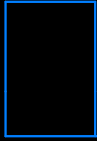


MEMORY SYSTEM PERFORMANCE



WARDEN

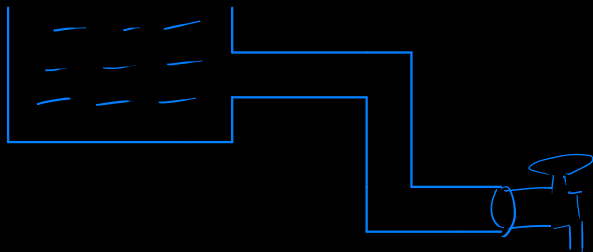


If corridor is getting mopped?

If student is sleeping?

If Corridor is only One-way?

Water Tank to Tap



Higher the tank?

Wider the pipe?

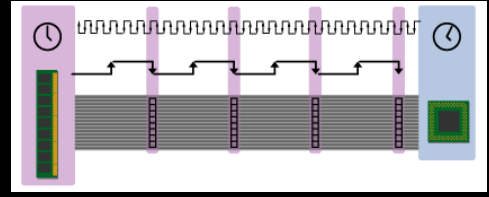
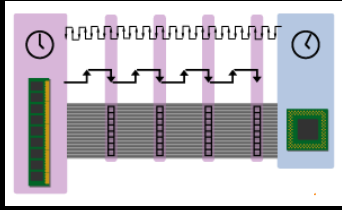
Pumping motor at top?

Conveyor Belt



Faster the movement?
(Quick slots)

Larger slots?



CPU: Request Control
Wait for grant
Send data
Relinquish Control

RAM: Receive Request
Prepare data
Request Control
Wait for grant
Send data
Relinquish Control

Time taken between CPU's Request for data on the bus to its arrival on the bus to the CPU

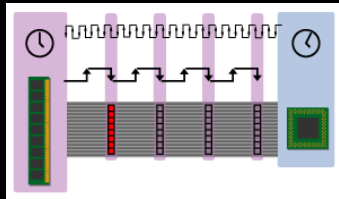
LATENCY (in units of time)

Address Note that it is char c[10];
5, 11, 3 bytes

If latency = 3 bus cycles \Rightarrow

CPU needs to minimum wait for 3 bus cycles
(Many CPU cycles)

0-7 - 1 Blk
8-15 -



$$(ID) \frac{100}{8} \times 8 = 96$$

CPU: Ask a byte

\rightarrow (Should I ask for more?)

RAM: Sends the block containing the byte

\rightarrow (Should I predict and send more?)

1 Block = 8 Bytes

100th byte \Rightarrow RAM: 8 Bytes
Which 8 bytes?

96, 97, 98, 99, 100, 101,
102, 103

STRUCTURE PADDING

Rate at which data can be moved between the CPU & RAM : BANDWIDTH (in units of bytes/sec)

100 MHz, 64-bit bus : At one cycle, 8 bytes sent
 $100 \times 10^6 \text{ in a sec} \Rightarrow 8 \times 100 \times 10^6 \text{ bytes/s}$
800 MB/s

133 MHz, 32-bit bus : At one cycle, 4 bytes
 $133 \times 10^6 \text{ in a sec} \Rightarrow \underline{532 \text{ MB/s}}$

Narrow Corridor Vs Wide Corridor
One Corridor Vs Multiple Corridors
32-bit Vs 64-bit
Single Block Vs Multiple Blocks

Bus width \rightarrow Bandwidth \checkmark
 \rightarrow Latency \times [TRAFFIC FLOW]
Memory blocks \rightarrow Bandwidth \checkmark
 \rightarrow Latency \times
Technology \rightarrow Bandwidth \checkmark
 \rightarrow Latency \checkmark

If Latency of an interconnect is l seconds and its bandwidth is b bytes/sec, then time taken to transmit n bytes

= First 8 $\rightarrow l \text{ sec}$, Next 8 $\rightarrow l + \text{one bus cycle}$
 $l + (n-8)/8 \times \text{one bus cycle}$

$$= l + n/b$$

100 bytes : (0-7) 8 bytes $\xrightarrow{0^{th} bc}$ l sec

(8-15) $\xrightarrow{1^{st} bc}$ ~~l~~ 1 bc (l+1)

$\frac{92}{8} = \underline{\underline{11}}$ (16-23) $\xrightarrow{2^{nd} bc}$ 2 bc (l+2)

\vdots

(80-87) $\xrightarrow{10^{th}}$ l + 10

Time taken (88-95) $\xrightarrow{11^{th}}$ l + 11

$$\approx \boxed{l + \frac{n}{b}}$$

$$l + \frac{100-8}{\textcircled{8}} \times \underline{\underline{1 \text{ bc}}}$$

DRAM latency 100 ns

No. of FMA units = 2

4 FLOPs in 1 cycle

$4 \times 10^9 \text{ FLOPS}$ $\overset{10^9}{\sim} 10^9$?

FMA operation completion = 1 cycle

$\gamma_{Block} = 1 \text{ word} = 4 \text{ Bytes}$

CPU $\xrightarrow{\text{REQ 1}}$ RAM $\xrightarrow{\text{BLK 1}}$ CPU $\xrightarrow{\text{REQ 2}}$ RAM $\xrightarrow{\text{BLK 2}}$...
 \Rightarrow STALLS \Rightarrow STALLS \Rightarrow

```
float a[100], b[100], ans = 0.0f;
```

```
int i = 0;
```

// a & b are set here

```
for ( i = 0 ; i < 100 ; i++ )
```

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

After 200 ns, 2 FLOPs can be done ^{≈ 200 ns}

$$\frac{10^7}{4 \times 10^9} \times 100 = 0.25\%$$

1 SEC,

$?? \Rightarrow 2 \times 5 \times 10^6$
 $= 10 \text{ MFLOPS} = 10^7$

$$\frac{10^9}{200} = \frac{1}{2} \times 10^7$$

$$200 = 5 \times 10^6$$

1 Block = 4 words (each word = 4 bytes)

Does the peak speed change? 1 RAM block = 16 Bytes
↓ Buswidth = 16 Bytes

Fetch $a[0]$: 100 ns \Rightarrow fetches $a[1], a[2] \& a[3]$

Fetch $b[0]$: 100 ns \Rightarrow " $b[1], b[2] \& b[3]$

In 200 ns: compute using $a[0] \& b[0]$: 2 FLOPs

$a[1] \& b[1]$: 2

$a[2] \& b[2]$: 2

200 ns: 8 FLOPs $a[3] \& b[3]$: 2

1 s: ?

8 FLOPs

$8 \times 5 \times 10^6 = \underline{\underline{40 \text{ MFLOPs}}}$