

CPU central processing unit of a computer. It is a synchronous unit (which is synchronized by an internal clock) It contains the ALU

Clock Rate: CR is the frequency that the internal clock runs at. It is measured in pulses or cycles per second.

CR is measures in Hz
eg. 60Hz means 60 cycles per second

Clock Cycle (CC) is the time length of one cycle, measured in seconds. Therefore

CR = 1/CC

CC can be measured in milliseconds (ms) = $1 * 10^{-3}$ seconds

microseconds (μ) = $1 * 10^{-6}$ seconds

nanoseconds (ns) = $1 * 10^{-9}$ seconds

For each of these CC measurements there is a corresponding Hz Measurement

1KHz = 103Hz

1MHz = 106Hz

1GHz = 109Hz

Instruction Count (IC) Number of assembly commands that were performed by the program

Cycles Per Instruction (CPI) number of clock cycles required to perform the entire program divided by the instruction count, i.e. average cycles per instruction.

CPI Number of cycles required to perform one instruction of type i

CPI Time (Execution Time) Time it takes for the processor to execute a certain program, measured in seconds.

CPU-Time = Clock Cycle length in seconds * Number of Cycles in the program or
CPU-Time = $CC * CPI * IC = 1/CR * CPI * IC$

$$MIPS = \frac{IC}{CPU\ Time_{(sec)} \times 10^6} = \frac{IC}{CC \times CPI \times IC \times 10^6} = \frac{CR}{CPI \times 10^6}$$

Amdahl's Law (Weinman's explanation):

$$ExTime_{new} = ExTime_{old} \times \left[(1 - \text{Fraction}) + \frac{\text{Fraction}}{\text{Speedup}} \right]$$

Amdahl's Law

$S_{\text{theory}} = 1/(1 - p/s)$
 S_{theory} is the theoretical speedup of the execution of the whole task
p (also called Fraction) is the proportion of execution time that the part benefiting from improved resources originally occupied.

s is the speedup of the part of the task that benefits from improved system resources

$$\text{Speedup}_{\text{overall}} = \frac{ExTime_{old}}{ExTime_{new}} = \frac{1}{(1 - \text{Fraction}) + \frac{\text{Fraction}}{\text{Speedup}}}$$

SAMPLE:

15% -2.5 faster 20%-3 faster 40%-1.5 slower

$$ExTime_{new} = ExTime_{old} \times \left[\left(1 - (0.15 + 0.20 + 0.40) \right) + \left(0.15 * \frac{1}{2.5} \right) + \left(0.20 * \frac{1}{3} \right) + (0.40 * 1.5) \right] =$$

$$= ExTime_{old} \times \left[\left(1 - (0.15 + 0.20 + 0.40) \right) + \frac{0.15}{2.5} + \frac{0.20}{3} + 0.40 * 1.5 \right] =$$

$$= ExTime_{old} \times 0.98$$

$$\text{Speedup}_{\text{overall}} = \frac{ExTime_{old}}{ExTime_{new}} = \frac{1}{0.98} = 1.02$$

sign(1) | Exp (8) | Mantissa (23)
add 127 to exponent **IEEE 754 format**

Floating Point Arithmetic: raise the smaller power to match the bigger power
The algorithm for addition

Note: This will work for subtraction if we note that
A-B = A+B' + 1
10
-01 => +10
+01

Given 2 floating point numbers:
=> 01

1. Equate the powers by increasing the power of the smaller number and moving the decimal point accordingly
2. Add the digits ignoring the power
3. Normalize the result into the form 1.ffff...₂ biased exponent
4. Check for overflow or underflow, if yes, create exception
5. Truncate (or round) the mantissa to 23 bits
6. If we get a negative result, we must make it positive and set the sign bit to 1
7. If there is a mathematical overflow, we place the overflow bit in the sign bit location

Multiplication

1. Write the numbers in the float format
2. add the exponents
3. multiply
4. normalize the result
 - a. if the normalization created a new exponent, add it to the exponent from step 2
5. round the result if needed
6. put the point in the right place

Multiplication example: 5 * -0.4375

1. $5 = 1.010 \cdot 2^3 - 0.4375 = 1.110 \cdot 2^2$
2. Sum of powers is $2-2=0$
3. Multiply: $1.010 \cdot 1.110 = 1.000110 \cdot 1000$
4. Normalize: $10.00110 = 1.000110 \cdot 2^1$
5. Round up: 1.001^2
6. put the sign on the number

Explanation of rounding: look at the digit after the desired precision. Add that digit to the last digit of precision.

Division

Same as multiplication except, in step 2 subtract the exponents instead of adding them

and in step 3 perform division

BOOTH:		action to be performed on the multiplicand
Right most two bits of multiplier		action to be performed on the multiplicand
00		left shift 1
11		left shift 1
01		add to product and left shift 1
10		subtract from product and left shift 1

Note: For the first iteration, we look at only the single rightmost bit of the multiplier and assume a zero to its right. Thus, we either have 00 or 10.

Example with the booth algorithm for multiplication

We call X the multiplier ($x_{23..0}$) and Y the multiplicand ($y_{15..0}$)

Let's say we want to multiply $6 * 3$
we have 0110 * 0011 (note that we add leading zeros so both numbers are the same length)

First let's make a product of twice as many bits 0000 0000

Now take the multiplier (the first number) and extend it to the size of the product by adding zeros on the left side.

This gives us 0000 0110

Now we are ready to start the loop

1. We read the rightmost bit $a_{23..0}$ of the multiplier, as we need to shift left one bit (using sl). This gives us 0000 0110. We only do this on first iteration.

2. Since the right two bits of the multiplier are 00 we are going to shift left the multiplicand (getting 000000110) and neither subtract nor add to the product.

3. Now shift right the multiplier and examine the right two bits. They are 10, so we subtract the multiplicand from the product:

0000 0000 - 0000 0110

Following rules of two's complement, we note that we can add the negative of the second number, giving:

+0000 0000
+ 1111 1011 (remember: to make a negative value, flip all bits and add 1)
= 1111 1010

Or in MIPS just use sub:

Now, as in step 2. we left shift the multiplicand, obtaining:

0000 1100

Next we right shift the multiplier and examine the rightmost bits. These bits are 11 so we do nothing except left shift the multiplicand, obtaining: 0001 1000

Again we right shift the multiplier and examine the rightmost bits. They are 01 so we add the multiplicand to the product, as on right.

0001 1000 + 0001 1010 = 10001 0010 = 18

We could continue shifting the multiplicand left and the multiplier right but the multiplier has become all 0s so we will never perform any additions or subtractions. We are done. Also, $i = 1$

RAM:

Sequential Access memory: we can only access one memory cell after the other. E.g. a tape drive. This means that accessing distant data will be slower.

Associative Access: This is a type of Random Access. We access by value instead of by address. We perform a parallel search.

Capacity: Maximum amount of data that can be stored in the system. E.g. 1TB 2.

Access Time: This parameter is defined differently for different memory systems.

Random Access: access time here as the time from when the address is supplied to the MAR until we can access it in the MIR.

Non Random Access: the time between the processor requesting the data from the controller and the controller supplying the data to the processor, not including the time it takes the controller to search the disk, just including the time to actually read the data.

Data Rate: Measured in BPS (Bits Per Second).

Cycle Time: Minimal time the processor must wait between the starts of two consecutive memory accesses.

Cycle Time = Access Time + Dead Time

Dead Time = Recovery Time + Transient Time

Transient Time: Time needed for the data to be electrically stable

Recovery Time: Time need to refresh the data, as necessary in DRAM

General Description of SRAM



Memory size is $2^n \times n$ where n is the number of addresses and n is the size of the data
RD (or RD') is the read control line
WR (or WE) is the write control line
When CS is 1 the data lines are TS
When CS is 0 the data lines are active
when CS=RD=0 and WR=1 then the value at address are read to the data lines
when CS=WR=0 and RD=1 then the data is written to the address in memory.

1D Organization of SRAM

is the memory space in bits
Using a decoder ($m_{n-1..0}$) we select desired word (n bits).
The line output of the decoder is connected to an AND gate for that memory location.
This requires 2^n AND gates.

Disadvantage: Needs a lot of logic gates, for a tiny $m=4$ you would need 2^4 AND gates.

Advantage: Simple to implement.

2D Organization of SRAM for Reading



2D SRAM explanation

The Address is m bits divided into two sections j and k such that $m=j+k$.
The Physical Memory is divided into n segments such that $s=n^k$.

Hence, the size of k determines the number of segments.
The width of each segment is the addressable word length, n bits ($n, m=k$).

The number of words (addresses) in each segment is 2. We note the number of segments, s .

We want a total of 2^n addresses in memory. $s=2^k \cdot 2^l$ and therefore $s=2^m$.

SRAM Memory Design
Goal: Enlarging the memory by using multiple small memories.

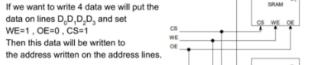
We will see methods for two types of memory expansion:

- 1) Increase the size of the data while keeping the address space constant.
- 2) Increase the number of addresses while keeping the data size constant.



Increasing the Size of the Data (word length)

We will create a data with size 4 by combining 2 of the previous memory units



If we want to write 4 data we will put the data on lines D0, D1, D2, and D3. Then this data will be written to the address written on the address lines.

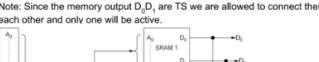
Read will be WE=0, OE=1, CS=1

Expand memory by Expanding the Address Space
The data will remain 2 bits. We want to expand the address space by 4 and still have a 2 bit data.

The new size will be $(2^n \times 4) \times 2$ or $2^{n+2} \times 2$. This means we need to add two bits to the address space.

Note: In the following diagram for the sake of simplicity, we omit the WE and OE. But in reality they will be there.

Note: Since the memory output D0, D1, D2, D3 are TS we are allowed to connect them to each other and only one will be active.



CACHE:

Hit: the datum is in the cache

Miss: the datum is not in the cache and must be brought from main memory

Hit-Rate: Number of access where the datum is in the cache/total accesses

miss-rate = 1 - Hit-rate

Block: A unit of bytes by which we divide the Cache and Main Memory. Identical for both. If we need to bring a single byte into the cache we will need to copy the entire block in which it resides.

Hit-Time: The time it takes to either read or write a datum to or from the Cache.

Miss-Penalty: The time it takes to bring a block from Main Memory to the Cache, or to write it to Main Memory from the Cache. As we have said, when the processor reads from cache it reads one datum, but when we copy data into the cache from main memory, we copy an entire block.

Average Memory Access Time (for one datum) = Hit-Time + Miss-Rate x Miss-Penalty

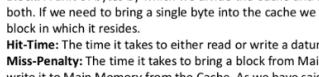
3 Methods for Cache Mapping:

1. Direct Mapped: Simple to implement but not flexible
2. Fully Associative: Any block can be mapped to any location in cache. Flexible, high cache utilization, Complex to implement.
3. n-way Set Associative: Compromise between other two methods

Direct Mapping: Block Number Mod Number of blocks in Cache \Rightarrow Block Number in cache.

example address:

Graphical representation of finding data



There are 4 options for this step based on the type of instruction:

- A. R type instructions: The ALU performs an action on operands A and B, formally: RESULT-ALU = A op B
- B. Accessing Memory (SW/LW): ALU calculates the effective memory address by adding the OFFSET to RS
- C. Branch Commands: (Reminder: In step 2 the jump address was stored in the TARGET register) The ALU determines if A and B are equal by subtracting A - B. If the subtraction result is equal to 0 then flag Zero is set to 1 and the address in TARGET will be transferred to PC (i.e. a jump will be performed). If the result is different from 0, the PC register will remain with the value PC + 4.
- D. Jump Commands: PC is 4 bits (PC+4) + 26 bits INDEX <> 2

3. Execute (EX)

We now know what command we have. We perform the required actions on the A and B values that were calculated in step 2.

There are 4 options for this step based on the type of instruction:

- A. R type instructions: The ALU performs an action on operands A and B, formally: RESULT-ALU = A op B
- B. Accessing Memory (SW/LW): ALU calculates the effective memory address by adding the OFFSET to RS
- C. Branch Commands: (Reminder: In step 2 the jump address was stored in the TARGET register) The ALU determines if A and B are equal by subtracting A - B. If the subtraction result is equal to 0 then flag Zero is set to 1 and the address in TARGET will be transferred to PC (i.e. a jump will be performed). If the result is different from 0, the PC register will remain with the value PC + 4.
- D. Jump Commands: PC is 4 bits (PC+4) + 26 bits INDEX <> 2

4. (Only for Load, Store or Rtype)

R-type: The result computed in step 3 which sits in the ALUOut register is written by the RF to register RD: Reg [IR (15-11)] = ALU Result

SW: The address computed in stage 3 gets sent to Write Address lines of the memory.

Contents of RF from the RF to Write Data of the memory. Control send the write command to memory: MEMORY = [ALU-Result] = B

LW: The address computed in stage 3 is sent to Read Address of memory unit, controller sends a read command to memory. The data read is put on the Memory Data lines of the memory Unit: Memory-Data = MEMORY[ALU-Result]

5. Write Back (Only for Load)

The memory data retrieved in stage 4 is written to register RT

Reg [IR (20-16)] = Memory Data

