

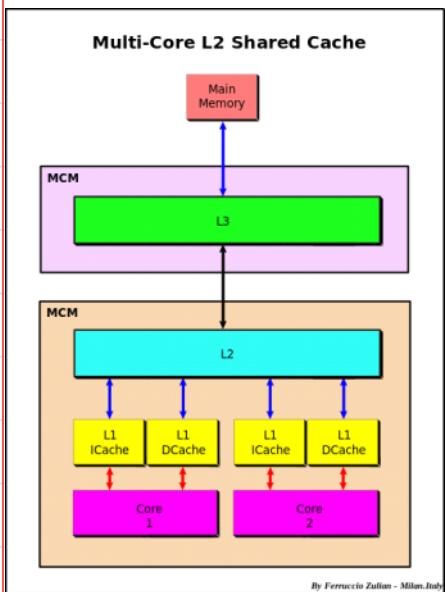
Memory models in a computer

(Simplified for mathematicians)

Friday, August 23, 2024

4:51 PM

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

RAM



Cache

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.

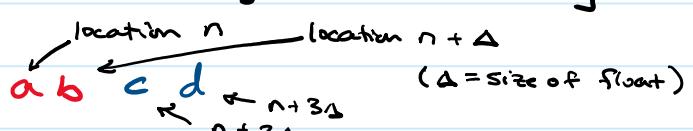
Effect

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

$$\begin{bmatrix} a & c & b & d \end{bmatrix}$$

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

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

we can choose the order

Say we're in Python

$$\sum_i \left(\sum_j A_{ij}^2 \right) \quad \begin{array}{ll} i=1 & j=1 \\ & j=2 \\ i=2 & j=1 \\ & j=2 \end{array} \quad \begin{array}{ll} a \\ b \\ c \\ d \end{array} \quad \begin{array}{l} \text{cache miss} \\ (\text{Request "a", get [a,b]}) \\ \text{cache miss} \\ (\text{Request "c", get [c,d]}) \end{array}$$

vs.

2 cache misses... good!

$$\sum_j \left(\sum_i A_{ij}^2 \right) \quad \begin{array}{ll} j=1 & i=1 \\ & i=2 \\ j=2 & i=1 \\ & i=2 \end{array} \quad \begin{array}{ll} a \\ c \\ b \\ d \end{array} \quad \begin{array}{l} \text{cache miss get [a,b]} \\ \text{cache miss get [c,d]} \\ \text{cache miss get [b,c]} \\ \text{cache miss get [d, ?]} \end{array}$$

4 cache misses, bad!

Even more important
(and complicated) for **sparse matrices**

Related to idea of **blocked algorithms** (= works on contiguous chunks of data)
and exploiting fast **BLAS** for vector and matrix operations (multiples)
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 9:38 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) and **heap** = the rest, used for large data ex. in C, any **malloc** or **calloc** calls

SIMD = Single Instruction Multiple Data

CPU
In most modern processors (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. 1000's) of parallel cores.

Like extreme SIMD, though at a higher level
it's "SPMD", Single Program M.D.

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

amc
a100 on CURE Alpine
are NVIDIA A100's

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

- OpenCL and AMD's ROCm

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

and Mac MPS (Metal Performance Shader)

are alternatives though not as widespread

am100 on CURE Alpine 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, 64bit)

Tensor Processing Unit (ASIC: application specific integrated circuit)

Google's hardware specialized for linear algebra

Uses 16 bit floating point numbers for speed,

with tiny mantissa... lots of catastrophic cancellation

but... implements a FMA (fused multiply add)

which does tricks so as to not loose much precision

$$a \leftarrow a + (b \times c)$$

I don't know how, but could be similar to Kahan summation

Uses an extra variable to
get extra accuracy

aka compensated summation

Goal: $S = \sum_{i=1}^n x_i$

Algo: $\text{sum} = 0$

$$c = 0$$

for $i = 1, \dots, n$

$$y = x_i - c$$

$$t = \text{sum} + y$$

$$c = (t - \text{sum}) - y \quad // = 0 \text{ in exact arithmetic}$$

$$\text{sum} = t$$

Return sum

Since ~2015, meant to work

w/ TensorFlow mostly.