

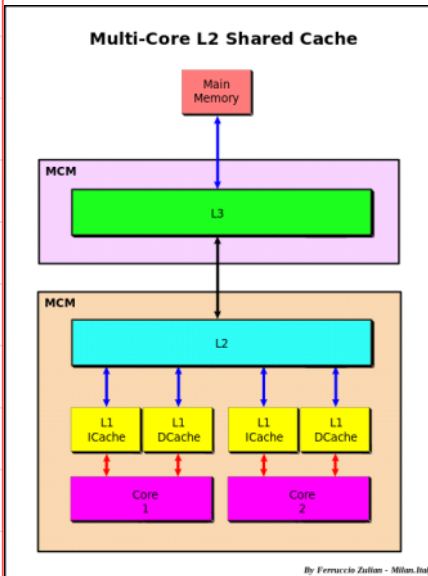
Memory models in a computer

(Simplified for mathematicians)

Friday, August 23, 2024 4:51 PM

- What scientists must know about hardware to write fast code, see <https://viralinstruction.com/posts/hardware/>
- https://book.sciml.ai/lectures/#optional_extra_resources

Memory Hierarchy



Slow (but large)

↑ transferring memory is a major slowdown
↓

Fast (but small)

Transferring memory to another computer or to a GPU is also slow

Cache misses

Transferring data is slow, has both a latency and bandwidth restrictions.

i.e. cost to move "x" bytes is $b \cdot x + l$
in time or clock cycles

So... due to latency, we don't like to transfer small amounts of data (since it's wasteful)

Instead, transfer a large block what you need in cache

RAM



Cache



whole block is transferred.

Computer tries to predict what memory it will use in the future.

Memory (page 2)

Wednesday, September 11, 2024

7:16 AM

whenever you need to move memory into cache (because it wasn't there already) it's a **cache miss** and really slows things down.

Effect

Ex: Storing a large matrix

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

What's this look like in memory?

Python does row-major order

Matlab/fortran does

column-major order

$$[a \ c \ b \ d]$$

Let's say our cache can fit 2 numbers at once.

We compute $\|A\|_F^2 = \sum_{i=1}^2 \sum_{j=1}^2 A_{ij}^2$ $i = \text{row}$
 $j = \text{column}$

$\underbrace{\sum_{j=1}^2}_{\text{we can choose the order}}$

Say we're in Python

$$\sum_i \left(\sum_j A_{ij}^2 \right)$$

$$i = 1 \quad \begin{matrix} j = 1 \\ j = 2 \end{matrix}$$

a
b

cache miss

(Request "a", get [a,b])

$$\begin{array}{cc} i=2 & j=1 \\ \hline & j=2 \end{array}$$

$$\frac{c}{d}$$

cache miss

(Request "c", get [c,d])

vs.

$$\sum_j \left(\sum_i A_{ij}^2 \right)$$

$$j=1 \quad i=1$$
$$i=2$$

cache miss get [a,b]

cache miss gut $[c, d]$

$$j=2 \quad i=1$$

cache miss get [b, c]

$i = 2$ d

cache miss für $[d, ?]$

Even more important

(and complicated) for sparse matrices

4 cache misses, bad!

Related to idea of **blocked algorithms** (= works on contiguous chunks of data)

and exploiting fast **BLAS** for vector and matrix operations (multiplies)

- Intel MKL, very optimized, already parallelized

and to **cache oblivious algorithms** (Frigo et al. 99) where you don't need to explicitly know cache size in order to get efficient performance

Memory and GPUs (page 3)

Wednesday, September 11, 2024

9:33 AM

In "main memory" (aka RAM), one physical location split into two (virtual) types: **stack** (includes call stack, and small amount of room for variables)
heap
= the rest, used for large data
ex. in C, any **malloc** **calloc** calls

SIMD = Single Instruction Multiple Data

In most modern processors (**CPU** mmx, sse, sse2, avx ...)

If you want to apply $f(x) = x^2$ to multiple data $\{x_i\}_{i=1}^{256}$
or anything else simple enough

a SIMD CPU does it efficiently, much less than $256 \times$ cost of doing it once. It's a hardware thing, less overhead

They excel at matrix multiplies...
the core of much Science and ML

GPU = Graphical Processing Unit, aka video/graphics card

Like CPU but can only do simpler things, slower clock speed, but huge number (eg. 1000s) of parallel cores.

Like extreme SIMD, though at a higher level

- NVIDIA's **CUDA** is dominant driver/language

amd a100 on CURC Alpha are NVIDIA A100's

Not all GPUs (even NVIDIA ones) are CUDA-compatible

- **OpenCL** and AMD's **ROCm**
and Mac **MPS** (Metal Performance Shader)
are alternatives though not as widespread

gaming GPUs often don't have much memory, or only support low-precision

ami 100 on CURC Alpha are AMD MI100's

- Data transfer from CPU to GPU is costly
- GPUs have small RAM compared to CPUs
- GPUs usually work in single (32bit) floating point precision, or less (vs. CPU / NumPy standard is double precision, 64 bit)