# Lecture 3: Single-processor Computing Summary
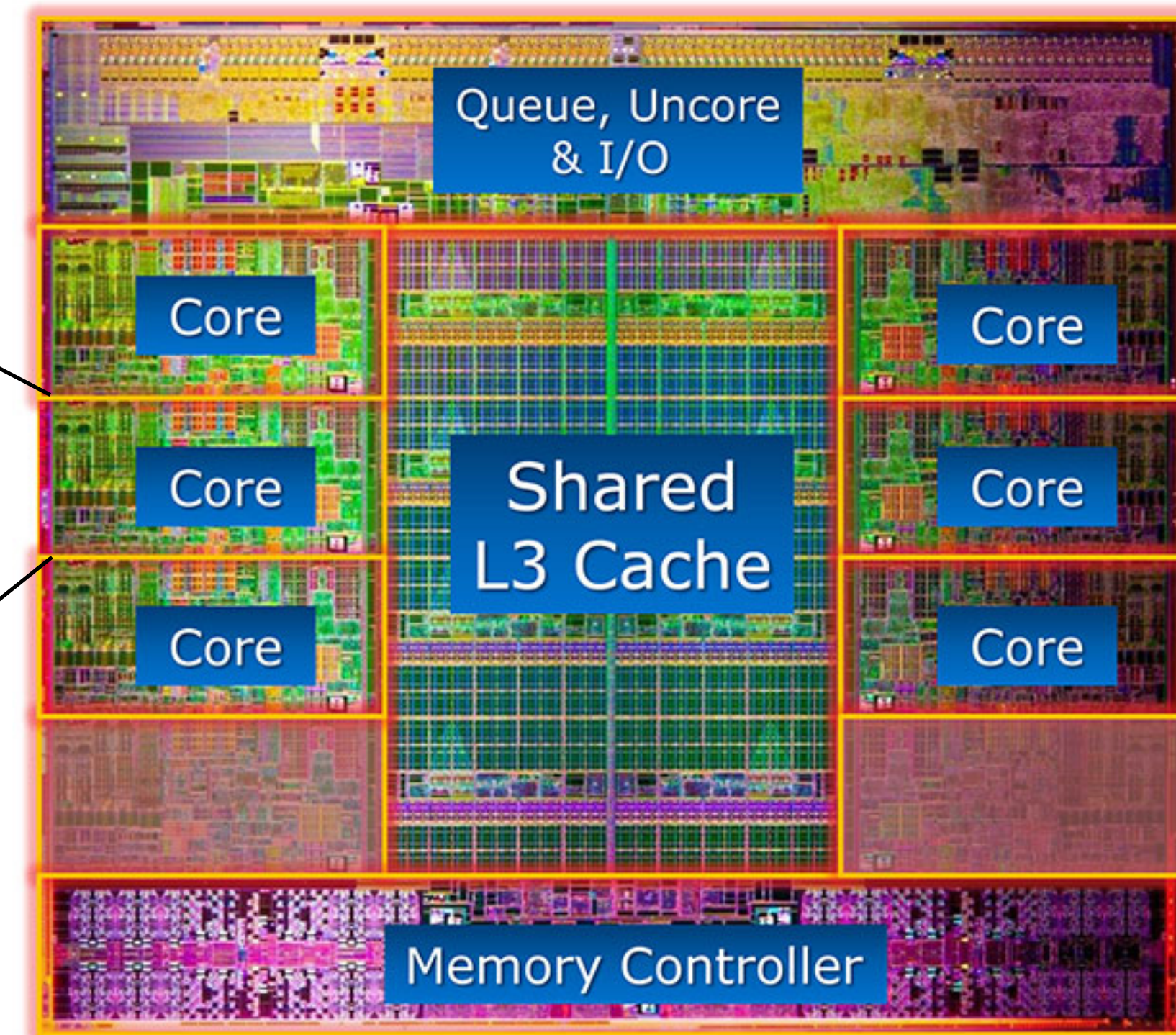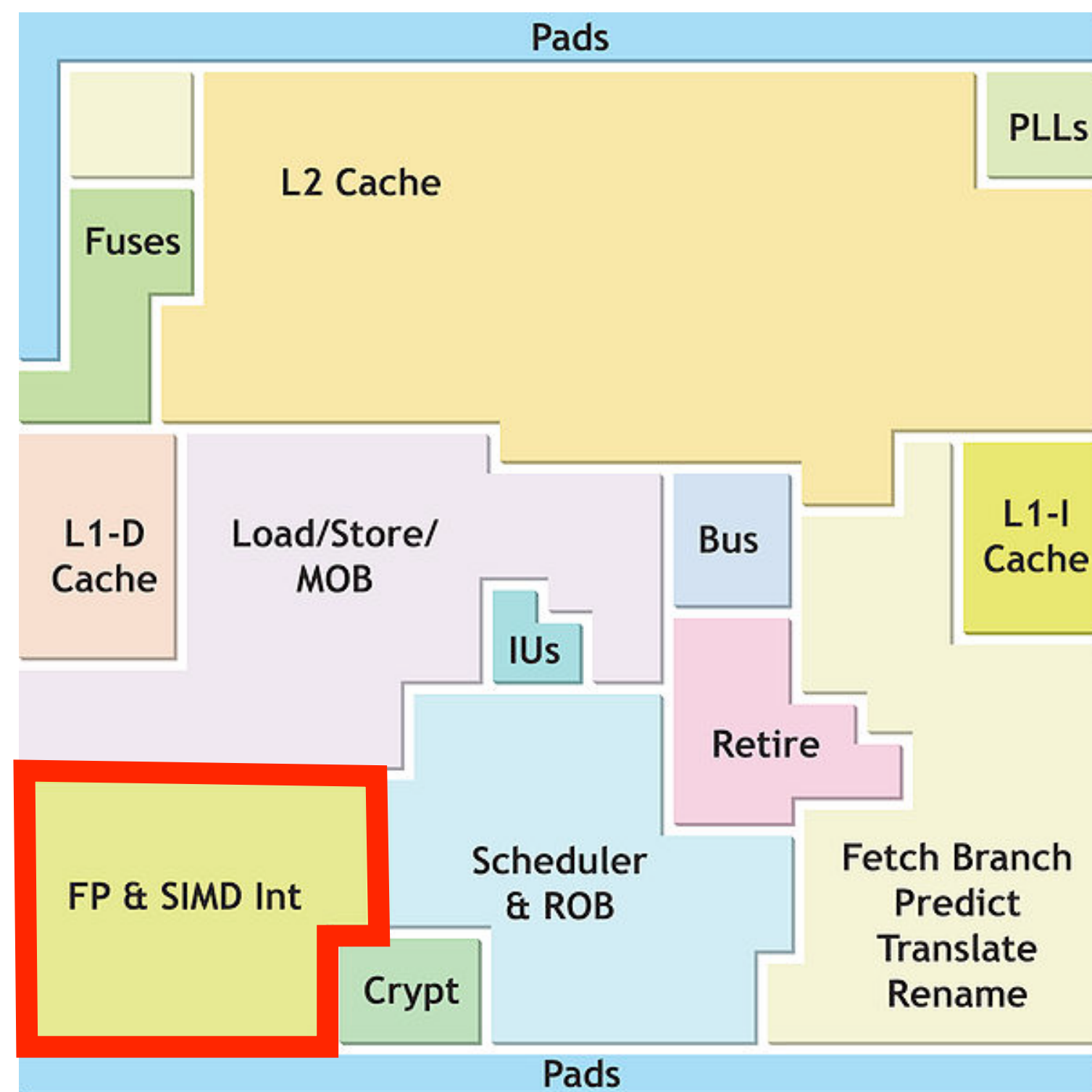
## CMSE 822: Parallel Computing
## Prof. Sean M. Couch

# Anatomy of a Computation
## A CPU

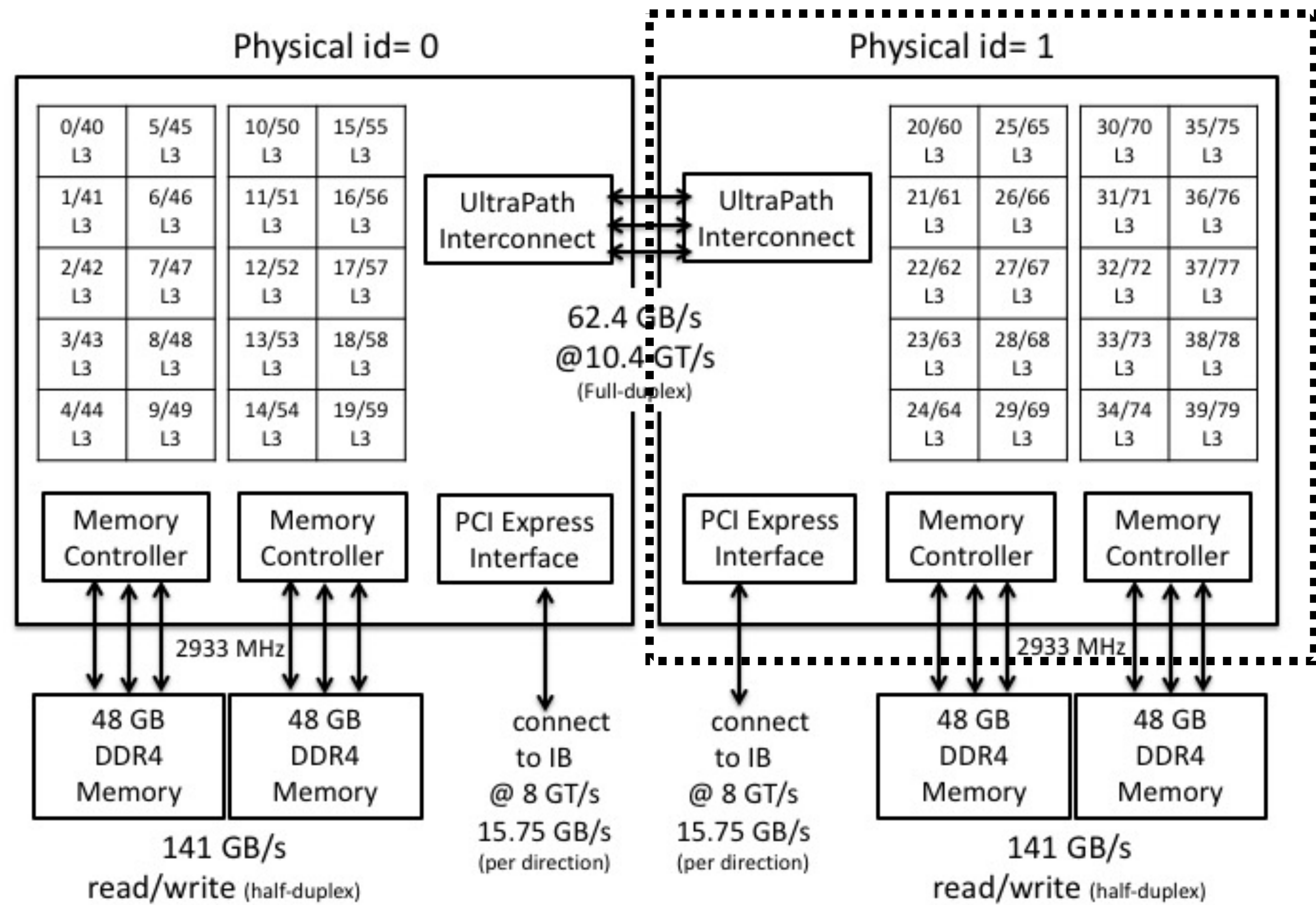Intel® Core™ i7-3960X Processor Die Detail



Single-CPU computing is parallel!

# Anatomy of a Computation
## A Node



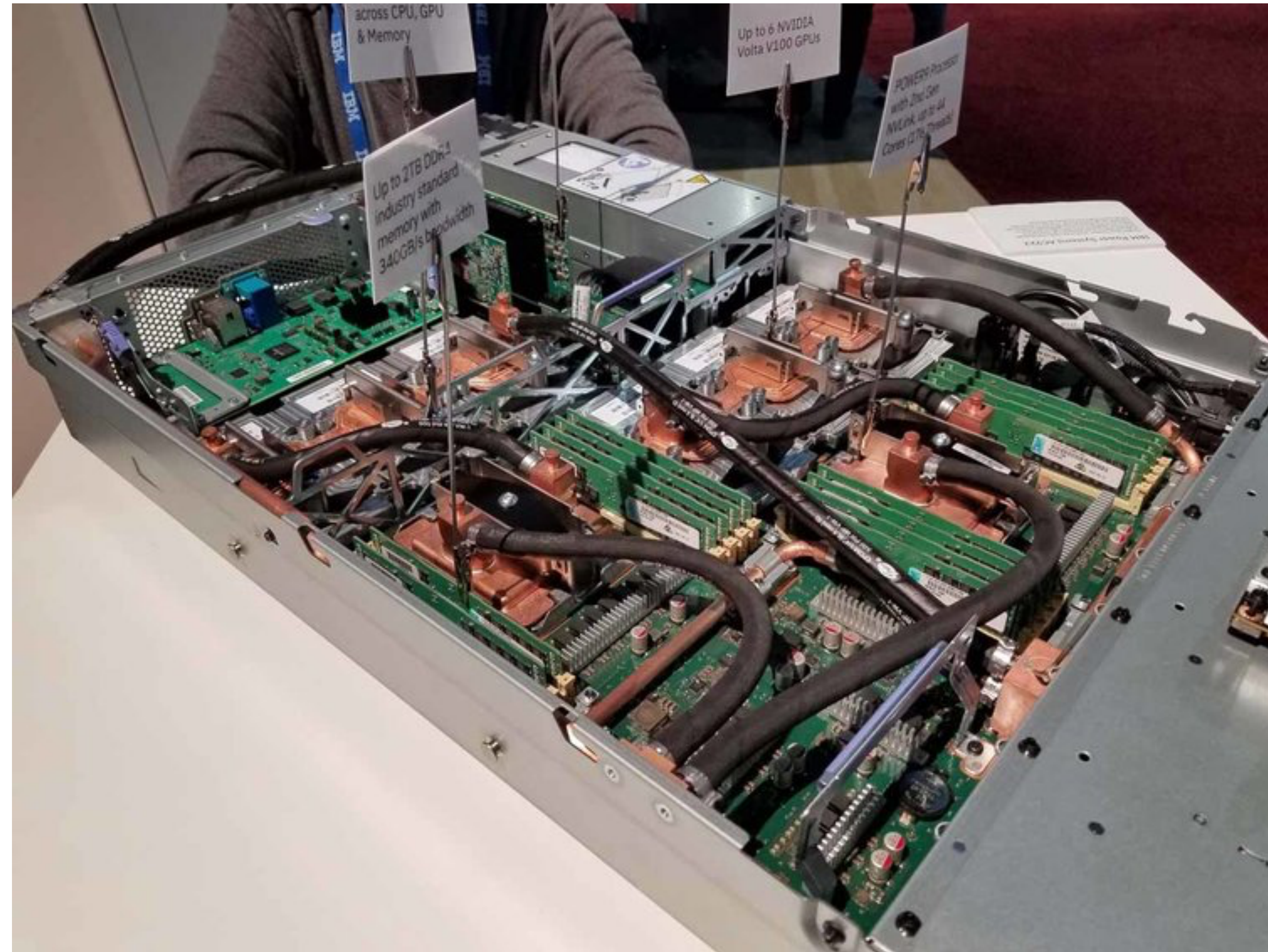Configuration of a Cascade Lake - SP Node

"socket"

# Anatomy of a Computation
## A Node



Summit, ORNL

# Anatomy of a Computation
## A Node



Summit, ORNL

# Anatomy of a Computation

## Node-to-node Interconnect
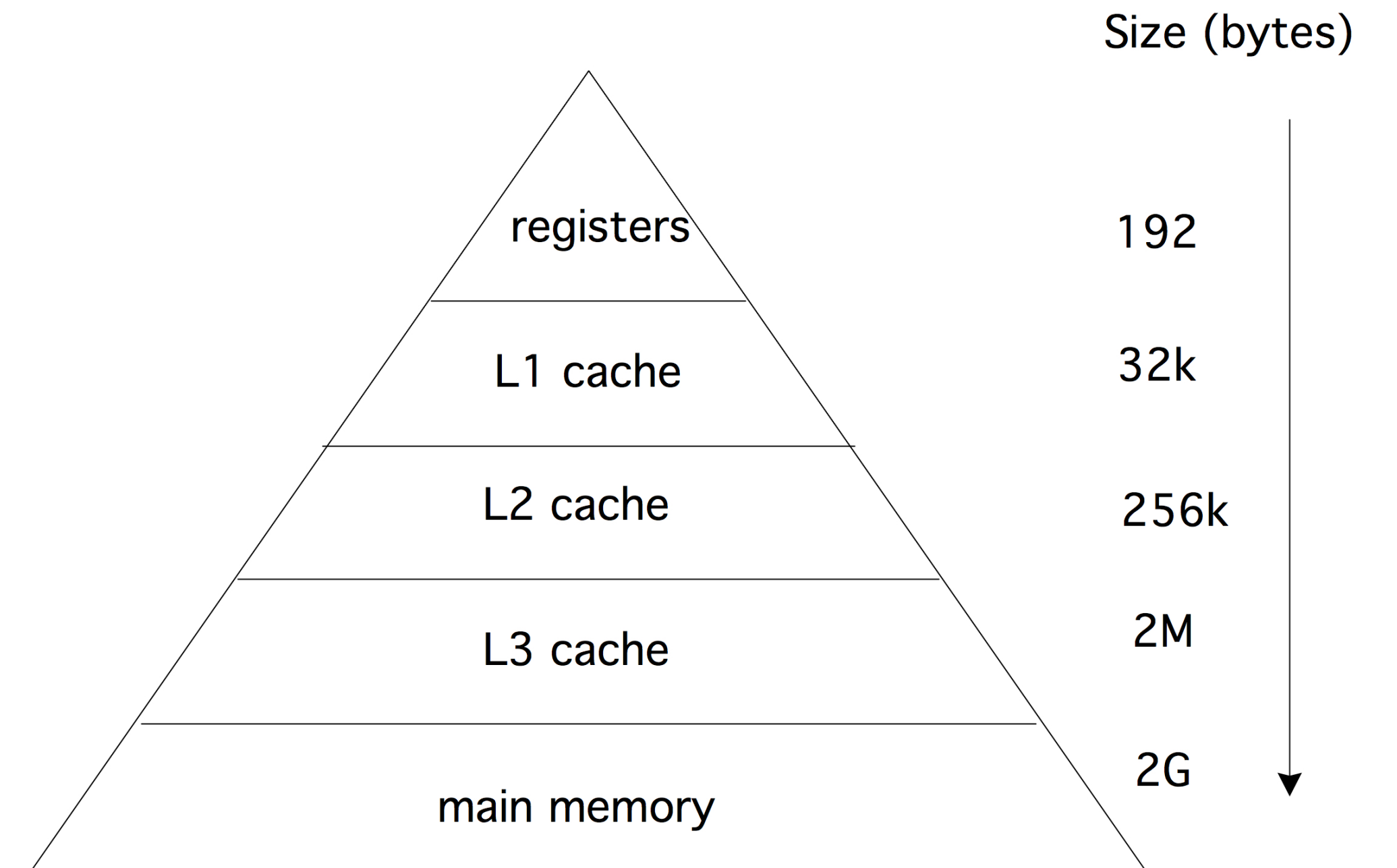
# Anatomy of a Computation
## Cluster



Summit, ORNL

# Memory hierarchy
## Often the limiter of performance…

Registers          Registers

L1 cache           L1 cache

L2 cache

Memory

Latency from next
level (cycles)

4

12

26

230-360

Size (bytes)

registers          192

L1 cache          32k

L2 cache          256k

L3 cache          2M

main memory       2G

# Memory hierarchy
## Need to feed the beast (er, CPU)



latency

bw

Mem

Proc

Independent data

Little's Law: Concurrency = Bandwidth x Latency

# Memory hierarchy
## Absolute unit



What is the fundamental unit of memory movement?

a. page

b. word

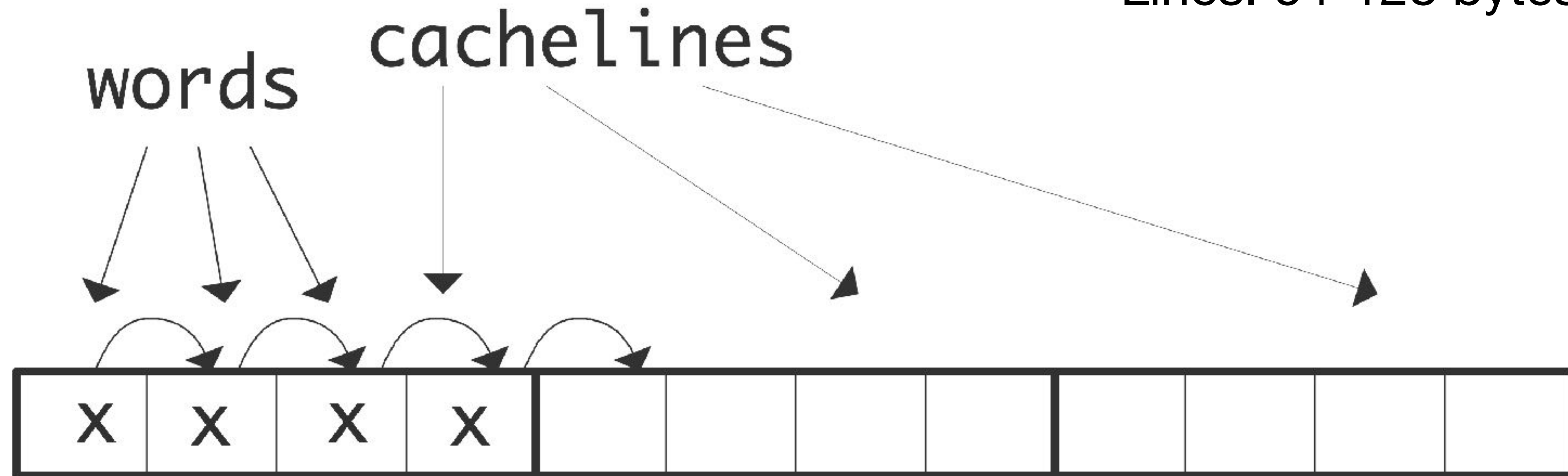c. line

d. byte

# Memory hierarchy
## Cache line

Lines: 64-128 bytes
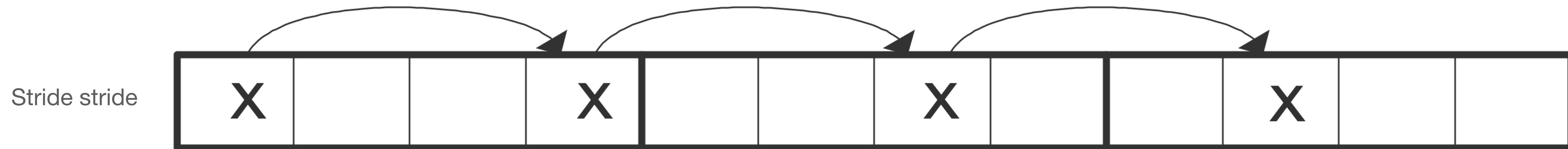
words    cachelines
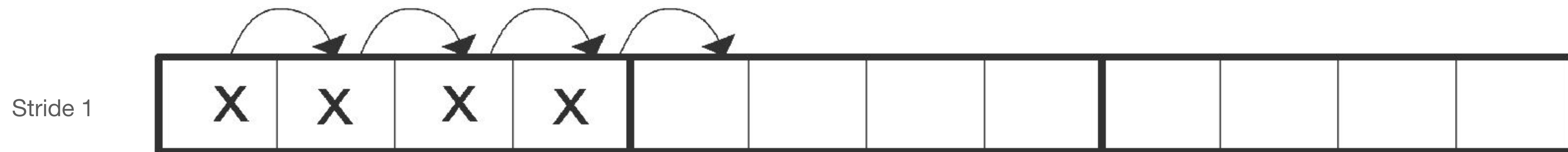
X   X   X   X

Words: usually 64 bits

# Memory hierarchy
## Strided access



Stride stride

```
for (i=0; i<N; i+=stride)
    ... = ... x[i] ...
```

Stride 1

```
for (i=0; i<N; i++)
    ... = ... x[i] ...
```

# Memory hierarchy
## Reuse is key to performance!

- Compulsory cache miss: first time memory is referenced

- Capacity cache miss: cache not big enough to fit problem

- Conflict cache miss: data mapped to same cache location as another

- Invalidation cache miss: another core changed value at memory address

# Memory hierarchy
## False sharing

```
    local_results = new double[num_threads];
#pragma omp parallel
{
  int thread_num = omp_get_thread_num();
  for (int i=my_lo; i<my_hi; i++)
    local_results[thread_num] = ... f(i) ...
}
global_result = g(local_results)
```

- Cores access and alter data in same *cache line*

# Exercise 1.14
## Matrix-matrix Multiply

**E**xercise 1.14.    The matrix-matrix product, considered *as operation*, clearly has data reuse by the above definition. Argue that this reuse is not trivially attained by a simple implementation. What determines whether the naive implementation has reuse of data that is in cache?

Caches can only hold a finite amount of data. Once a row of A and a column of B take up more than the size of the cache, their elements will be flushed between iterations of the outer loop.

# Project 1
## Group work