

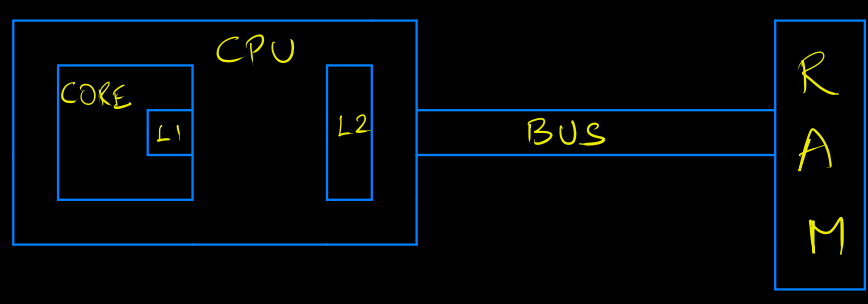
MEMORY SYSTEM PERFORMANCE - CACHING

Cache : Collection of memory locations

Resides On-core , On-chip

what? High speed access^{~1ns} (Vs RAM^{100ns})

→ Cannot be large size (Vs RAM)



Instruction Cache & Data Cache

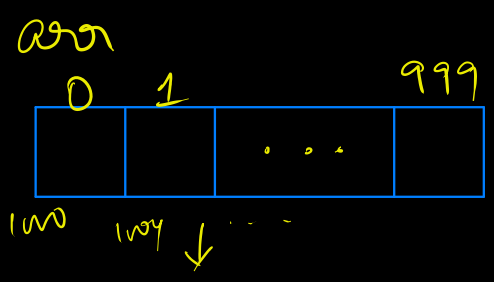
→ Principle of Locality

why?

→ Sequential Execution

Branch Instructions?

→ Iterative data access



```

float arr[1000], sum = 0.0f;
int i = 0;
// arr values are computed here
for(i = 0; i < 1000; i++)
    sum += arr[i];
    
```

Contiguous Memory locations

⇒ Access of one location is followed by an access of a nearby location ⇒ Spatial locality,
in the nearby future ⇒ Temporal locality

How?

One Cache line = Multiple Memory blocks
(e.g. 8 to 16)

Let Memory block = 4 bytes (= bus width)

Let Cache line = 8 blocks

⇒ One cache line = 32 bytes

For single level cache,

When `arr[0]` is accessed

↳ check if available in Cache

HIT
(Proceed)

MISS
→ Fetch Cache line
Containing `arr[0]`
from memory into
Cache

→ Extend the idea
to multi-level hierarchy

→ Execution stalls

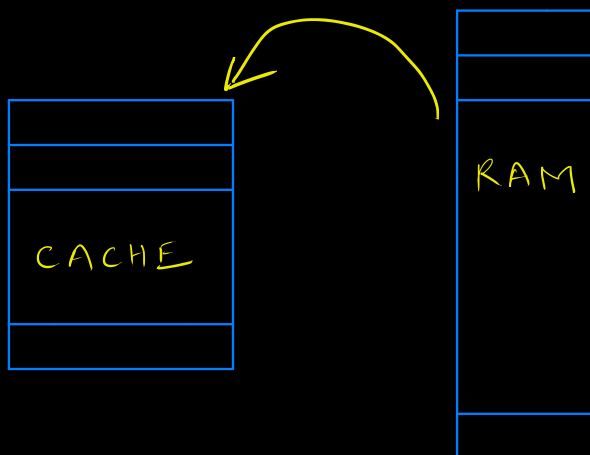
→ But what about
the next access?

Issues in Cache Design

I → Writes can cause inconsistency.

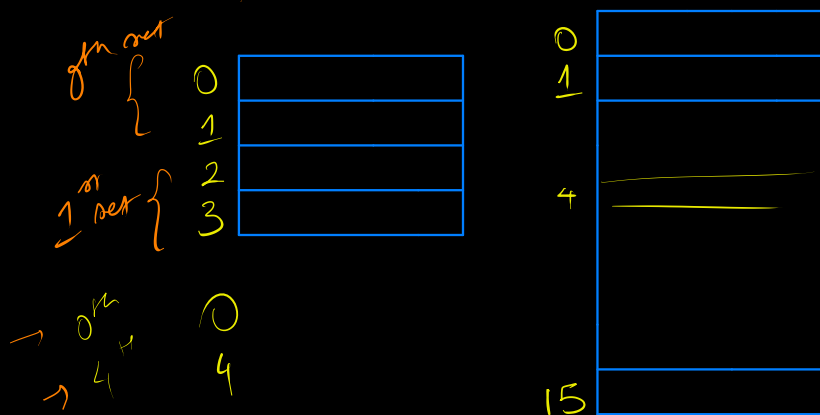
Write-through Write-back

II →
EVICTION?



Memory blocks to
Cache line mapping

Assumption: 1 cache line = 1 Memory block



Direct Mapping:

$$CL = MB \% 4$$

New line placed at CL

n-way Set Associative Mapping:

$$NS = \text{No. of CL} / n$$

$$SN = MB \% NS$$

(No. of sets)

Fully Associative Mapping:

New line can be placed
at any location in cache

↓
Cache Replacement
Strategy
(LRU, LFU, ...)

I. Processor operating at 1 GHz

DRAM latency 100 ns

No. of FMA units = 2

FMA operation completion = 1 cycle

Peak performance = 4 GFLOPS

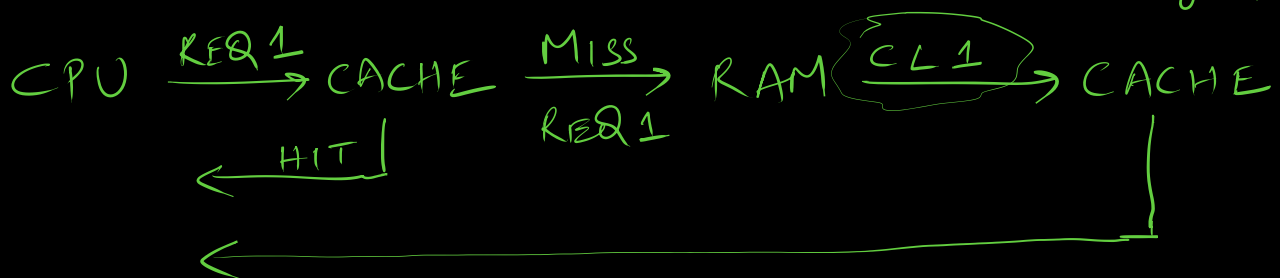
Cache size = 1 KB

Cache line size = 16 bytes

Cache latency = 1 ns

Data bus speed = 100 MHz (clock = 10 ns)

Data bus width = 1 word (4 Bytes)



```
float a[96], b[96], ans = 0.0f;
```

```
int i = 0;
```

```
// a & b are set here
```

```
for (i = 0; i < 96; i++)
```

```
    ans = ans + a[i] * b[i];
```

How many FLOPS can this program achieve?
or what is the peak speed of this algorithm?

CPU accesses $a[0] \rightarrow$ HIT or MISS?

\downarrow

Fetches $a[0], \dots, a[3]$
from RAM to CACHE

100 ns
 $+ 10 \text{ ns} \times 3$
 $= 130 \text{ ns}$

CPU accesses $b[0], \dots, b[3] \Rightarrow 130 \text{ ns}$

\Rightarrow In 260 ns $\Rightarrow a[0], \dots, a[3]$
& $b[0], \dots, b[3]$

NO CACHE MISS till $a[4]$ or $b[4]$

2 FMA units : 2 cycles to complete computation
(cache access time generally ignored) $= 2 \text{ ns}$

\Rightarrow 8 FLOPs in total 262 ns

[Same pattern continues for $4, \dots, 7; \dots; 92, \dots, 95 - 24 \text{ times}$]

Peak Computation Rate of the algorithm $= \frac{24 \times 8}{24 \times 262} = \frac{8}{262}$

% of peak CPU rating $= \frac{30.53 \times 10^6}{4 \times 10^9} \times 100 = 30.53 \text{ MFLOPS}$
 $= 0.76 \% \quad (\text{better than } 0.25 \%)$

II. Processor speed = 1 GHz (+ 2 FMA in 1 cycle)

DRAM Latency = 100 ns

Cache Size = 32 KB

Cache Latency = 1 ns

Assumption:
[Cache line size = 1 word
Bus width = 1 word]

Let us look into calculation of $C[0][1]$:

Already 0^{th} row of A is in cache!

\Rightarrow After $100\text{ns} \times 32$ ($B[0][1], \dots, B[31][1]$), we can compute $C[0][1] \Rightarrow 32 \text{ FMA}$ can be done.

Similarly for each $C[0][j]$, for $j \leftarrow 0$ to 31

In general, fetch 0^{th} row of A : 3200 ns

fetch 0^{th} col of B : 3200 ns

\vdots

fetch 31^{st} col of B : 3200 ns

$$\text{Total} = 3200 + 3200 \times 32$$

After that time, we have 0^{th} row of C with 32×32 FMA operations.

Does the pattern change for calculating 1^{st} row of C ?

Fetch 1^{st} row of A : 3200 ns

All cols of B are already available in Cache!

\Rightarrow After 3200 ns , 1^{st} row of C can be computed with 32×32 FMA operations.

\Rightarrow In general, for i^{th} row of C , $i \leftarrow 1$ to 31 , can each be calculated after 3200 ns each.

$$\therefore \text{Overall time} = 3200 + 3200 \times 32 + 3200 \times 31$$

$$32 \times 32 \times \frac{100}{2}$$

$$= 3200 \times 64 = 204800 \text{ ns}$$

$$\begin{aligned}\text{Total Computations} &= 32 \text{ FMA} \times (32 \times 32) \\ &= 32^3 \text{ FMA} = 32^3 \times 2 \text{ FLOPs}\end{aligned}$$

$$\text{Time taken for computation} = \frac{32^3 \times 2 \times 1 \text{ ns}}{4 \text{ (2 FMA per cycle)}} = 16384 \text{ ns}$$

$$\therefore \text{Fetch Time + Compute time} = 221184 \text{ ns}$$

(cache access time ignored)

$$\therefore \text{Peak Computation Rate of the algorithm} =$$

$$\frac{32^3 \times 2}{221184 \times 10^{-9}} = 296.3 \text{ MFLOPS}$$

$$\% \text{ of Peak performance} = \frac{296.3 \times 10^6}{4 \times 10^9} \times 100 = 7.4 \%$$

How would this change if Cache line size = 16 bytes?

Another way to calculate : Cache HIT Ratio

$$\text{I. 1 DRAM access for every 4 words : CHR} = \frac{3}{4}$$

$$\Rightarrow \text{Average Memory Access time} = \frac{3}{4} \times \frac{1}{\text{(ns)}} + \frac{1}{4} \times \frac{100}{\text{(ns)}} = 25.75 \text{ ns/word}$$

Dot product : 8 FLOPs after fetching 8 words

$$\Rightarrow \text{Average Computation per word} = \frac{8}{8} = 1 \text{ FLOPs/word}$$

$$\Rightarrow \text{Average Computation Rate} = \frac{1 \text{ FLOPs/word}}{25.75 \text{ ns/word}} = 38.83 \text{ FLOPS}$$

II. 1 DRAM access for every word : $CHR = \frac{0}{1} = 0$

$$\begin{aligned}\text{Average Memory access time} &= 0 \times \frac{1}{(\text{ns})} + 1 \times \frac{100}{(\text{ns})} \\ &= 100 \text{ ns/word}\end{aligned}$$

Matrix Multiply : $2n^3$ FLOPs after fetching $2n^2$ words

$$\text{Average computation per word} = n \text{ FLOPs/word} = 32 \text{ FLOPs/word}$$

$$\therefore \text{Average Computation Rate} = \frac{32}{100} = 320 \text{ MFLOPS}$$

\Rightarrow Comparing dot product Vs Matrix Multiply,

0.76% Vs
(16 byte cache line)

7.4%
(even with 4 byte cache line)

What is increasing an algorithm's
Peak Computation Rate?

Data Re-use?