

~200 GB – 1 TB   $0.1/**GB**    Solid State (flash) Drive   ~.1 **ms**

2-5 TB           $0.007/**GB**  Disk Drive             ~3 **ms**

Bigger

Faster

THE UNIVERSITY OF BRITISH COLUMBIA

Computer Science
Faculty of Science

# The Memory Hierarchy -- Goal

Get performance more like this

Bigger

Faster

Registers

~2 KB (~128 B/core)    0.3 ns

640 KB (80 KB/core)    L1 Cache    1.25 ns

10 MB (1280KB/core)    L2 Cache    3.7 ns

30 MB    L3 Cache    17 ns

32 GB    $3.15/GB    Main Memory    80 ns

~200 GB – 1 TB    $0.1/GB    Solid State (flash) Drive    ~.1 ms

2-5 TB    $0.007/GB    Disk Drive    ~3 ms

At size and cost more like this

# Implications

How might this affect how machines are built or how we program? What are the tradeoffs/priorities?

- We saw many factors of 10 or 100 in:
  - Size
  - Performance
  - Price

- "When you see a factor of 100, it's going to affect how you program."
  – E. Kohler

- As the ratios between different parts of the system change, so do our priorities.
  - 1956:
    - $/MB(mem) : $/MB(disk)=> $411M : $9200 => 44,673 X
  - 2024:
    - $/MB(mem): $/MB(disk) => $0.00315 : $0.000007 => 450 X

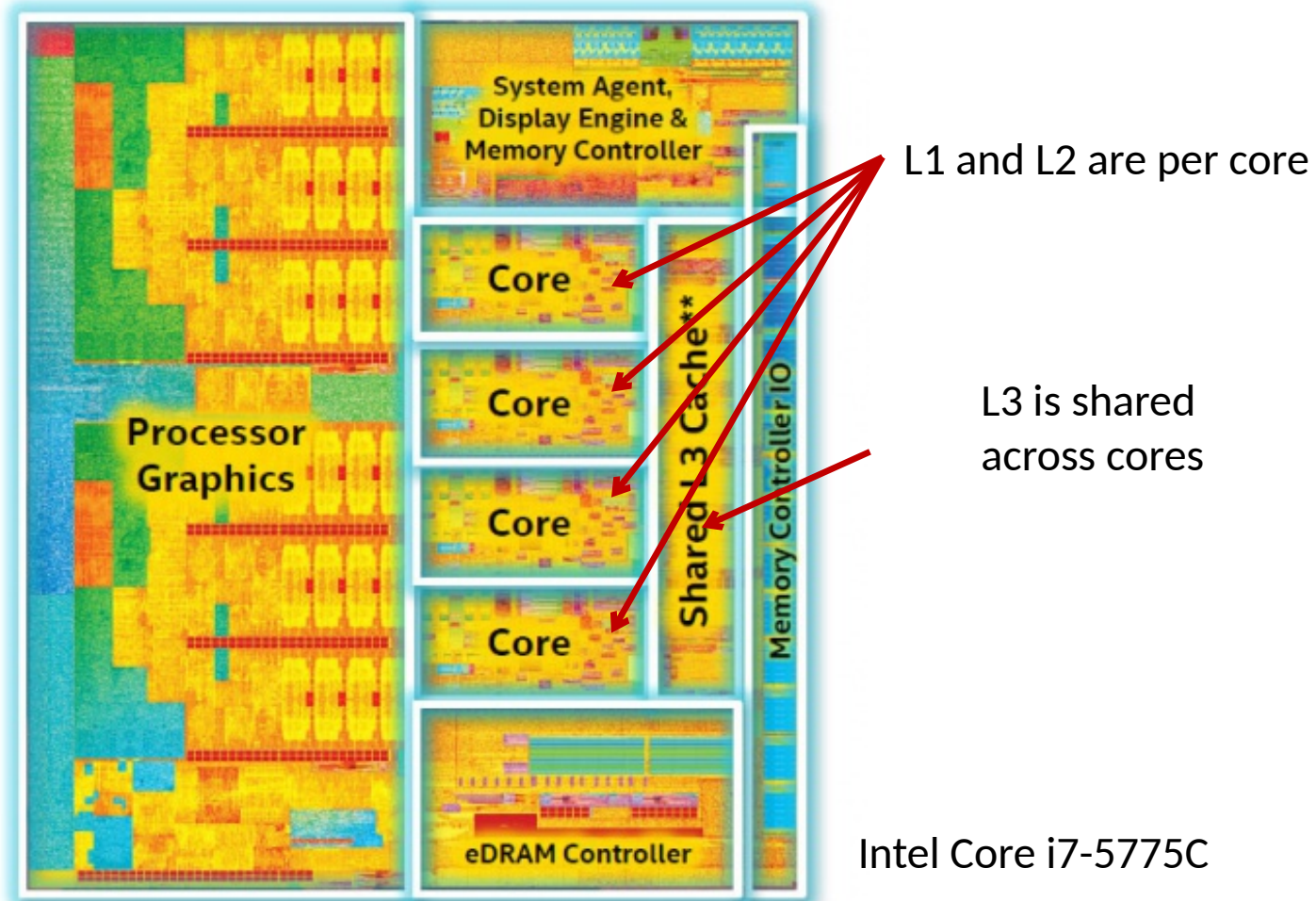THE UNIVERSITY OF BRITISH COLUMBIA

Computer Science
Faculty of Science

# Price of memory over time

- Cost of memory in 1956 versus 2024:
  - 1956 $ / MB : 2024 $ /MB => $411M/.00315 => 130 trillion times cheaper
- Cost of disk in 1956 versus 2015:
  - 1956 $ / MB : 2024 $ /MB => $9200/.000007 => 13 trillion times cheaper

.

# Caching

- Definition:
  - Colloquially: store away in hiding or for future use
  - Applied to computation:
    - Placing data somewhere it can be accessed more quickly
- Examples:
  - In assembly language: we move data from memory to registers.
  - In the hardware: we move data from main memory into memory banks that live on the processor (more on this in a moment).
  - In software: we read things into our program's local buffers and manipulate them there.

THE UNIVERSITY OF BRITISH COLUMBIA

**Computer Science**
Faculty of Science

# Processor Caches



L1 and L2 are per core

L3 is shared across cores

Intel Core i7-5775C

THE UNIVERSITY OF BRITISH COLUMBIA

Computer Science
Faculty of Science

# Inclass Exercise: Caching is *Everywhere*
**(This is _not_ an example of the HW memory caches we'll focus on this month!)**

- You will want to use the student machines, not a workspace!

- Learning objectives:
  - Use both **file system system calls** (open/close/read/write), and **standard IO calls** (fopen/fclose/fread/fwrite)
  - Explain why the following things have enormous impacts on performance:
    1. Buffer size (e.g., how many bytes you send to write or fwrite)
    2. Use of write versus fwrite.

  Use UNIX (Linux) man pages
  - Accessible either via the CLI man command, or
  - Googling "man read"

# Wrapping Up

- Caching is ubiquitous throughout our computing systems:
  - In the processor
  - In the operating system
  - In databases
  - In middleware
  - In applications

- Writing efficient and correct software requires a deep understanding of caching and its implications.

THE UNIVERSITY OF BRITISH COLUMBIA
Computer Science
Faculty of Science