

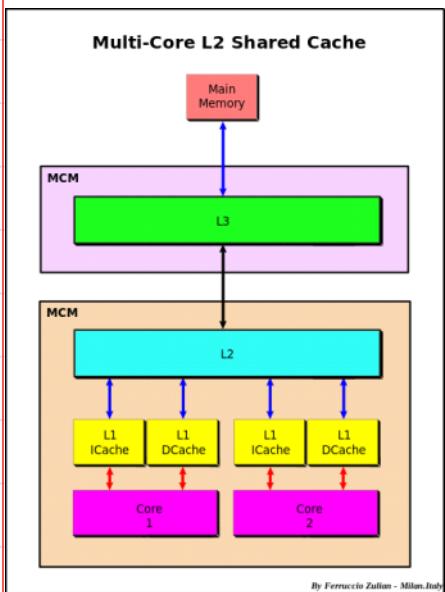
Memory models in a computer

(Simplified for mathematicians)

Friday, August 23, 2024

4:51 PM

Memory Hierarchy



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

Learning Objectives: Memory, not flops, often determines speed of an algorithm

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

What's this look like in memory?

Python does **row-major order**

Matlab/Fortran does

column-major order

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

location n location $n + \Delta$
 a b c d
 $n + 2\Delta$ $n + 3\Delta$
($\Delta = \text{size of float}$)

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

(RAM is linear, so must make the 2D array linear)

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

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

we can choose the order

Say we're in Python

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

vs.

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

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

2 cache misses... good!

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 (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.

Sparse matrices

Wednesday, September 3, 2025 9:22 AM

Mathematically, a **Sparse matrix** is one with "a lot" of zeros.

$$m \boxed{A}$$

"nnz" = number of nonzeros, so if $\text{nnz}(A) \ll m \cdot n$, we say A is sparse.

This is vague of course!

On a computer, if your matrix is sparse but you don't do anything

special, the computer treats those zeros like any other number. If you want to exploit sparsity, you must explicitly use a sparse data structure.

⚠ You need a matrix to be very sparse before it's worth it to exploit!

Sparse matrix data structures

• Baseline: Store $m \cdot n \cdot \text{sizeof}(\text{float})$ bits ($+ 2 \cdot \text{sizeof}(\text{int})$ for metadata)

Location is implicit. row-major: $A = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]$
via order

• Naive Format: "COO", Coordinate, 1-based indexing

Store (i, j, A_{ij}) triplets

e.g. $(2, 1, a), (4, 1, b), (1, 2, c), (2, 2, d), (3, 2, e), (3, 3, f)$.

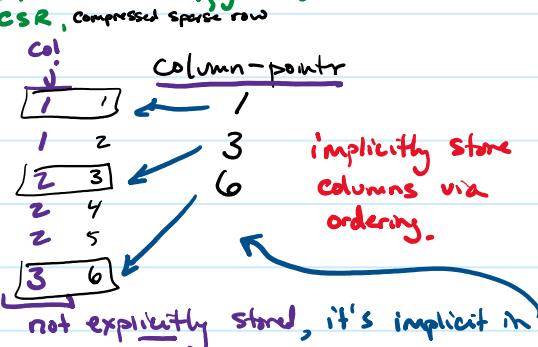
$$A = \begin{bmatrix} 0 & c & 0 \\ a & d & 0 \\ 0 & e & f \\ b & 0 & 0 \end{bmatrix}$$

Cost: $\text{nnz}(A) \cdot (\text{sizeof}(\text{int}) + \text{sizeof}(\text{float}))$

• Compressed Sparse Column (CSC) synonyms
• Compressed Column Storage (CCS) or CRS, compressed sparse row

Store

value A_{ij} :	row i :
a	2
b	4
c	1
d	2
e	3
f	3



See Demo

Cost: $\text{nnz}(A) \cdot (\text{sizeof}(\text{int}) + \text{sizeof}(\text{float})) + n \cdot \text{sizeof}(\text{int})$ (+ metadata)

Since usually $\text{nnz}(A) > n$, this is less storage.

Implementation (due to data locality + structure)

Accessing a column of a CCS matrix is way faster than accessing a row.

In Python's `scipy.sparse` package, easy to convert formats

Ex. of Storage

$M = n = 10,000$, double precision dense storage is 760 MB

if 10% sparse, CCS storage is 145 MB

(if 100% sparse, i.e. not sparse at all, it'd be 1450 MB... worse than dense!)

HDF5 (Hierarchical Data Format 5)

Wednesday, September 3, 2025 9:22 AM

Saving a dataset to disk.

- Considerations:
- **documentation, metadata** (variable names...)
Often nice to store several variables in 1 file
 - **reusable by others** Use a standardized format, not homegrown, not proprietary
 - **allow compression**
exploits redundancies. Almost all data can be compressed.
Compression can be lossless, though lossy compression
obviously gives even greater savings
 - **be quick to read/write** binary (not text) files,
and allow you to load the data starting in the
middle!
i.e. $x = np.ones(10^9)$ stored in your file.
If you just want $x[1000:1001]$ you don't
want to have to load in all the data
(may even not fit in RAM!)

HDF5 is an example of a well-known, standardized data format
that addresses the above considerations. There are python libraries
that make it easy to use (better than "pickle")

Matlab has its own format which is a variation of HDF5, and
is generally interoperable w/ HDF5