Von - Neumann Architecture (1945): Princeton Architecture

Fetch - Decode - Execute Cycle:

a) Program counter points to the present instruction to be fetched.
b) Bits in the register "control" the subsequent actions.
c) Fetch the next instruction & continue.

Stored program concept.

a) Same physical memory to save instructions and data
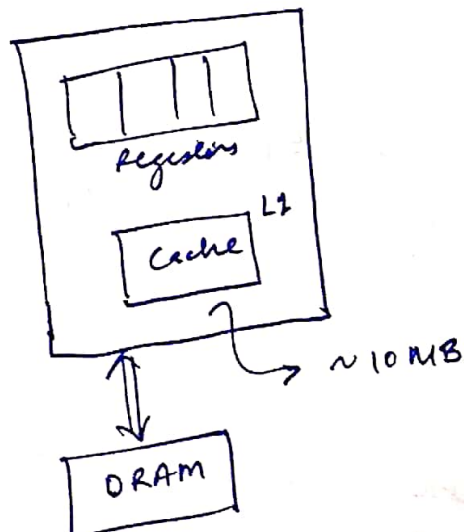b) Instruction fetch & data transfer cannot be done simultaneously need 2 clock cycles.

# NOTE: Hard disk is treated as an I/O device & NOT memory. as there are computers without hard disk.

A basic block of code has usually no jump or branch instructions averages size = 6-7 instructions in other words every 6th or 7th instruction is a branch or jump.

Von - Neumann Bottleneck:

Time spent in memory access limits performance.
(as same memory is used for program & data)

CPU



note an 18Mb cache
is not preferred over

2, 10 cache and 8 cache
as the more the cache
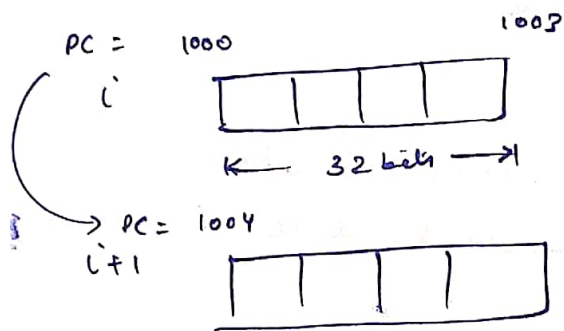size, the more the
time to search

~ 10 MB

To avoid the bottleneck, later architectures restrict operands to registers

MIPS processor has 32-bit registers.

* MIPS does not have cache memory

MIPS every instruction is 4 bytes.
∴ program counter increases by 4.

PC = 1000                              1003

i



K— 32 bits —→|

→ PC = 1004
i+1

by default pc (program counter) advance by steps of 4.
but if it encounters a jump, the control unit
takes it to 5000.

$$1000 \xrightarrow{+4} 1004 \xrightarrow{\text{Control Unit}} 5000$$

Branchy does not take an extra clock cycle

instructions on signed integers DIFF from unsigned integer

2 sources       one destination

add     a, b, c     # a gets b + c
All arithmetic operations have this form

compiled MIPS code :

add t0 , g , h          # temp t0 = g + h
add t1, i, j            # temp t1 = i + j
sub f , t0, t1          # temp f = t0 - t1

c code :

$$f = (g+h) - (i+j)$$

operations occur on data already present in
registers. ∴ we need to do register binding
MIPS has 32 × 32 - bit register file
Numbered 0 to 31
32 - bit data called a "word"

Assembler names:

1) $t0, $t1 ..., $t9 for temporary values
2) $s0, $s1 ... $s7 for saved variables

Processor → Memory   LOAD
memory   → Processor  STORE

Direct transfer of data from memory to memory (NOT allowed)
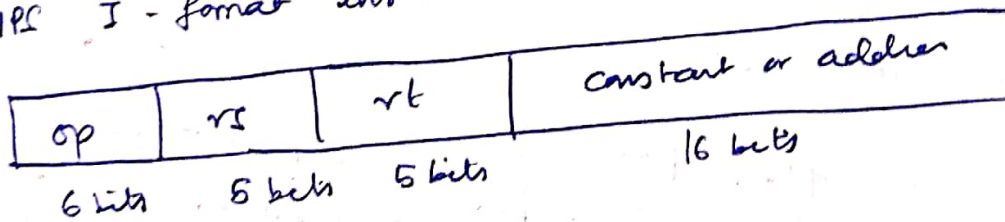operation on data in memory (NOT allowed)

Register numbers:

$t0 — $t7 are reg's 8-15
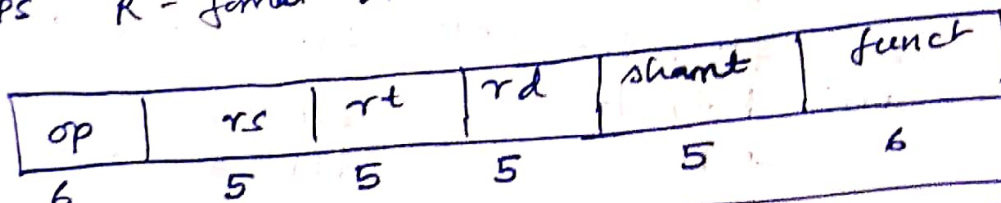$t8, $t9 are reg's 24-25
$s0 - $s7 are reg's 16-23

Immediate Operands:

Constant data specified in an instruction
addi  $s3, $s3, 4   } entire instruction must
                      AMIN be encoded in
   5 bit    5 bit           32 bits.
        32 bit

MIPS  I - format instruction

| op | rs | rt | Constant or address |
|----|----|----|---------------------|
| 6 bits | 5 bits | 5 bits | 16 bits |

MIPS  R - format instruction

| op | rs | rt | rd | shamt | funct |
|----|----|----|----|-------|-------|
| 6 | 5 | 5 | 5 | 5 | 6 |

op   : operation code
rs   : first source register number
rt   : second source register number
rd   : destination register number
shamt : shift amount (00000 for now)
funct : function code (extends opcode)

Barrel Shifter

special shift register!!

Loop — A1 (address)
$\vdots$

$A_2$ — bne  $t_1$, $zero, $\underline{\underline{Loop}}$
J

this contains $A_2 + 4 - A_1$

and **not** the entire address of Loop

b/c Loop address = 32 bits
space in beg for Loop ~ 16 bits

NOTE:  shamt (shift amt.) has 5 bits
as we can shift by maximum 0-31 positions

shift right might cause problem in sign!

j LL
unconditional jump to instruction labelled LL

C code:
while ( save[i] :: k)  i += 1;
( checks whether all elements of array "save" are
equal to k )

compiled MIPS Code :

i = 0          i in $S8
$S3 = 0        k in $S5

| Loop: | sll  | $t1, $S3, 2 | $t1 = 0×4=0  address of |
|       | add  | $t1, $t1, $S6 | $t1 = $t1 + $S6  save in |
|       | lw   | $to, 0($t1) | = & save[0]  $S6 |

bne  $to, $S5, Exit        & to = Mem[$t1 + 0]
addi $S3, $S3, 1           $t1 = $S3×4 = 4
j    Loop                  $t1 = 4 save[i]

Exit.                      $\vdots$
and so on

To check whether 2 32-bit no's are equal,

bitwise XOR the 2 numbers:

Take NOR of all bits in the 32 bit XOR value,

if NOR = 1 then XOR value = 0
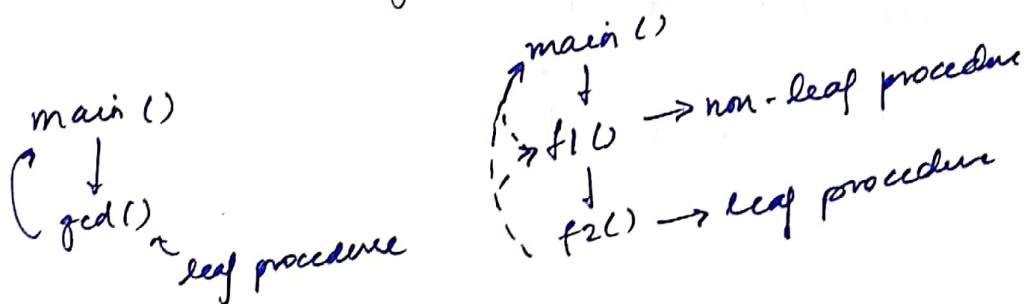
and both no's are equal.

Note: hardware for $<$, $\geq$ is slower than $=$, $\neq$

that is why no blt, bgt etc.

→ Small negative numbers in 2's complement look like large numbers in unsigned.

let array bound = 15

index = 16

bound − index = −1

$\underbrace{\qquad}_{\text{signed subtraction}}$

main ()



main ()

└→ gcd()

↳ leaf procedure

main ()

↓

→ f1 () → non-leaf procedure

↓

f2() → leaf procedure

jal f2

↑

jump and link

jr $ra      stores address of instruction

next to the one it has to jump to

in $ra (return address)

instruction next to jal f2 )

```
fact:
        addi $sp, $sp, -8
        sw   $ra, 4($sp)
        sw   $a0, 0($sp)
        slti $t0, $a0, 1
        beq  $t0, $zero, L1
        addi $v0, $zero, 1
        addi $sp, $sp, 8
        jr   $ra


L1:     addi $a0, $a0, -1
        jal  fact
        lw   $a0, 0($sp)
        lw   $ra, 4($sp)
        addi $sp, $sp, 8
        mul  $v0, $a0, $v0
        jr   $ra
```
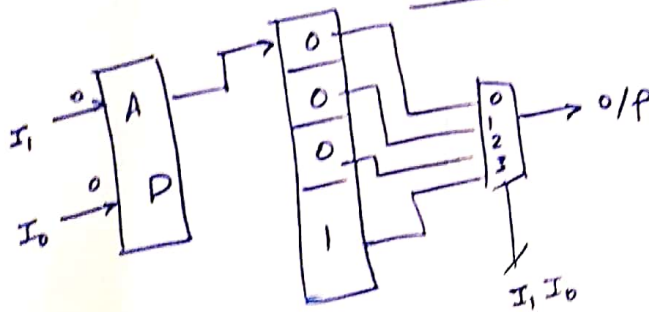
- void strcpy (char *dest, char *src)

```
{   while (* dest ++ = *src ++);

}
```

to use a register, first save its previous value somewhere and then after all computations are done, restore its value.
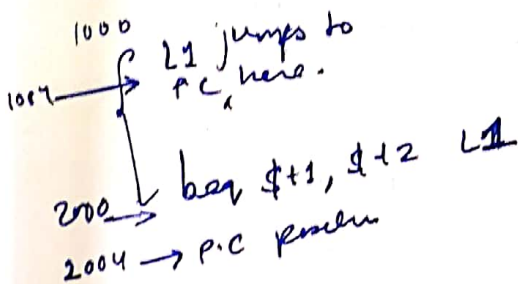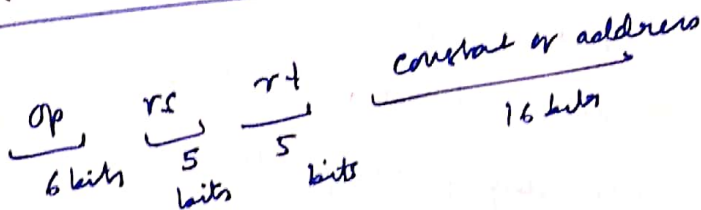
After end of a function, when stack pointer returns to original position, the memory stored in the registers is **NOT** destroyed.

Lookup table (LUT)



$I, I_0$

Electrically, all registers are exactly the same, except $ zero

ori $t1, $t2, 0x2121

## Branch Addressing :

| op | rs | rt | constant or address |
|---|---|---|---|
| 6 bits | 5 bits | 5 bits | 16 bits |

NOTE: mips label addresses all end with 00.

∴ all address are divisible by 4

1000
1004 — { ℓ1 jumps to PC here.

2000 → beq $t1, $t2  L1
2004 → P.C goes here

offset = $\frac{1000 - 2004}{1000 - 2004}$

$= \frac{-996}{4} = 224$

## Jump addressing : ??

```
Loop:  sll  $t1, $s3, 2          80000
       add  $t1, $t1, $s6        80004
       lw   $t0, 0($t1)          80008
       bne  $t0, $s5, Exit       80012
       addi $s3, $s3, 1          80016
       j    Loop                 80020    | 2 | 2000 |
                                 80024
Exit:
```

Addressing mode

a) Immediate addressing
b) Register addressing
c) lw , sw, lh , sh  ; Base addressing
d) PC relative addressing   branch conditionals
e). Pseudodirect addressing ; jump (unconditional)
                                    or jal
26 bit → left shifted,
centrally
                ↓
        concatenate with 4 bits of P.C.


MIPS ∈ RISC : Reduced Instruction set Computers.
x86  ∈  CISC : Complex Instruction set computers.

Major difference :  MIPS has lesser addressing modes
Interestingly  CISC came before RISC

        (but they were bulky & slow)

• pipelining / superscalar , multicore

NOTE: MIPS has only global registers

∴ for recursion, you need to save on a stack.

NOTE:

arguments: $a0 - . $a3
if more arguments, use an array to
store & retrieve from
memory.
DO NOT use $t
(it will work but
is not recommended)

Q. Implement the following C code in MIPS:

```
int fib (int n) {

    if (n <= 1)
        return n;

    else
        return fib(n-1) + fib(n-2);

}
```

$a1 = array
$a2 = argument
$a3 = 1

→ fib:
```
    addi $sp, $sp, -12
    sw   $ra , 0($sp)
    sw   $a2 , 4($sp)
    bgt  $a2, $a3, recurse &
    move $v0 , $a2
    jr   $ra
```

recurse:
```
    lw   $t0, 4($sp)
    addi $t0 , $t0 , -1
    move $a2, $t0
    jal fib
    sw   $v0, 8($sp)
    lw   $t0, 4($sp)
    add  $t0, $t0, -2
    mov  $a2, $t0
    jal fib
```

```
        lw   $t1 , 8($sp)
        add  $vo, $vo, $t1
        lw   $ra,  0($sp)
        addi $sp, $sp, 12
        jr   $ra
```

(2) :

```
   func:
        addi $to , $zero, 1
        addi $vo, $zero, 1
```
                                        → set less than or equal to
```
   Loop:
        sle  $t1, $to, $a0
        beq  $t1 , $zero, Exit
        mul  $vo , $vo, $to
        addi $to, $to, 1
        j  Loop
```

```
      Exit:
        jr  $ra
```

```
int func (int a0)
{ int to = 1
  int vo = 1

    for ( to = 1 ; to ≤ a0 , ++to )

        {   vo = vo * to ;

        }
    return vo;
}
```

# Instructions per Clock Cycle (IPC)

IPC average number of instructions executed per clock cycle.

For processors with a single execution unit, $IPC|_{max} = 1$

In practice, IPC is always $< 1$

$$cPI = \frac{1}{IPC} \quad ; \quad CPI|_{ideal} = 1 \ , \ CPI|_{actual} > 1$$

↓
clock cycles per instruction

$IPC = 0.8$
∵ $CPI = 1.25$

goal of CoA:
increase IPC (and thus decrease CPI)

If we have super-scalar processors with $IPC|_{max} > 1$
we will get better processors!!

## Execution time of a program (P):

a) P has machine level instructions (IC — instruction count)
b) CPI: average clock cycles per instruction
c) CCT: clock cycle time.

∴ CPU-time = $\boxed{IC \times CPI \times CCT}$ (units (equal to CCT is unit))

humans can distinguish b/w only events which are atleast
10 milliseconds apart.

$En|_{CPU} : IC \times CPI \times CCT$

CCT → completely dependent on Technology
IC → better Algorithm / compiler
CPI → reduce clock cycles

In CMOS IC Tech:
$Pnev = Capacitive\ load \times Voltage^2 \times Frequency$

MIPS (Millions of Instructions per second)

$$= \frac{\# \text{ of Instructions executed}}{\text{Ex}|_{cpu} * 10^6}$$

Rejected as $\text{Ex}|_{cpu} \rightarrow 5$ times

$\# \text{ of instruction} \rightarrow 5$ time

$\Big\}$ MIPS reap same

Each instruction of $i^{th}$ class : $CPI_i$ # of clock cycles to execute

Note: Just total no. of instructions does not help, as some instructions may take up more time than others.

Amdahl's law : ("Law of diminishing Returns") $\rightarrow$ (as speedup is not $\propto$ number of CPUs)

Time taken to execute a program using one CPU $= T_1$
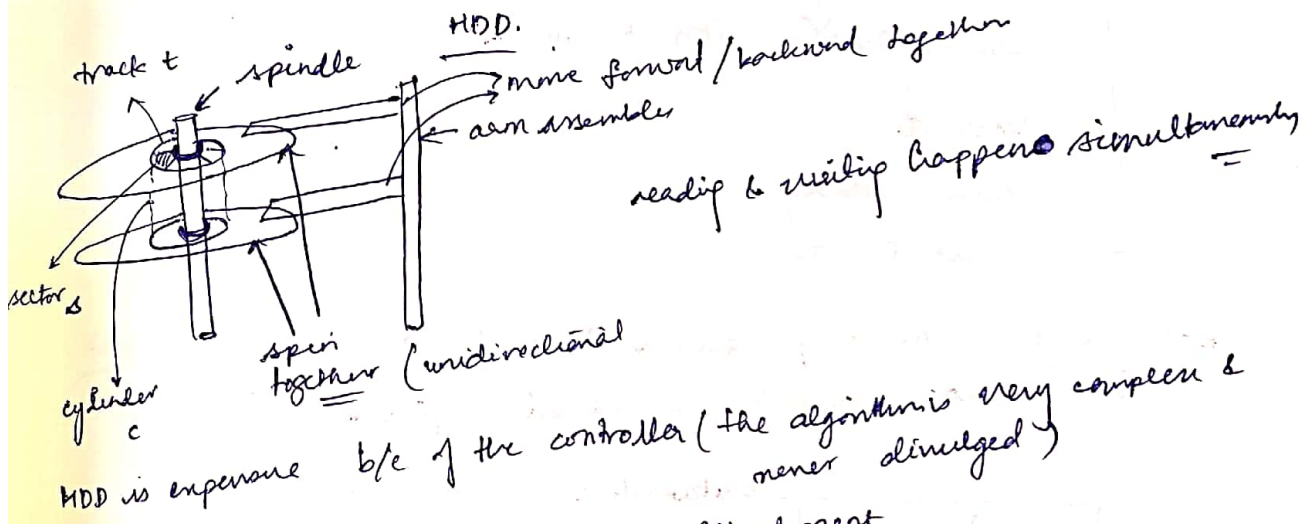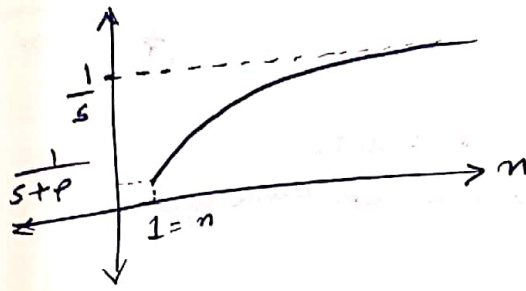Time taken to execute the program using 'n' CPUs $= T_n$

$\therefore$ Speedup $= \dfrac{T_1}{T_n}$ ;  Time to execute parallelizable part $= P$
Time to execute non parallelizable part $= S$

for one processor situation, let (normalized) $S + P = 1$

$$T_1 = P + S ; \qquad T_n = \frac{P}{n} + S$$

$\therefore$ speedup$_n$ : $\dfrac{P+S}{\dfrac{P}{n}+S} = \dfrac{1}{S+\dfrac{P}{n}}$  if $(s+P = 1$ (normalized)

$$\therefore \lim_{n \to \infty} \text{Speedup}_n = \frac{1}{S}$$

Graph: vertical axis marked $\frac{1}{s}$ (dashed line) and $\frac{1}{s+p}$ at origin region, horizontal axis $n$, with point $1=n$, curve rising.

---

**HDD.**

track t    spindle

→ move forward/backward together
← arm assembly

reading & writing happens simultaneously

sector s

cylinder
c

spin together (unidirectional)

HDD is expensive  b/c of the controller (the algorithm is very complex & never divulged)

SATA : (serial) Advanced Technology Attachment
SCSI : Small computer system Interface.
skuzy → I scazee (pronunciation !)

On an average; half the length of a track to be rotated to access data.
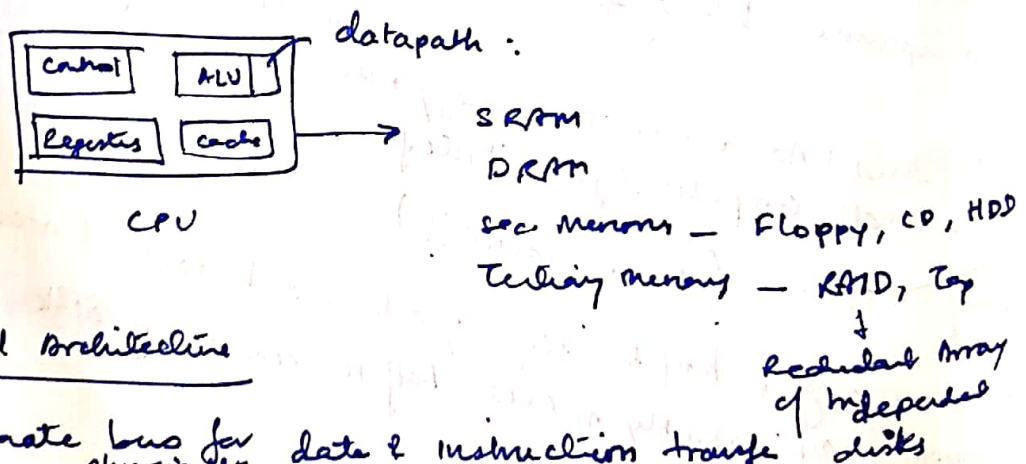On an average ; arm assembly moves half the radius of disk

Scanned by CamScanner

# Von Neumann

Physical Memory to save instructions & data
Instruction fetch & data cycle, 2 sep. Clock cycle

→ Instruct
fetch ↓

Ins. decode
operand fetch ↓    Von Neumann

Execute ↓

Result store ↓

Next Instruction

Bottleneck, memory access takes times
∴ nano, operands stored in registers

datapath :

| Control | ALU |
|---------|-----|
| Registers | Cache |

CPU

SRAM
DRAM
Sec Memory — Floppy, CD, HDD
Tertiary memory — RAID, Tap
                                    ↓
                              Redundant Array
                              of Independent
## Harvard Architecture                        disks

separate bus for data & instruction transfer
if not physically separate, then also

## Performance affected by :

Algo:— no. of operation ( IC & possibly CPI )
no. of Language / Compiler / Architecture:
   No. of instructions executed per speler
Processor & Memory system
   How fast instruction executed
I/o system (including OS)
   How fast I/o opns executed

Response time , time for one operate incl· I/O, O·S, CPU time

Throughput : work done / unit time

$$performance \propto \frac{1}{CPU\ time}$$

M/c level program P.

$$Exec.\ time = IC \times CPI \times CCT$$
(= CPU time) | Instruction | |
                Count    | avg. Clock cycles | clock cycle 'tms
                          per instr

$$CPU\ time = CPU\ clock\ cycles \times CCT$$

Reduce :

Reduce no. of clock cycles
Increase clock rate
Hardware designer must often trade off clock rate against
cycle count

$$CPU\ time = \frac{CPU\ clock\ cycle}{clock\ Rate}$$ — How fast

Clock cr :

10x : $\frac{2}{?}$ $\frac{2\ Ghz}{}$

6x : $1.2 \times 2 \times C.T'$

$\frac{6}{2.4}$ $\frac{6}{2 \times 2}$ $5 : 2.5$

/ instruction count

$$No\ Clock\ cycles = IC \times Cycles\ per\ instruction$$

$$\therefore CPU\ time = \frac{IC \times Cycles\ per\ instruction}{clock\ rate}$$

Avg. cycles P.I.

$$clock\ cycle = \sum_{i=1}^{n} (CPI_i \times IC_i) \qquad i = i^{th}\ class\ instruction$$

$$CPI = \sum_{i=1}^{n} \left( CPI_i \times \frac{IC_i}{IC_{total}} \right)$$

CPU time: $IC \times CPI \times CCT$

Algo: @ IC & CPI

Programming language: IC & CPI

Compiler: IC & CPI

Instruction set architecture: IC, CPI, CCT.

CPI: memory hierarchy, pipeline

CCT: logic design, technology

MIPS: Millions of instructions per second

$$= \frac{IC \text{ of program } P}{\text{Exec. time of P in sec} \times 10^6}$$

Power = Capacitive Load $\times$ Voltage$^2$ $\times$ Freq

Exec. time: best performance measure

use parallelism to improve performance as power is limited

Power wall: min. voltage ⟶ limit
heat removal max limited

$$\text{Ex Time}_{new} = \text{Ex Time}_{old} \left[ (1 - \text{fraction enhanced}) + \frac{\text{fraction enhanced}}{\text{speedup enhanced}} \right]$$

$$\text{speedup (overall)} = \frac{\text{Ex time old}}{\text{Ex time new}} = \frac{\text{Perf. new}}{\text{Perf. old}}$$

MIPS : big endian
MSB at least address of word

R- format

| op | rs | rt | rd | shamt | func |
|---|---|---|---|---|---|
| 6 | 5 | 5 | 5 | 5 | 6 (extends opcode) |

I . format
$\longrightarrow -2^{15} / 2^{15} - 1$

| op | rs | rt | const /address |
|---|---|---|---|
| 6 | 5 | 5 | 16 |

$7 fffffff_{hex}$

| Stack |
| Dynamic data |
| Static data |
| Text |
| Reserved |

sll : shift left logic
srl : shift right logical

$\longrightarrow$ mask bits

eg.
and $t0, $t1, $t2

| t2 | 0 0 1 1 | 0 1 1 1 0 0 |
| t1 | 1 1  1 1 | 0 0 0 0 0 0 |
| t0 | 0 0 1 1 | |

or $\longrightarrow$ include bits

nor  $t0,  $t1,  $zero  = not

basic block :
no branch (except at end)
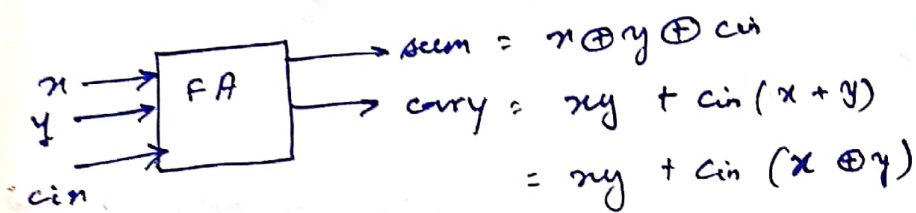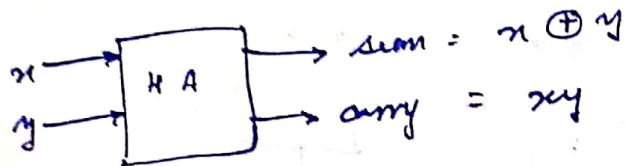no branch target (except at begin)

slt  rd, rs , rt
   rd = 1 if rs < rt
      else rd = 0

lui
ori

slti  rd , rs, constant

slt  $t0, $s1, $s2    if s1 < s2

bne  $t0, zero, L    branch to L

( better than
    slt )

# Binary (Integer) Addition :



Half Adder:

sum $= x \oplus y$

carry $= xy$



Full Adder:

sum $= x \oplus y \oplus c_{in}$

carry $= xy + c_{in}(x+y)$

$= xy + c_{in}(x \oplus y)$



$xy$

$x \oplus y$
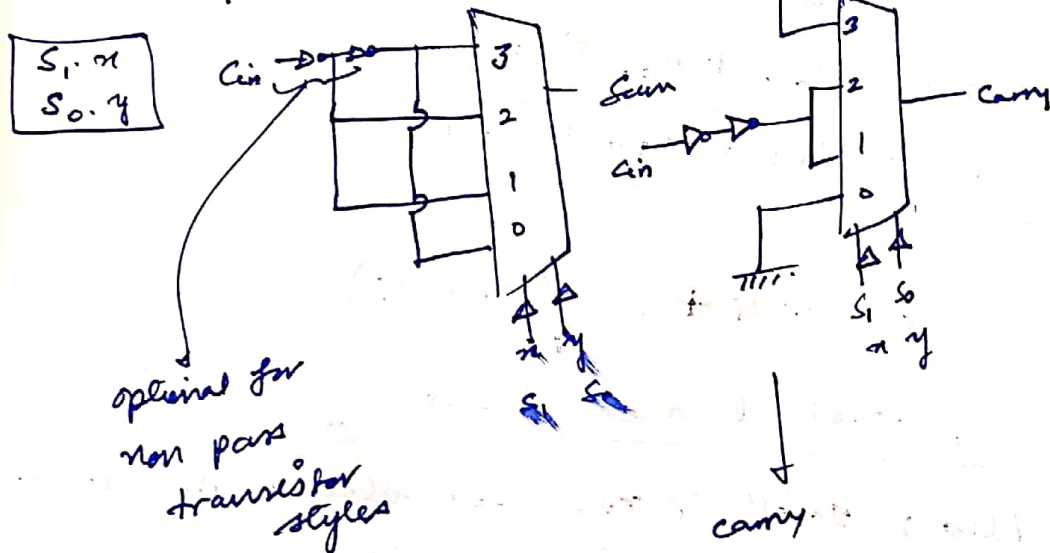
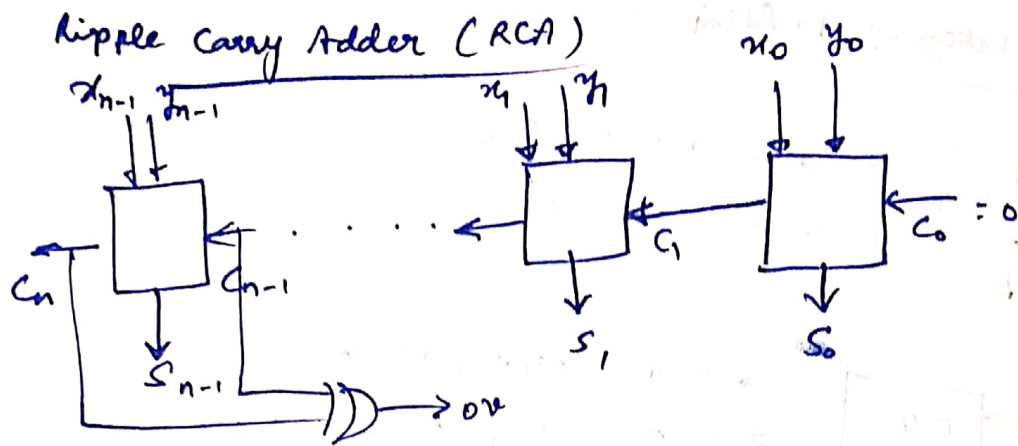$xy + c_{in}(x \oplus y)$

sum $= x \oplus y \oplus c_{in}$

## FA with MUXes :

O Muxes can be implemented very efficiently using
. pass transistor logic CMOS circuits.

Sum $= (xy) c_{in} + (x\bar{y})\bar{c_{in}} + (\bar{x}y)\bar{c_{in}} + (\bar{x}\,\bar{y})c_{in}$

$S_1 \quad S_0$

$S_1 \cdot x$
$S_0 \cdot y$

$V_{DD}$

$c_{in}$

Sum

$c_{in}$

Carry

optional for
non pass
transistor
styles

carry

carry

$= (xy)1 + (x\bar{y})c_{in} + (\bar{x}y)c_{in} + (\bar{x}\,\bar{y})0$

## Ripple Carry Adder (RCA)



## Overflow Bit

for 2's complement number addition, usually $C_{out}$ is discarded. However, $C_{out}$ is useful in detecting overflow

overflow : 2 no's of same sign are added but the result of opposite sign.

$$ov = x_{n-1} \, y_{n-1} \, \overline{S_{n-1}} + \overline{x_{n-1}} \, \overline{y_{n-1}} \, S_{n-1}$$

$\downarrow$
to detect overflow

Prove this

Text book definition $\boxed{ov = C_n \oplus C_{n-1}}$

sometime $C_{n-1}$ is not available outside the chip

then we can do the following :

$$S_{n-1} = x_{n-1} \oplus y_{n-1} \oplus C_{n-1}$$

$$\boxed{C_{n-1} = S_{n-1} \oplus x_{n-1} \oplus y_{n-1}}$$

$$C_i = x_i y_i + C_{i-1} (x_i + y_i)$$

## Carry Lookahead Adder (CLA)

Idea ; unroll $C_i$ recurrence relation to remove dependency on $C_1, C_2, \ldots, C_{i-1}$.
Ideally leads to $O(1)$ delay addition of two. $n$-bit numbers. But not feasible practically.

**Carry Generate** $\quad g_i \overset{def}{=} x_i \, y_i \qquad g_i = 1 \Rightarrow$ i-th stage generated a carry of its own

**Carry propagate** $\quad P_i = x_i \oplus y_i$

$C_{i+1} = g_i + C_i P_i$ ← always

$P_i = 1 \Rightarrow$ ith stage will propagate the i/p carry $C_i$ to $C_{i+1}$

Absorb. $a_i = \overline{x_i} \, \overline{y_i}$

On an average carry propagation occurs in length 2, doesn't occur when $x_i, y_i = 0, 0$
$\qquad \qquad \qquad \qquad \qquad x_i, y_i = 1, 1$

<u>carry select Adders / Conditional sum Adders.</u>

Motivation $O(\lg(n))$ delay for $n$-bit addition

Brent Kung (early 1980's) ⎤
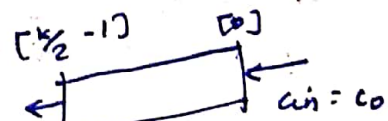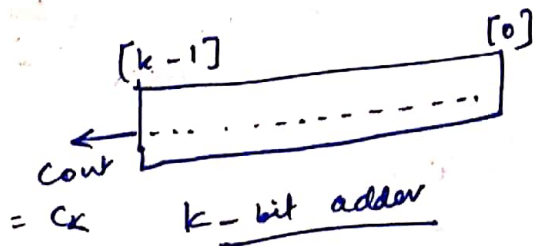Kogge Stone (1970s) ⎦ very old

$S = a \oplus b \oplus C_{in}$

$C_{out} = \underline{ab} + C_{in}(a \oplus b)$
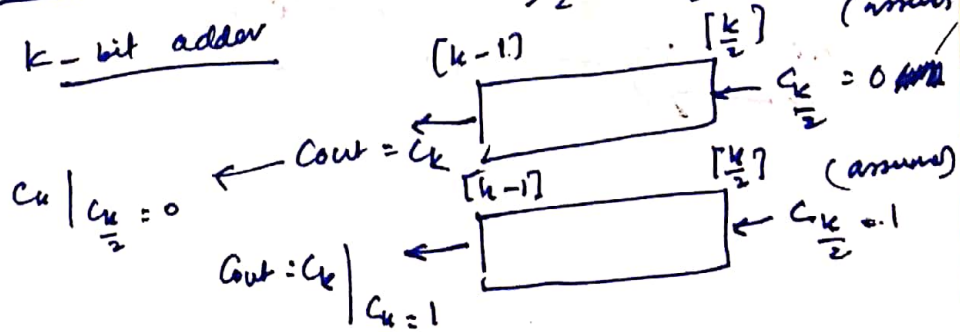$\qquad \qquad$ re-express in terms of $a \oplus b$ to enable logic sharing with $s$.

$ab = (a'b' + ab) \cdot a$
$\quad = (\overline{a \oplus b}) \cdot a$

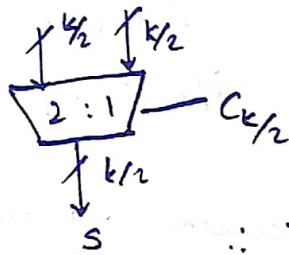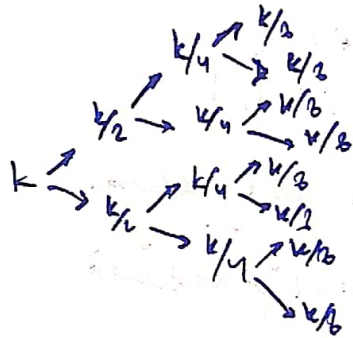$C_{OUT} = (\overline{a \oplus b}) \, a + C_{in}(a \oplus b)$

generate $C_k$ for both $C_{\frac{k}{2}} = 0$ & $C_{\frac{k}{2}} = 1$

finally select the correct o/p using a $2:1$ mux

and $C_{\frac{k}{2}}$ as select line.



$\therefore T(k) = T(k/2) + c$



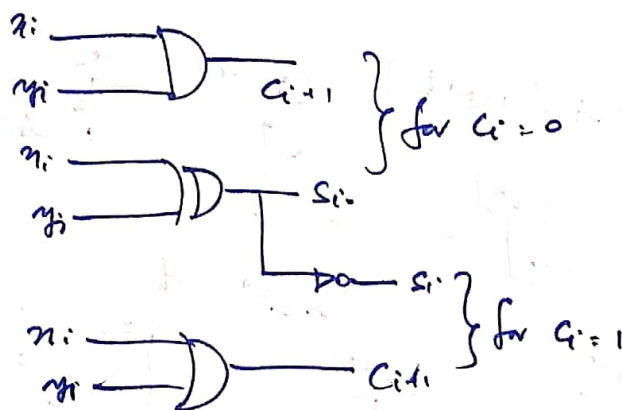till $n$ leaf nodes, we'll find specialised full adders

have $S_i = z_i \oplus y_i \oplus C_i$

$\Rightarrow S_i|_{C_i=0} \xrightarrow{?} z_i \oplus y_i$

$\Rightarrow S_i \oplus|_{C_i=1} : \overline{z_i \oplus y_i}$

$C_{i+1} = z_i y_i$ if $C_i = 0$

$= x_i + y_i$ if $C_i = 1$



**Hardware cost**

$O(k \lg k) \cdot (\text{for mux})$
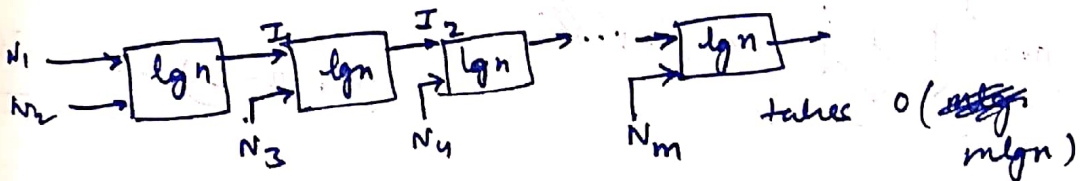
it dominates the total time taken

exercise

perform 2 level unrolling with
4 $\frac{k}{4}$-bit adders

---

NOTE: all this is done as ripple carry takes a lot of time

---

## carry save adders:

Goal: Assuming I have a fast carry-lookahead adder
with logarithmic delay which I call "logarithmic CLA"
$$(LCA) \leftarrow$$
(gives $O(\lg n)$ delay
for $n$-bit
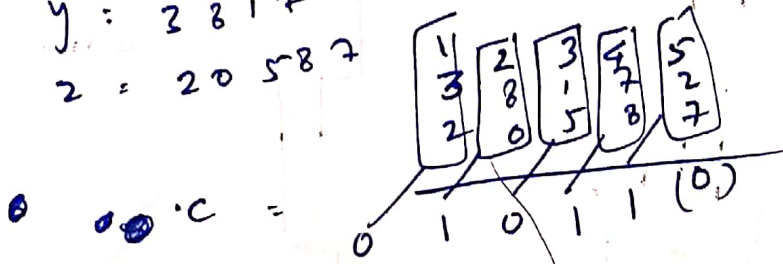addition)

I will add $m$ numbers each of $n$-bits in time $\underline{O(\lg m \lg n)}$



$N_1 \rightarrow \boxed{\lg n} \xrightarrow{I_1} \boxed{\lg n} \xrightarrow{I_2} \boxed{\lg n} \rightarrow \cdots \rightarrow \boxed{\lg n} \rightarrow$
$N_2 \rightarrow$
$\quad\quad N_3 \quad\quad N_4 \quad\quad N_m$  takes $O(m \lg n)$

$$S = (x + y + z) \% 10$$

$$C_{i+1} = (x_i + y_i + z_i) / 10$$
$\quad\quad\quad\quad\quad\quad\quad\uparrow$ integer division

e.g   $x = 1\ 2\ 3\ 4\ 5$
$\quad\quad y = 3\ 8\ 1\ 7\ 2$
$\quad\quad z = 2\ 0\ 5\ 8\ 7$

$\quad\quad\quad\quad\quad\quad \boxed{\frac{4}{3}}\ \boxed{\frac{2}{8}}\ \boxed{\frac{3}{1}}\ \boxed{\frac{4}{7}}\ \boxed{\frac{5}{2}}$
$\quad\quad\quad\quad\quad\quad\quad 2\quad 0\quad 5\quad 8\quad 7$

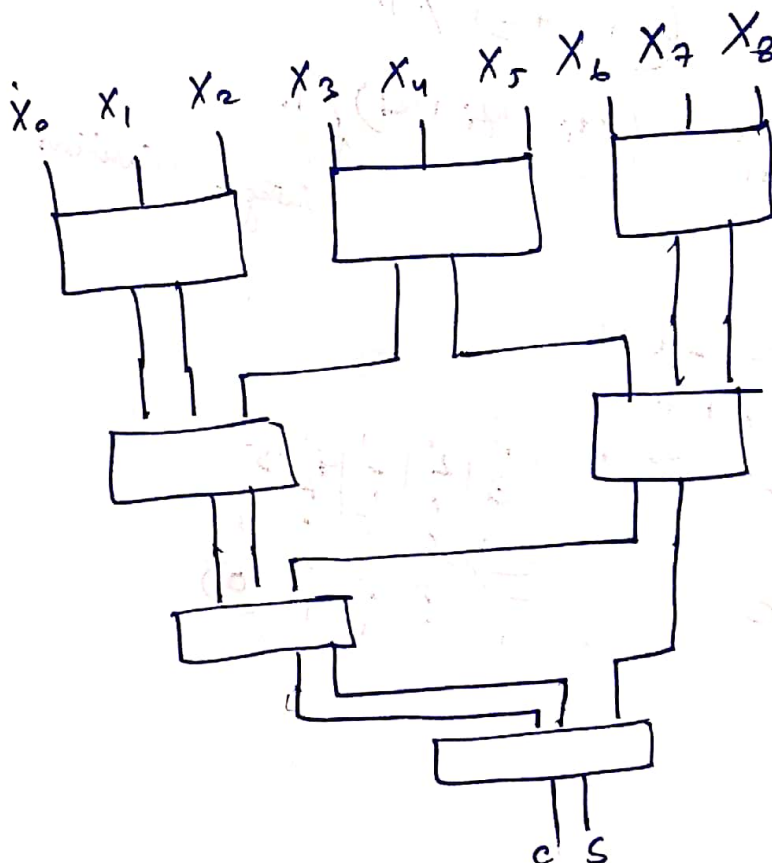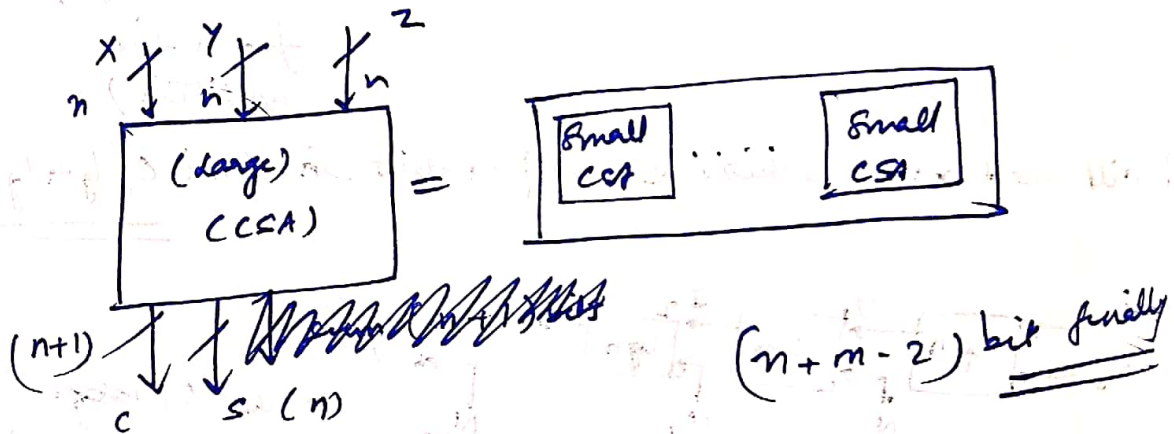$\quad\quad\quad\quad \cdot C = \quad 0\quad 1\quad 0\quad 1\quad 1\ (0)$

## 3-to-2 encoding:

note:

i) $\{c, s\}$ generation for 3 binary numbers with any no. of bits can be done in $O(1)$ time.

ii) $c + s = $ sum can happen. in $O(\lg n)$ time, where $n = $ no. of bits in each number.

iii) 3 numbers each $n$-bit long can be added in $O(\lg n)$ time

Now to add $m$ binary numbers each "$n$" bit long.



$(n+1)$ ... $c$ ... $s$ ($n$)

$(n + m - 2)$ bit finally

Quiz $3^{rd}$ September



$c$ $s$

# Booth's multiplication :

↑ see theory later

example :

4-bit representation  $X = -5 = 1011$

$Y = 0111 = 7$

2 registers required   $Q(n+1 \text{ bit})$

A accumulator $(n$ bit$)$

Accumulator

$0000$    $Q$

$1011\underline{0}$ $= \{X, -1'b0\}$
    (appended)

## step

**0 (init)**

**1**   $A - Y$
        $R.S.$    ~~$1001$~~  $\rightarrow 1 1 0 1 1 Q$
        ~~$R.S. \quad 1100$~~
                $\uparrow$
            sign extension

**2.**   $R.S$    $1100 \rightarrow 0 1 1 0 1 X$
                $1110$
                ___
                $0101 \rightarrow 1 0 1 1 0 X$
                $0010$
**3.**   $A + Y$   ___
        $R.S.$    $10 01) \rightarrow 1 0 1 1 0$
                $110 0 \quad 1 1 0 1 X$
**4.**   $A - Y$
        $R.S$    $| 10 0 \quad $ product !!

$X : 2 = 0010$              John Hayes

$Y = 6 = 0110$

$-Y = 1010$

1
2
3
4