

NUMERO: 2419A

ANNO: 2019

# **A P P U N T I**

STUDENTE: Ricci Francesco

MATERIA: Microelectronics Systems - Prof. Mariagrazia  
Graziano

Il presente lavoro nasce dall'impegno dell'autore ed è distribuito in accordo con il Centro Appunti.

Tutti i diritti sono riservati. È vietata qualsiasi riproduzione, copia totale o parziale, dei contenuti inseriti nel presente volume, ivi inclusa la memorizzazione, rielaborazione, diffusione o distribuzione dei contenuti stessi mediante qualunque supporto magnetico o cartaceo, piattaforma tecnologica o rete telematica, senza previa autorizzazione scritta dell'autore.

**ATTENZIONE: QUESTI APPUNTI SONO FATTI DA STUDENTIE NON SONO STATI VISIONATI DAL DOCENTE.  
IL NOME DEL PROFESSORE, SERVE SOLO PER IDENTIFICARE IL CORSO.**

# MICROELECTRONICS HISTORY

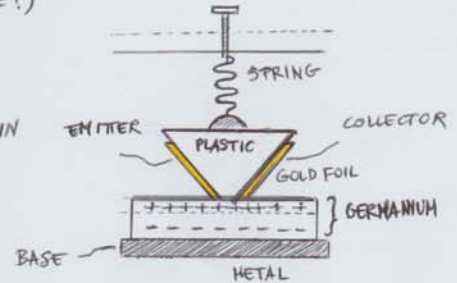
DIGITAL ELECTRONICS BORN WHEN RELAYS WERE SUBSTITUTED BY VACUUM TUBES. (1906)

1946: **ENIAC**: ELECTRICAL NUMERIC INTEGRATOR AND CALCULATOR.

COMPOSED BY: 17000 VACUUM TUBES, 7000 RESISTANCES, 10000 CAPACITORS, 6000 SWITCHES, 30 TONS, POWER CONSUMPTION OF 160 KW (HUGE!)

ABLE TO COMPUTE 5000 ADDITIONS IN 1 SECOND.

1949: BIPOLAR TRANSISTOR INVENTION: SHOCKLEY, BARDEEN, BRATTAIN (POINT-CONTACT TRANSISTOR)



1958: KILBY INTEGRATED CIRCUIT (TEXAS INSTRUMENTS)

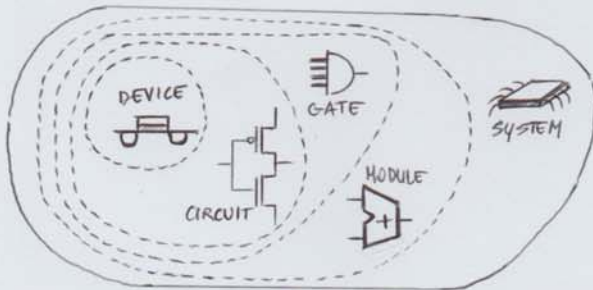
1971: 4004 INTEGRATED MICROPROCESSOR.

FIRST USAGE OF CMOS TECHNOLOGY (SILICON GATE): 10 μm NMOS LOGIC  
 COMPOSED BY 2300 TRANSISTOR, ONLY 3mm x 4mm (AREA)  
 2 METAL LAYERS, 4 BIT CPU, 100 KHz

2000: PENTIUM 4 INTEGRATED PROCESSOR

COMPOSED BY 54 MILLION OF TRANSISTORS, 1.5 GHz, 0.18 μm TECH

1965 MOORE'S LAW: INTEGRATION COMPLEXITY GROWS EXPONENTIALLY - THIS LAW SEEMS CONFIRMED UNTIL 2000.... WHY? WHAT DID ALLOW SUCH A PERFORMANCE IMPROVEMENT? **KEYWORD IS ABSTRACTION**



FIRST APPROACH: all circuit was done directly at the layout level from teams of designers which knew everything

ABSTRACTION: various HIERARCHY  
 SMALL BLOCKS ARE DEFINED CAREFULLY AT FULL CUSTOM LEVEL, AFTER CHARACTERIZED. PARAMETERS ARE DEFINED AND USED BY NEXT LEVEL TO BUILD BIGGER BLOCKS, .... AND SO ON....

**EACH LEVEL SOLVES IT'S OWN PROBLEMS!**  
**HIGHER LEVELS JUST USE A LOWER LEVEL REPRESENTATION**

TODAY THIS IS NOT ABSOLUTELY TRUE: STRIKING DOWN TRANSISTORS, NEW PROBLEMS OCCURS ACROSS MULTIPLE LEVELS... SO SOLUTIONS ON ONE LEVEL AFFECT ALSO OTHER LEVELS...

I.T.R.S (international technology roadmap of semiconductors) make some predictions, like:

YEAR	2010	2013	2018
NODE [nm]	45	32	18
METAL LAYERS	15	16	18
MAX FREQ [MHz]	15079	22480	53207
POWER [W]	120	138	158

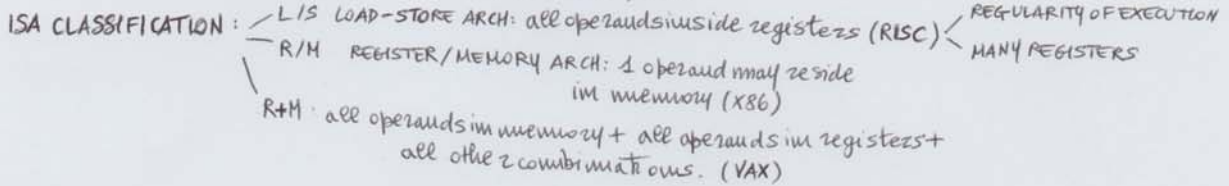
**SOI** SILICON ON INSULATOR is a huge improvement in tech: AVOID CAPACITANCE = INCREASE SPEED WHILE REDUCING POWER  
 VERY FREQUENTLY USED TODAY!



COMPLEX MICROELECTRONIC SYSTEMS

• ASIC: application specific integrated circuits ; composed by data memory, hierarchical finite state machine CU, arithmetic datapath with few dedicated registers.

• ISA: instruction set architectures  $\Rightarrow$  general purpose machines.  
 composed by a data memory and an instruction memory -  
 control unit update the program counter (PC): store the address of the next instruction that should be executed. That instruction is loaded into the IR (instruction register), used by the data path.



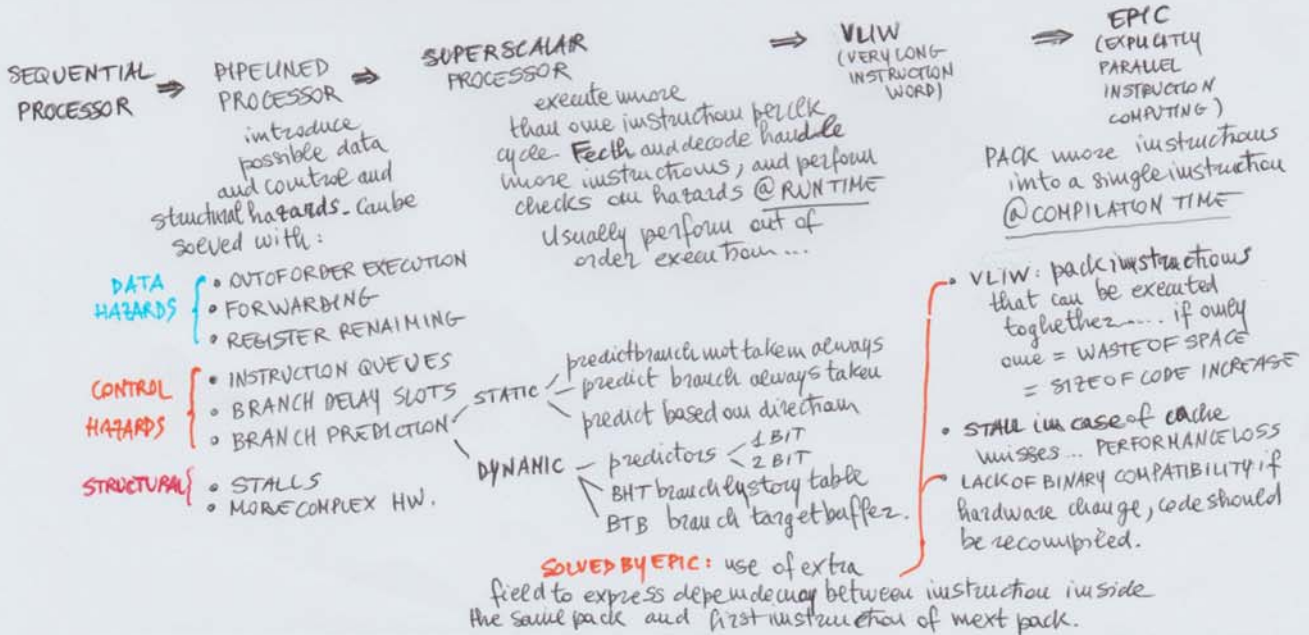
RISC VS CISC: cisc processors use a variable length instructions: in order to identify the correct ~~access~~ addressing mode, more complex logic is necessary. Decoding multiple instructions at the same time means prediction on the starting point and ending point of each. Extra logic means extra area and an higher power consumption.

- risc processors are based on two principles:
1. KEEP IT SIMPLE: a reduced instruction set is available because real programs spend 80% of the time executing 20% of possible type of instructions
  2. L/S: only registers as operands means fixed length instructions and simple addressing modes.

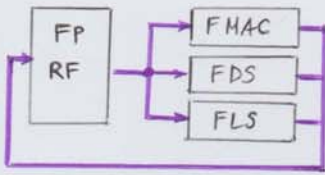
Historically cisc give the possibility to build better processors because of

- poor silicon budget (maximum number of transistors: if use too many for registers, you don't have enough for ALU, ...)
- high memory cost  $\Rightarrow$  because cisc allows usage of less number of instructions, total code is shorter.

Today risc processors are clearly better, but binary compatibility stops them to spread all over the market. cisc nowadays implements a transition inside, to exploit risc advantages, while maintaining binary compatibility.



## ARM VFPV4 COPROCESSOR



- FMAC: floating point multiply and accumulate
- FDS: floating point division and square root
- FLS: floating point load and store



SPARSE TREE

Carry look ahead (CLA) adder is the basis for every tree adder:

$a_i$	$b_i$	$c_i$	$S_i$	$C_{out}$	$p$	$g$	
0	0	0	0	0	0	0	GENERATED
0	1	0	1	0	1	0	PROPAGATED
1	0	0	1	0	1	0	PROPAGATED
1	1	0	0	1	0	1	GENERATED
0	0	1	1	0	0	0	GENERATED [KILLED]
0	1	1	0	1	1	0	PROPAGATED
1	0	1	0	1	1	0	PROPAGATED
1	1	1	1	1	0	1	GENERATED

$$s = a \oplus b \oplus c_{in}$$

$$C_{out} = a \cdot b + a \cdot c + b \cdot c$$

$$p = a \oplus b$$

$$g = a \cdot b$$

Let's do some math:

$$s_1 = a_1 \oplus b_1 \oplus c_0 = p_1 \oplus c_0$$

$$c_1 = a_1 \cdot b_1 + a_1 \cdot c_0 + b_1 \cdot c_0 = a_1 \cdot b_1 + (a_1 \oplus b_1) c_0 = g_1 + p_1 \cdot c_0$$

GENERATED OR PROPAGATED

↑ NB: I can write XOR instead of OR because in this particular case the outcome is the same:

a	b	OR	XOR
0	0	0	0
0	1	1	1
1	0	1	1
1	1	1	0

WHEN THIS HAPPEN, THE FIRST TERM ( $a_1 \cdot b_1$ ) IS EQUAL TO 1. 1 OR SOMETHING = 1, NO MATTER WHAT ( $a_1 + b_1$ )  $c_0$  IS EQUAL TO...

$$s_2 = a_2 \oplus b_2 \oplus c_1 = p_2 \oplus c_1$$

$$c_2 = g_2 + p_2 \cdot c_1 = g_2 + p_2 g_1 + p_2 p_1 c_0$$

SUBSTITUTE

$$s_3 = p_3 \oplus c_2$$

$$c_3 = g_3 + p_3 c_2 = g_3 + p_3 g_2 + p_3 p_2 g_1 + p_3 p_2 p_1 c_0$$

SUBSTITUTE

NOTES:

1. Virtually all carries can be calculated in function of  $c_0$  (carry in).
2. However this require a very complex net, that is not easily scalable
3. It introduce a little congestion (not a regular structure)  $\Rightarrow$  HIGHER DELAY
4. Each stage drive a different amount of gates: sizing is more complex
5. Different size means different delay and a more complex design phase.
6. Different configurations are possible: maybe  $c_5$  is written in function of  $c_1$  instead of  $c_0$  ...

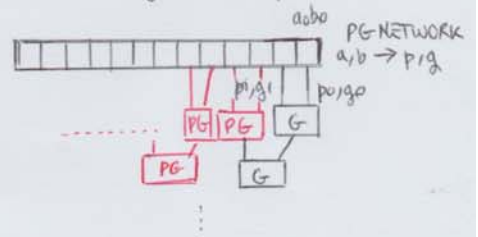
GENERAL PROPAGATE AND GENERAL GENERATE

By modifying values of  $i, j$  and  $k$ , the network can be constructed and carries can be computed.

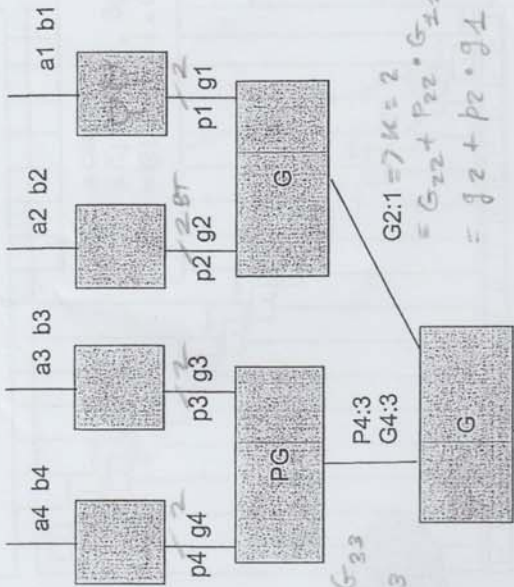
$$P_{i:j} = P_{i:k} \cdot P_{k-1:j}$$

$$i > k > j \quad G_{xx} = g_x \quad P_{xx} = p_x \quad g_0 = c_{in} \quad p_0 = 0$$

$$G_{i:j} = G_{i:k} + P_{i:k} \cdot G_{k-1:j}$$



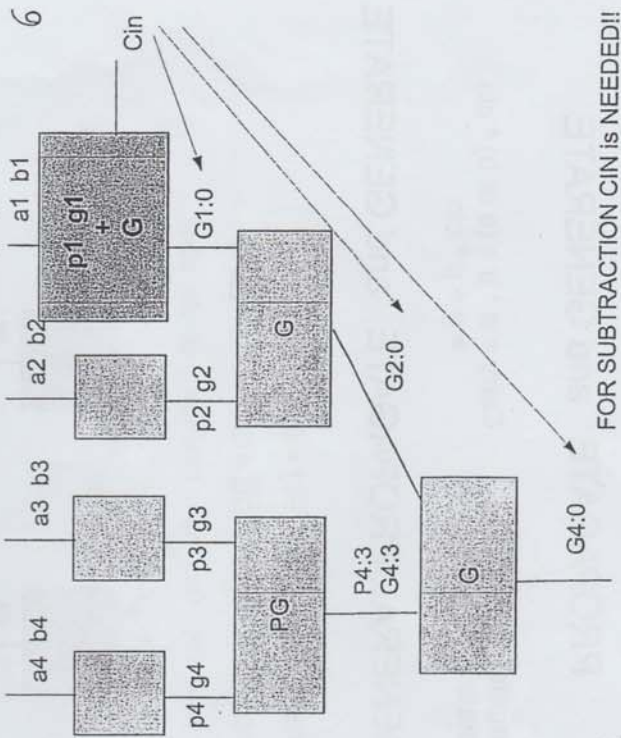
5



first rule  
 $P_{4:3} = P_{4,4} \cdot P_{3:3}$   
 $= P_{4,4} \cdot P_3$   
 $G_{4:3} = G_{4,4} + P_{4,4} \cdot G_{3:3}$   
 $= g_4 + p_4 \cdot g_3$

$G_{2:1} = g_2 = 2$   
 $= G_{2,2} + P_{2,2} \cdot G_{1:1}$   
 $= g_2 + p_2 \cdot g_1$

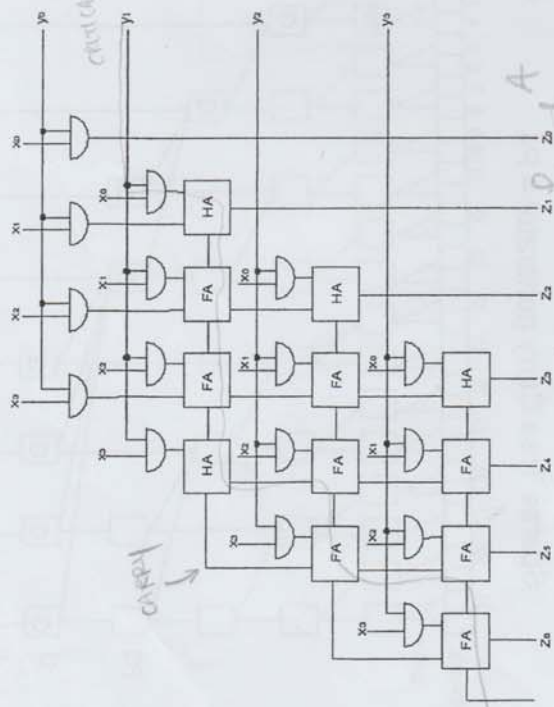
$G_{4:1} = G_{4,3} + G_{3:1}$   
 $P_{4:3} \cdot G_{2:1} = g_4 + p_4 \cdot g_3 + p_4 \cdot p_3 \cdot g_2 + p_4 \cdot p_3 \cdot p_2 \cdot g_1$   
 OK for ADDER



FOR SUBTRACTION CIN IS NEEDED!!

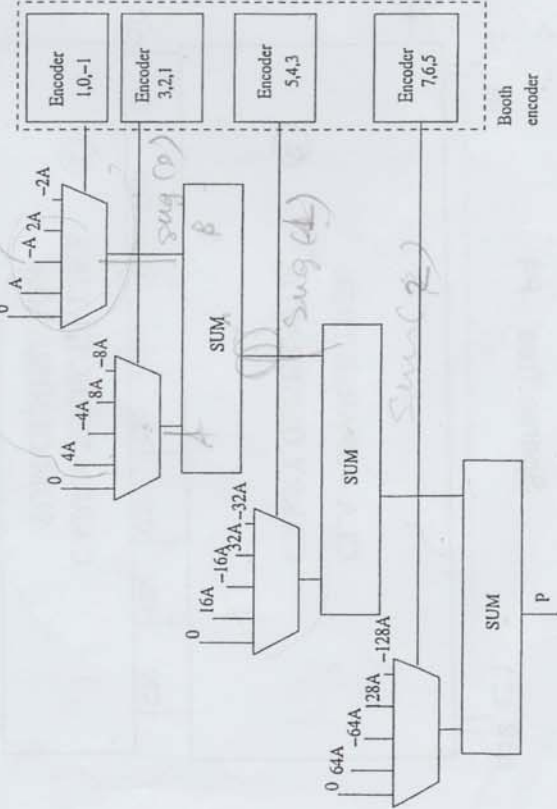
7

### Array multiplier



8

### Booth's multiplier



1231

### BOOTH'S ALGORITHM

Suppose  $M$  is even; by convention  $B_{-1} = 0$ . NB:  $M$  is the parallelism of multiplicand  $B$   
algorithm:

```

i = 0;
P = 0;
WHILE (i <= M-2) LOOP
    P ← P + Vp(Bi+1, Bi, Bi-1);
    A ← 4 * A;
    i ← i + 2;
END LOOP;
    
```

$y_{i+1}$	$y_i$	$y_{i-1}$	$V_p$
0	0	0	0
0	0	1	+A
0	1	0	+A
0	1	1	2A
1	0	0	-2A
1	0	1	-A
1	1	0	-A
1	1	1	0

A: MULTIPLIER

Example:  $B = 00011000$  (bits 7 6 5 4 3 2 1 0 -1)

$(i=0)$   $V_p(B_1, B_0, B_{-1}) = V_p(000) = 0$   
 $P \leftarrow P + 0 = 0$   
 $A \leftarrow 4A = 2^2A$   
 $i \leftarrow i + 2 = 2$

$(i=2)$   $V_p(B_3, B_2, B_1) = V_p(100) = -2A$   
 $P \leftarrow P + (-2A) = -2^3A$   
 $A \leftarrow 4A = 2^4A$   
 $i \leftarrow i + 2 = 4$

$(i=4)$   $V_p(B_5, B_4, B_3) = V_p(011) = 2A$   
 $P \leftarrow P + (2A) = -2^3A + 2 \cdot 2^4A = 2^5A - 2^3A$   
 $A \leftarrow 4A = 2^6A$   
 $i \leftarrow i + 2 = 6$

$(i=6)$   $V_p(B_7, B_6, B_5) = V_p(000) = 0$   
 $P \leftarrow P + 0 = P = 2^5A - 2^3A$   
 :

RIGHT!

WHAT HAPPENS IF  $B = 1001100$  (1 IN MSB...)?

In this case only last iteration changes:

$(i=6)$   $V_p(B_7, B_6, B_5) = V_p(100) = -2A$   
 $P \leftarrow P + (-2A) = -2^7A + 2^5A - 2^3A = -10^4A$

MULTIPLICATION STILL WORKS, BUT ONLY IF YOU CONSIDER NEGATIVE (2' COMPL.) THE MULTIPLICAND.

Key concept: any binary number can be written expressing addition or partial shift + subtraction:

EX:  $B = 00011000 = 2^3 + 2^4 = 48$   
 $\quad \quad \quad \uparrow \quad \uparrow = 2^5 - 2^3 = 48$

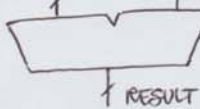
EX:  $B = 01011000 = 2^7 - 2^6 + 2^5 - 2^3 = 196$



27 MARZO  
PTZ  
REC CELL

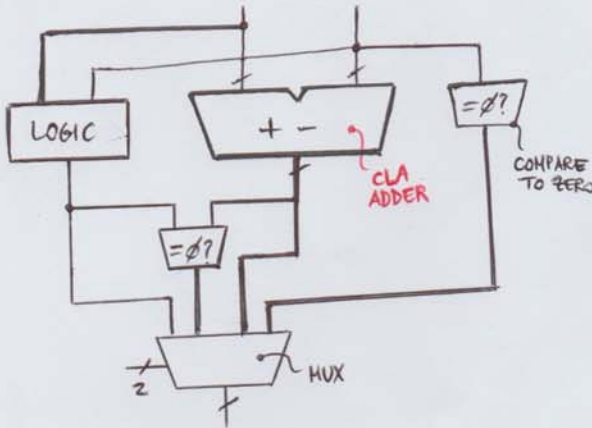
### GENERAL ALU

High level view: OP. 1 OPERAND 2



- ADD, SUB
- MUL, DIV
- LOGIC
- SHIFT, ROT
- COMPARISON (= != >> <<)
- MODULOS, SQUARE, ROOT

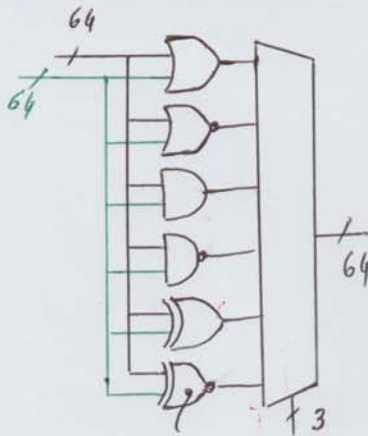
### ALU OF T2



In the T2 processor, the alu is separated from the shifter unit

Logic unit : it is the last piece missing from the above circuit.

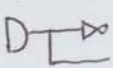
Obvious implementation :



they are ~~tree~~ of gates, or an array: in the second case I use 64x64 gates + a huge mux

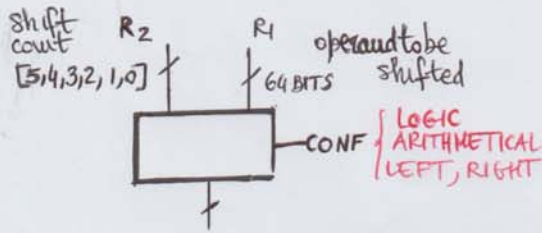
This is a very bad solution!

- 1) Mux has a tonne of logic!  $\leq$  AREA DELAY
- 2) Every gate has a different delay, for this reason I need to consider the slowest part (NAND is the fastest) (XOR is slowest)
- 3)

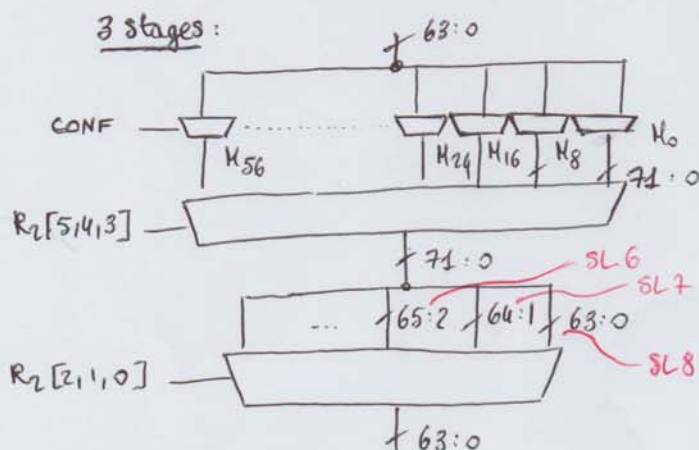
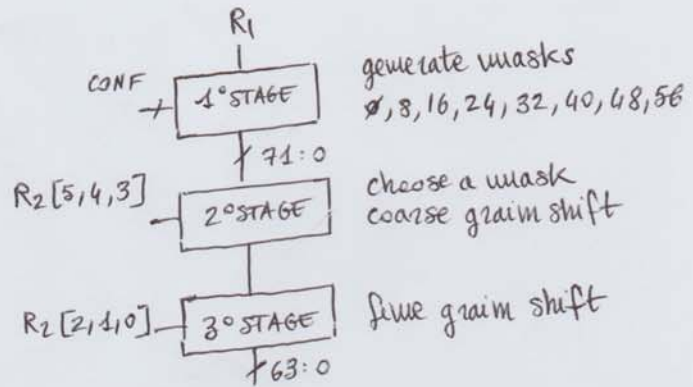
Why reducing the number of transistors using  is a bad solution? Because the and gate has to drive a lot (twice) the electronics! And for this reason it must be larger!

Oral question: implement this solution and discuss why it is bad.

## T2 SHIFTER.



Composed of 3 stages:



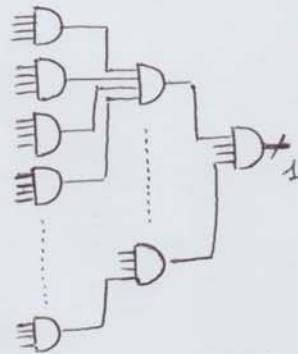
Correct implementation:

A good implementation is a multilevel (tree) one.

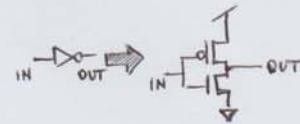
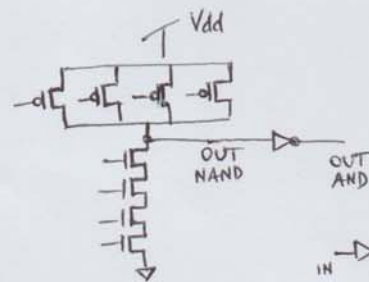
High level view:



→ using an AND gate with 4 inputs



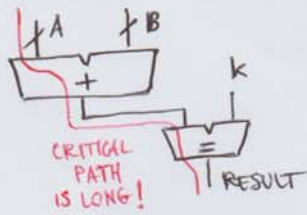
16 ANDS 4 ANDS 1 AND = 21 gates = 210 transistors



## EVEN MORE COMPLEX COMPARATOR

Now I want to perform  $A+B=k$  operation. Result = 1 if  $A+B=k$

Obvious implementation



A big disadvantage is the delay due to carry propagation inside the adder.

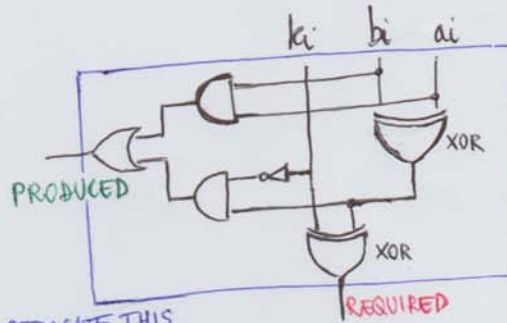
Better implementation

$a_i$	$b_i$	$k_i$	REQUIRED CARRY IN $C_{i-1}$	PRODUCED CARRY OUT $C_i$
0	0	0	0	0
0	0	1	1	0
0	1	0	1	1
0	1	1	0	0
1	0	0	1	1
1	0	1	0	0
1	1	0	0	1
1	1	1	1	1

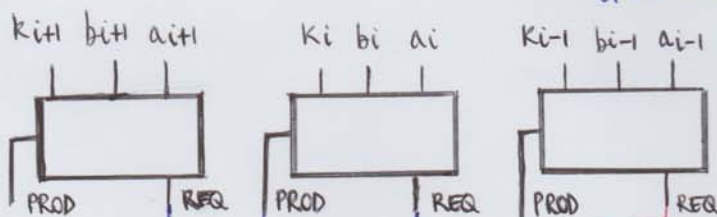
We don't demonstrate it, but we can write:

$$C_{i-1} = a_i \oplus b_i \oplus k_i \quad \text{REQUIRED}$$

$$C_i = (a_i \oplus b_i) \cdot \bar{k}_i + a_i \cdot b_i \quad \text{PRODUCED}$$



REPLICATE THIS BLOCK FOR EVERY BIT OF THE TWO OPERANDS



SHOULD BE CONSIDERED DURING THE EVALUATION!

if all produced carries are equal to the next required (XNOR tests equality)

should be a tree!  $\downarrow$  if  $A+B=k$

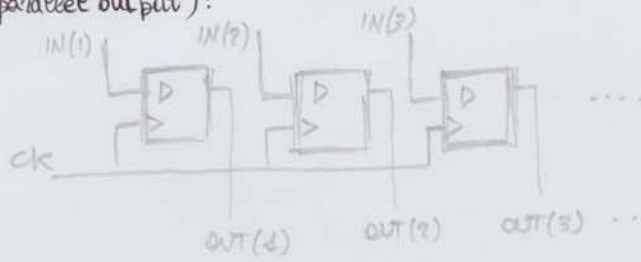
Advantage: I process every bit in parallel. A ripple carry adder process each bit in a sequential way.



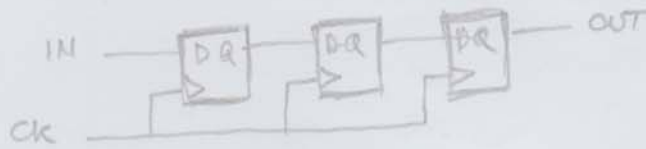
## REGISTERS

Based on how input and outputs are connected, we can classify registers as 4 groups:

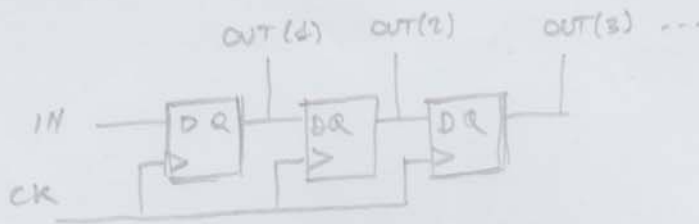
- PIPO (parallel input, parallel output):



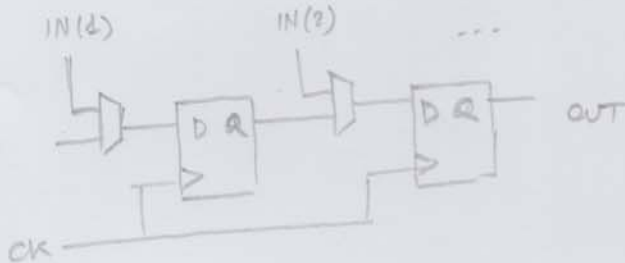
- SISO (serial input, serial output):



- SIPO:



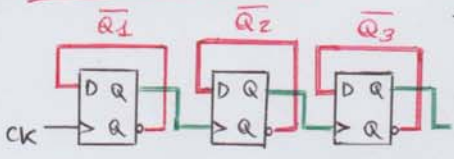
- PISO:



**COUNTERS**

- ASYNCHRONOUS
- SYNCHRONOUS
- JOHNSON
- HA COUNTER

**ASYNCHRONOUS COUNTER**

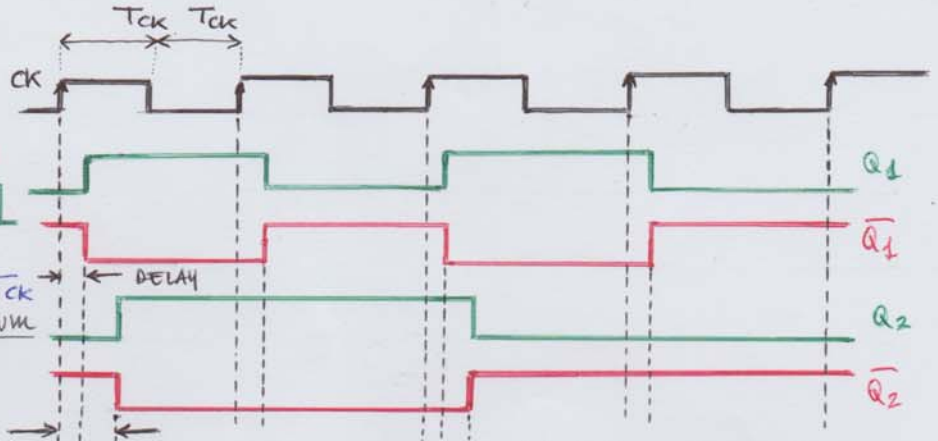


if use in parallel  $Q_1 Q_2 Q_3$   $\uparrow$  CK you obtain an asynchronous down counter!

**BEWARE OF DELAYS!**

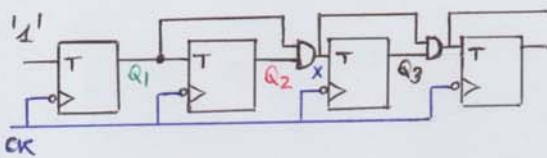
$$T_{Q1} = 2T_{CK}$$

$$T_{Q2} = 2T_{Q1} = 4T_{CK}$$

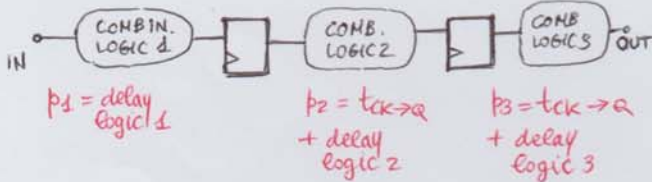


**SUM OF DELAYS!** **WARNING:** by adding bits to the counter structure, the delay sum at the end will increase. Maybe, if the clock period is too short, there will be a setup time violation!  $\Rightarrow$  NOT EASY TO SCALE  $\Rightarrow$  EASY TO BUILD

**SYNCHRONOUS COUNTER**



In general:

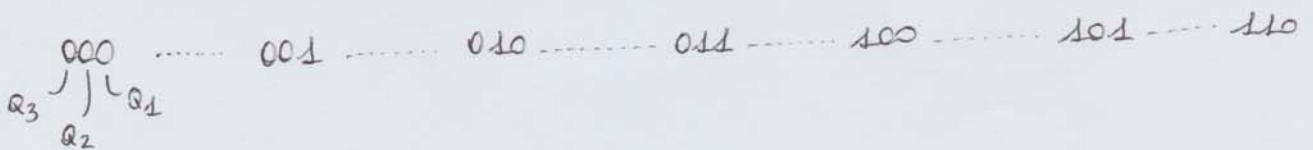
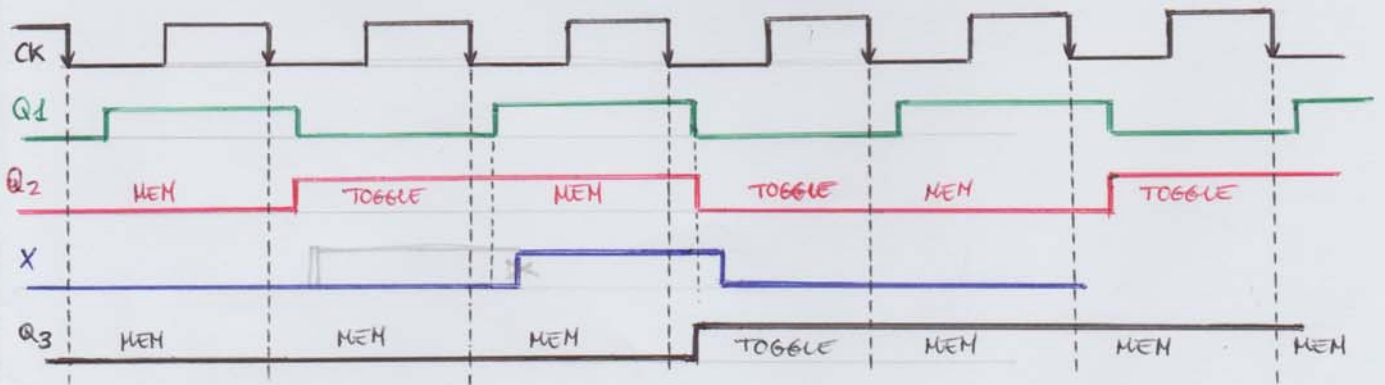


**SEARCH THE CRITICAL PATH:** a path should always start from the out of FF to the IN of another FF (or the same one, if there are paths).

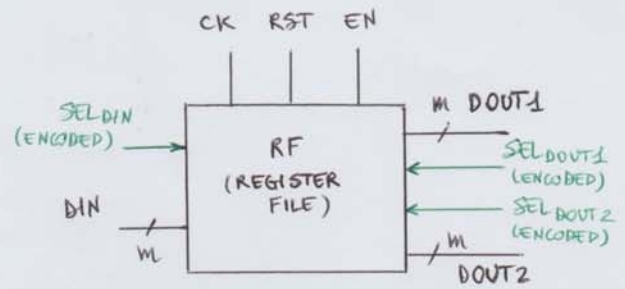
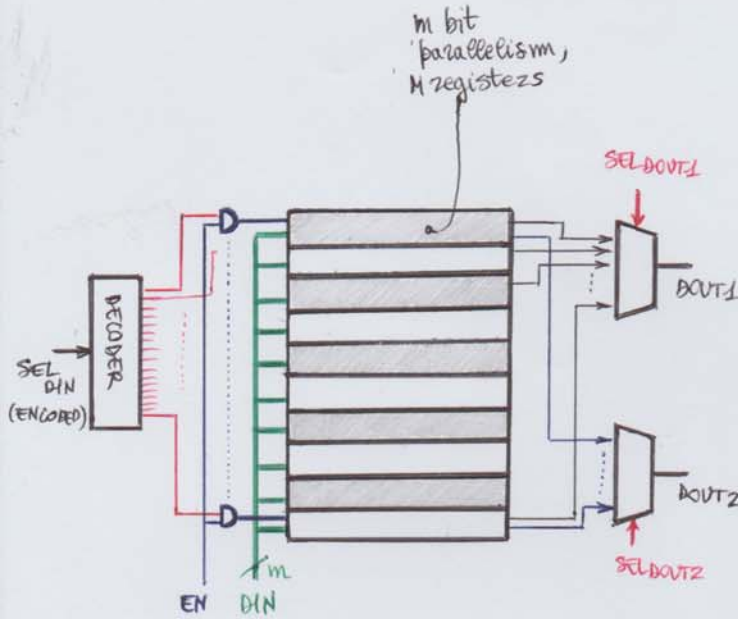
Supposing  $t_{AND} = 3\mu s$ ,  $t_{FF}^{CK \rightarrow Q} = 5\mu s$ , paths are:

- $t_{p1} = 5\mu s$
- $t_{p2} = 5\mu s + 3\mu s = 8\mu s$
- $t_{p3} = 5\mu s + 3\mu s = 8\mu s$
- $t_{p4} = 5\mu s + 3\mu s = 8\mu s$
- $t_{p5} = 5\mu s + 3\mu s + 3\mu s = 11\mu s$  **CRITICAL PATHS!**
- $t_{p6} = 5\mu s + 3\mu s + 3\mu s = 11\mu s$  **CRITICAL PATHS!**

$$MAX\_delay = MAX\{p_1, p_2, p_3\} \Rightarrow MIN\_freq = \frac{1}{MAX\_delay + SETUP\_TIME}$$



## REGISTER FILE



NB: a real RF is very different. The structure is heavily optimized to be very fast!

### NOTES:

1. Conflict between input and output path! When is requested to read and write the same register, which is the order to follow? Is important to define a **policy**!  $\Rightarrow$  WAR or RAW  
Another solution could be a **double data rate**  $\Rightarrow$  half clock cycle read, half clock cycle write.
2. Decoder recall: it is a combinational logic circuit that converts binary information from  $m$  coded inputs to a maximum of  $2^m$  unique outputs.
3. Problem of context switching: if a procedure is called (by using assembler instruction CALL/RET) two actions should be performed:
  - a) move all registers from RF to the stack  $\Rightarrow$  CALL or move all registers from stack to the RF  $\Rightarrow$  RET
  - b) save in the stack the value of flags  $\Rightarrow$  CALL or restore from stack the value of flags  $\Rightarrow$  RET

NB: context switching creates a long "bubble" in the pipeline.

Often a procedure requires few registers to accomplish the task, but every time CALL/RET are used, the entire RF should be copied to the stack. A solution is windowing.

WINDOWING  $\rightarrow$  FIXED: size of windows are fixed  
 $\rightarrow$  FLEXIBLE: size of windows variable

Instead of moving the entire RF content every time, I just move a small portion of it (a window).



**REGISTER FILE FULL :** In case the RF is full (all windows are occupied) and a CALL arrive, a **SPILL** should be performed. A spill moves the oldest window in the stack, in order to free space for the new procedure.

When a RET arrive, if one or more SPILL were performed, a **FILL** should be done, in order to restore the replaced window from the stack.



In order to understand when perform a SPILL/FILL we use 2 pointers:

- CWP: current window pointer
- SWP: saved window pointer.

SUB0:  
operations  
CALL SUB1 → SHIFT CWP  
SHIFT SWP IS NOT NEEDED! ONLY AFTER A SPILL!

SUB1:  
operations  
CALL SUB2 → SHIFT CWP  
SHIFT SWP IS NOT NEEDED!

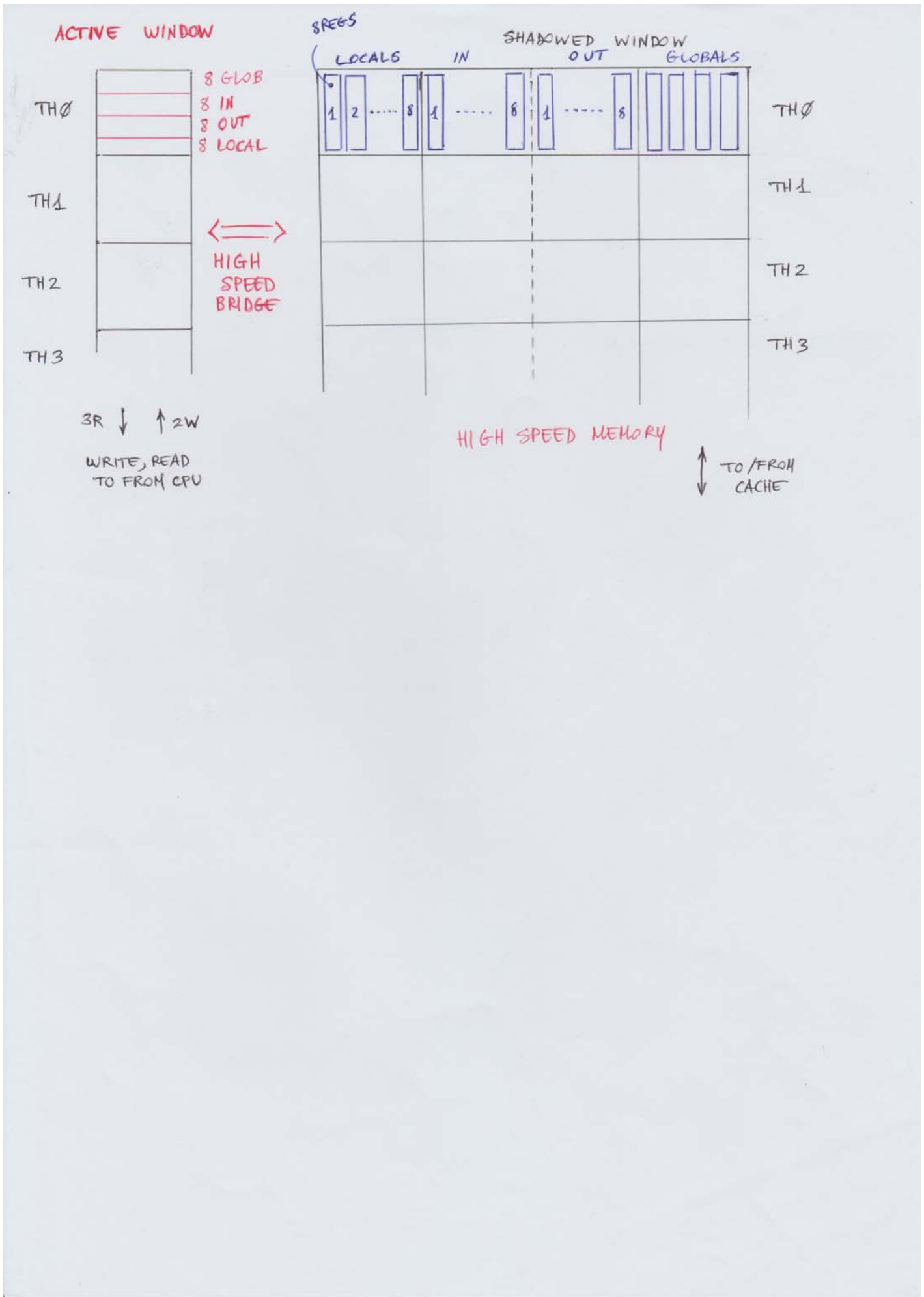
SUB6:  
operations  
CALL SUB7 ← FREE SPACE IS NOT ENOUGH! NEED A SPILL  
SPILL W0 IN and W0 LOC in STACK (or MEMORY)  
SHIFT CWP TO SUB7 ⇒ NO RF USE W7 W7 W7 IN LOC OUT  
SHIFT SWP TO REMEMBER THAT W0 IN and W0 LOC ARE IN MEMORY

SUB7:  
operations  
CALL SUB8 ← NO FREE SPACE! SPILL REQUESTED  
SPILL W1 IN and W1 LOC IN MEMORY  
SHIFT CWP TO SUB8 (OLD SUB0)  
SHIFT SWP TO SWP1

SUB8:  
operations  
RET FROM SUB8 (TO SUB7)  
SHIFT BACK CWP, DON'T TOUCH SWP!  
SUB7:  
RET  
SUB6:  
SUB3:  
RET FROM SUB3  
SHIFT BACK CWP.  
CHECK: SWP = CWP? YES, BECAUSE CWP2 = SWP1  
→ IF YES NOTHING TO DO NOW.  
SUB2:  
RET FROM SUB2  
⚠ SWP1 >> CWP1 ← perform a FILL!  
FILL FROM STACK W1 IN and W1 LOC  
SHIFT BACK CWP TO CWP1  
SHIFT BACK SWP TO SWP0

**SUMMARY:** I need those operations:  
1) SHIFT BACK AND FORTH CWP & SWP  
2) COMPARISON BETWEEN CWP & SWP  
3) SPILL / FILL ⇒ COMMUNICATE WITH A MMU (memory management unit)!



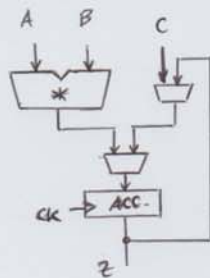


### ARM VFP11 COPROCESSOR

As already said before, it is composed by

- LOAD AND STORE UNIT
- F.P. MULTIPLY AND ACCUMULATE
- F.P. DIVIDE AND SQUARE ROOT

#### MULTIPLY AND ACCUMULATE : working principle :



$$z_0 = A_0 \times B_0$$

$$z_1 = z_0 + A_1 \times B_1$$

$$\vdots$$

$$z_{i+1} = \sum_{i=0}^m A_i \times B_i$$

#### DIVIDER : ARM COPROCESSOR IMPLEMENTS A RADIX 4 SRT ALGORITHM.

HOWEVER LET'S TAKE SOME STEPS TO UNDERSTAND THAT ALGORITHM:

- A) BASIC PRINCIPLE : ITERATIVE (THEORY, NO HW IMPL. POSSIBLE)
- B) RESTORING ALGO
- C) NON RESTORING + REDUNDANT ALGO
- D) SRT
- E) RADIX-4
- F) ARM IMPLEMENTATION...

#### (A) BASIC : ITERATIVE

DIVIDEND  $Z = 14_{10} = 1110_2$   
 DIVISOR  $D = 3_{10} = 11_2$

**BASE 10**

$$\begin{array}{r} 14 \overline{) 3} \\ 0 \phantom{0} \\ \hline 14 \phantom{0} \\ 12 \phantom{0} \\ \hline 2 \end{array}$$

PARTIAL REMINDER: 14  
 REMINDER: 2

QUOTIENT: 04

$$\frac{Z}{D} = 4 + \frac{2}{D} \cdot 10^0$$

**BASE 2**

$$\begin{array}{r} 1110 \overline{) 11} \\ 11 \phantom{00} \\ \hline 001 \phantom{00} \\ 00 \phantom{00} \\ \hline 10 \phantom{00} \\ 00 \phantom{00} \\ \hline 10 \phantom{00} \end{array}$$

P.R.: 001  
 P.R.: 10  
 REMINDER: 10

QUOTIENT: 100

$$\frac{Z_2}{D_2} = 100 + \frac{10}{D_2} \cdot 2^0$$

ITERATIVE ALGO : GUESS Q. DIGIT, MULTIPLY, FIND P.R., GO BACK TO POINT 1.  
 REALLY HARD IN HW      EASY IN HW

FORMALLY THE ALGORITHM IS:  $R^{j+1} = \tau R^j - q_{j+1} \cdot D$

$$R^{j+1} = 2 \cdot R^j - q_{j+1} \cdot D$$

BINARY ALGORITHM.

WHERE  $R^j$  PARTIAL REMAINDER ( $R^0 = Z = \text{DIVIDEND}$ )  
 $D$  DIVISOR  
 $q_j$  QUOTIENT DIGIT  
 $\tau$  #DIGIT OF THE ALPHABETH (BINARY  $\tau = 2$ )

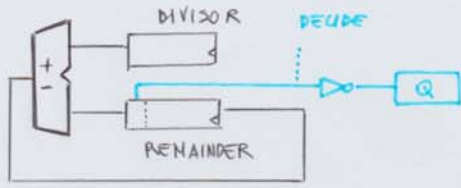
TO SIMPLIFY ASSUME:

$Z, D$  POSITIVE AND NORMALIZED AND  $Z < D$  (OBTAIN SO 0. ....)

$$Z = R^0 \rightarrow \cdot r_1^0 r_2^0 r_3^0 r_4^0 \dots r_m^0$$

$$D \rightarrow \cdot d_1^0 d_2^0 d_3^0 \dots d_m^0$$

$$Q \rightarrow \cdot q_1 q_2 \dots q_m$$

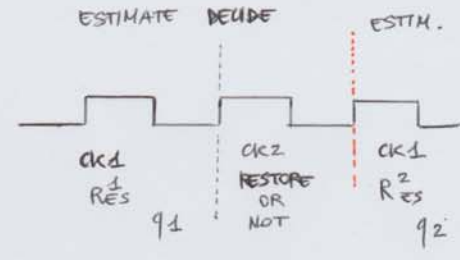


if SIGN 0  
MEANS  
 $R_{ES} > 0$   $q_i = 1$

---

if SIGN 1  
MEANS  
 $R_{ES} < 0$   $q_i = 0$

RESTORE  
BY SUMMING D!



let's verify:  $A: 0.01010110 = R^0 = Z$

$B: 0.1100 = D \Rightarrow \bar{D} \begin{array}{r} 1.0011 \\ + 1 \\ \hline 1.0100 \end{array}$

You have to subtract D, that in electronics means sum the 2's complement  $\bar{D}$

$2R^0 = 0.1010110 +$   
 $-D \quad 1.0100$   
 $R^1_{ES} \quad \textcircled{1} 1110110$

ASSUME  $q_1 = 1$

DECIDE  $\rightarrow$  IT'S NEGATIVE! RESTORE  
CORRECT  $q_1 = 0$

$R^1 = 2R^0$  AND NOT  $2R^0 - D \dots$

$2R^1 = 1.010110 +$   
 $-D \quad 1.0100$   
 $R^2_{ES} \quad \textcircled{0} 100110$

ASSUME  $q_2 = 1$

DECIDE  $\rightarrow$  IT'S POSITIVE! NO RESTORE  
 $q$  IS OK:  $q_2 = 1$

$R^2 = R^2_{ES} = 2R^1 - D$

$2R^2 = 1.00110 +$   
 $-D \quad 1.0100$   
 $R^3_{ES} \quad \textcircled{0} 01110$

ASSUME  $q_3 = 1$

DECIDE  $\rightarrow$  IT'S POSITIVE! NO RESTORE  
 $q$  IS OK:  $q_3 = 1$

$R^3 = R^3_{ES} = 2R^2 - D$

$2R^3 = 0.1110 +$   
 $-D \quad 1.0100$   
 $R^4_{ES} \quad \textcircled{0} 0010$

ASSUME  $q_4 = 1$

DECIDE  $\rightarrow$  IT'S POSITIVE! NO RESTORE  
 $q$  IS OK:  $q_4 = 1$

$R^4 = R^4_{ES}$

$Q = 0.0111$   
 $R = 0.0010 \cdot 2^{-4}$  # ITERATIONS

Let's verify:  $A: 0.01010110 = z$   
 $B: 0.1100 = D \Rightarrow \bar{D} = 1.0100$

$R^0$   $\textcircled{0} 01010110$  → IT'S POSITIVE  $q_1 = 1$   
 $2R^0 = 0.1010110 +$   
 $-D \quad 1.0100$

$R^1$   $\textcircled{1} 1110110$  → IT'S NEGATIVE COMPENSATE  $q_2 = \bar{1}$

$2R^1 = 1.110110 +$   
 $+D \quad 0.1100$

$R^2$   $\textcircled{0} 100110$  → IT'S POSITIVE  $q_3 = 1$

$2R^2 = 1.00110 +$   
 $-D \quad 1.0100$

$R^3$   $\textcircled{0} 01110$  → IT'S POSITIVE  $q_4 = 1$

OF COURSE WITH ANOTHER OPERATION, YOU OBTAIN THE SAME RESULT OF THE PREVIOUS EXAMPLE

$Q = 0.1\bar{1}11 \Rightarrow 0.1011 - 0.0100 = 0.0111$   
 $R = 0.0111 \cdot 2^{-3}$

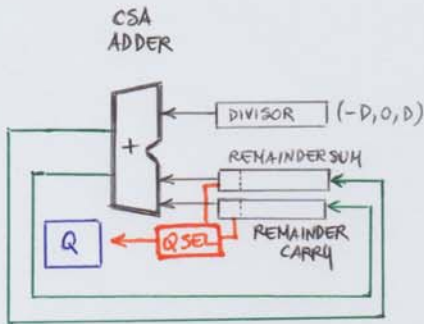
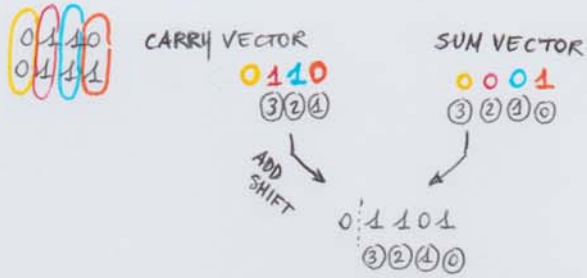
SAME RESULT



Basic CSA algorithm:

$$\begin{array}{r}
 A: 0110 \\
 B: 0111 \\
 \hline
 A+B = 1101
 \end{array}$$

LET'S CALL IT CARRY VECTOR. THE PRINCIPLE IS SAVE CARRY BEFORE RUN TIME, KEEPING THE FINAL SUM THAT REQUIRES CARRY PROPAGATION ONLY AT THE END.



(E) RADIX 4

Operands on  $m$  bits

$$\begin{array}{l}
 r = 2^b \\
 \text{RADIX} \quad K = \frac{m}{b}
 \end{array}$$

WHERE  $Q$  comprises (?)  
 $K$  is radix  
 $r$  are digits

GENERAL NAME: RADIX- $r$  ALGORITHM, NEEDS  $K$  ITERATION TO COMPLETE

IN OUR CASE RADIX-4:  $r = 2^2$ ,  $K = \frac{m}{2}$

RADIX-8: faster but much more complex

# C.U. IMPLEMENTATION

3

1) CLASSIC F.S.M. (ENCODED STATES)

2) HARDWIRED

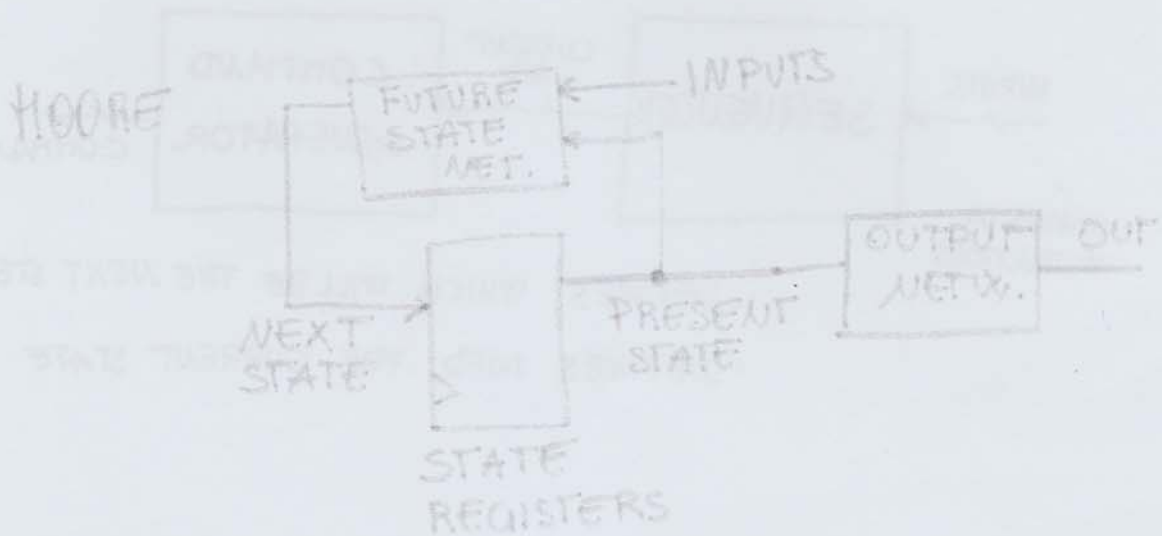
A FSM (NOW ENCODED STATES) WHERE THE SEQUENCER IS A COUNTER

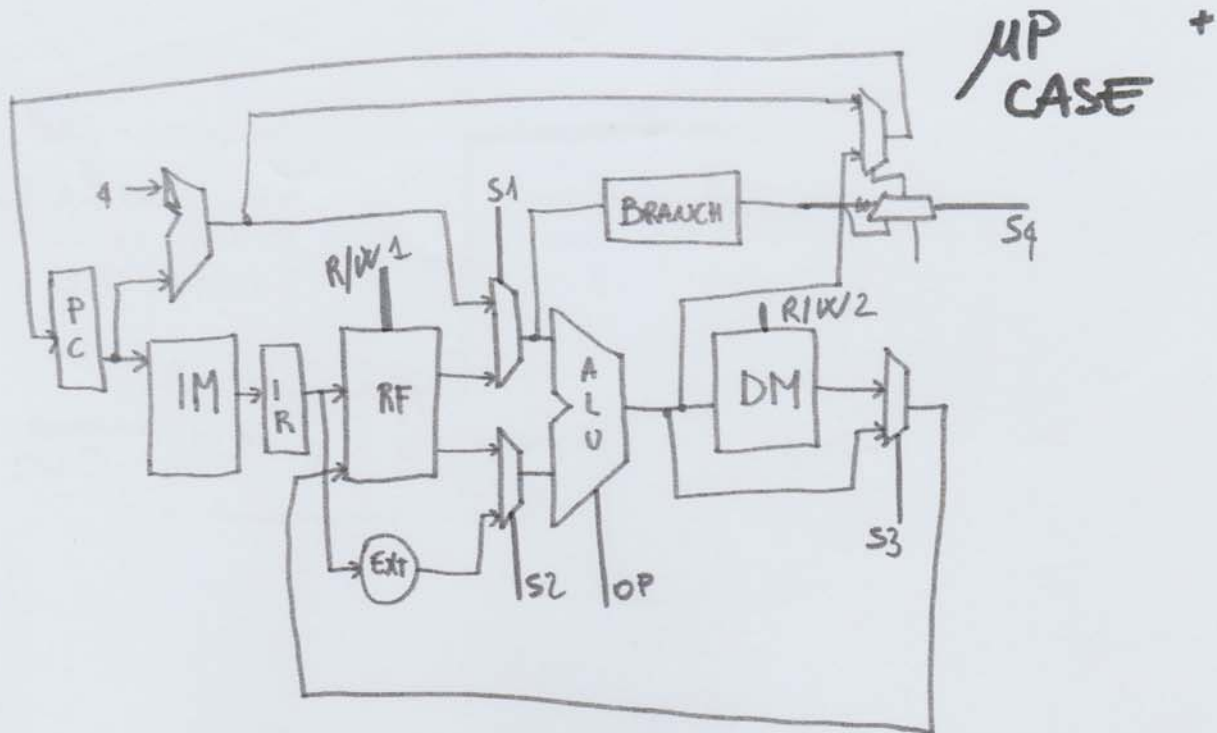
3) MICROPROGRAMMED

## A FLASH BACK TO F.S.M.

9

2 TYPES : { MOORE  
MEALY



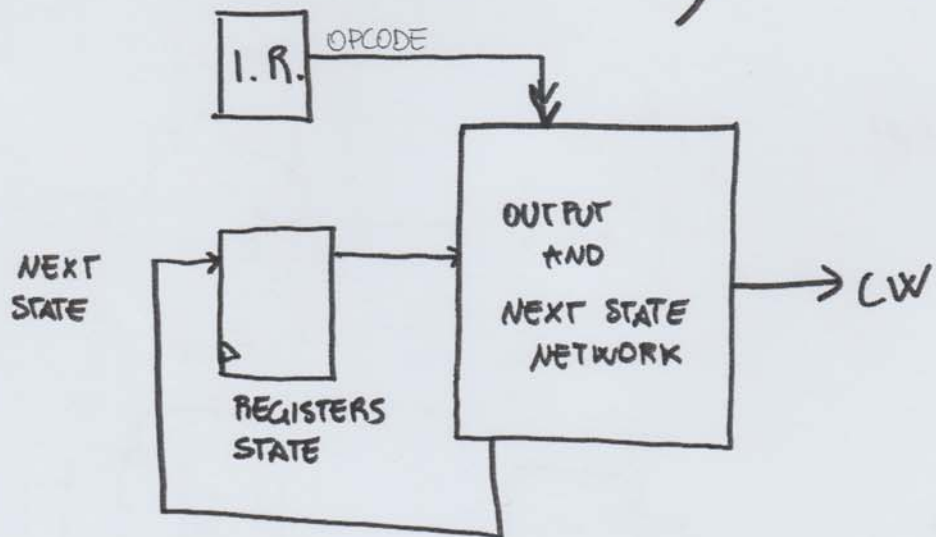


DLX CONTROL SIGNALS → CW CONTROL WORD

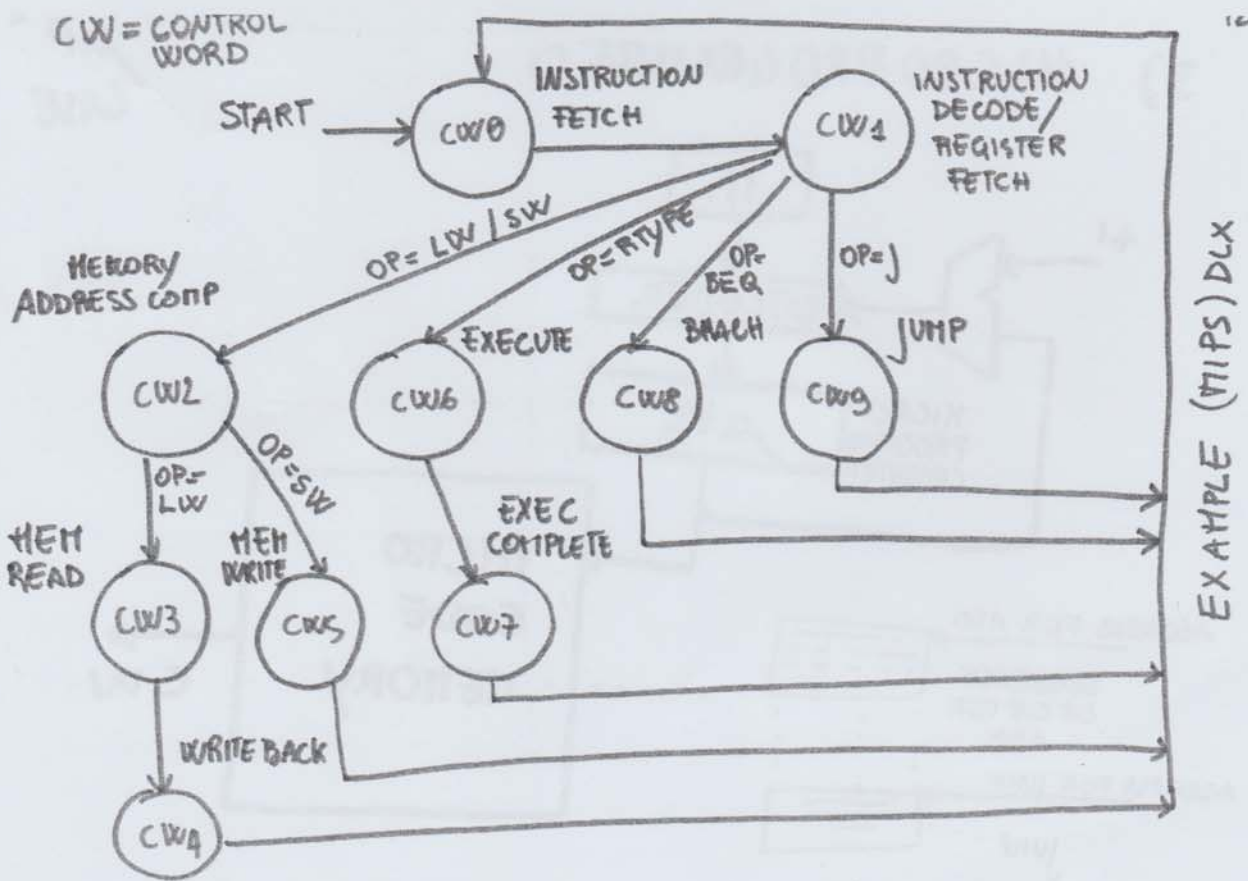
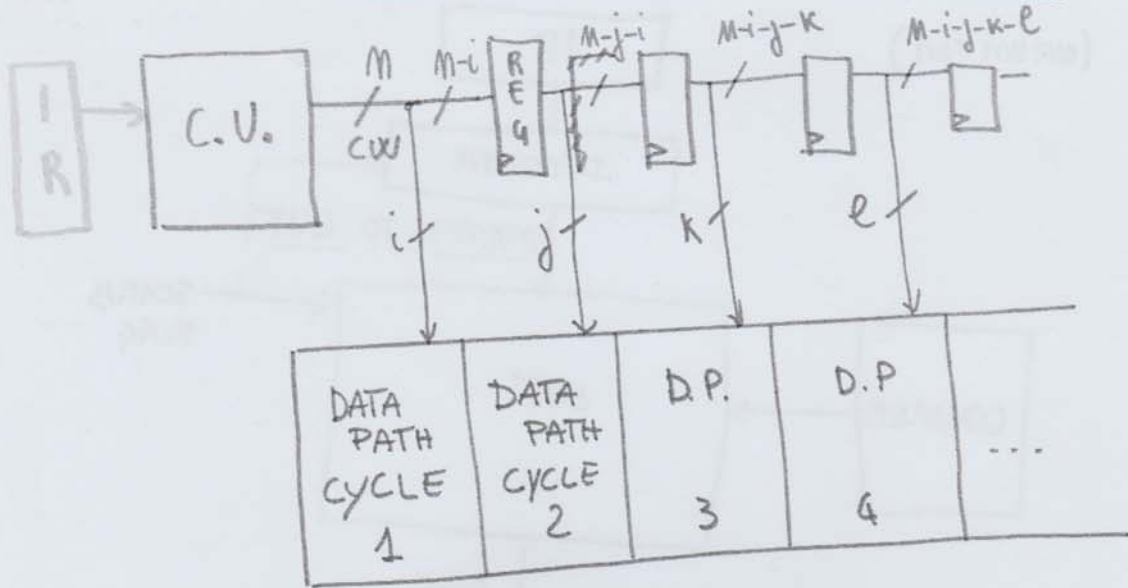
CW 01100101100.....  
 01 10 01 01 10 0 .....  
 S1 S2 R/W1 S3 S4

1) FSM

MP CASE :



# WHAT IF PIPELINED?

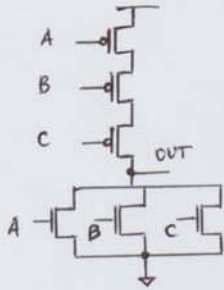




Homework: using CMOS logic family, implement the following logic functions:

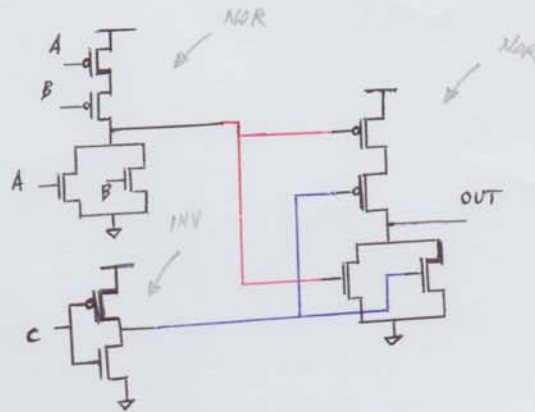
①  $OUT = \overline{A+B+C}$  ②  $OUT = (A+B)(C+D)$  ③  $OUT = A + (B + \overline{C}D)$  <sup>at most</sup>  
 by using a single stage and then a cascade of stages (each one with 2 inputs).

①  $OUT = \overline{A+B+C}$  it is a 3 input NOR



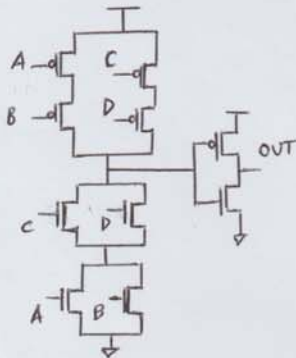
Using de Morgan I can manipulate the expression:

$$\overline{A+B+C} = \overline{A \cdot B \cdot C} = \overline{A+B} \cdot \overline{C} = \overline{\overline{A+B+C}}$$

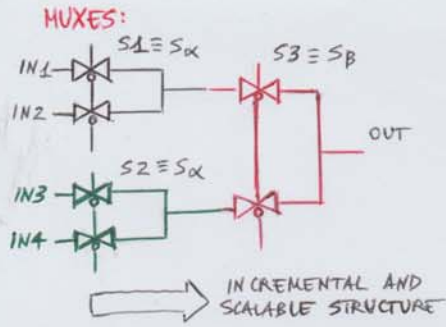


②  $OUT = (A+B)(C+D)$  using de Morgan:  $\overline{\overline{A+B}} \cdot \overline{\overline{C+D}} = \overline{\overline{A+B+C+D}}$

Since CMOS is inverting, I need an inverter at the end!



76/04



**OPTIMIZATION:** we can reduce the number of selection signals:

$$S1 = S2 = Sx$$

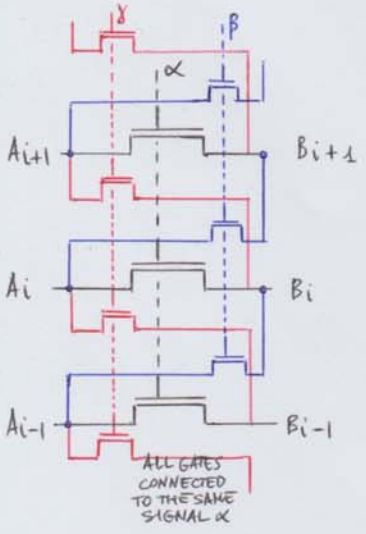
$$S3 = Sp$$

**IMPORTANT:**  $Sx$  and  $Sp$  should be inverted and connected to the other side of TG. The inversion is characterized by a small delay! Here isn't so important, but in LATCHES AND FFs based on TG it matters a lot!

**SHIFTERS:** three different implementations will be proposed:

- BINARY SHIFTER
- BARREL SHIFTER
- LOGARITHMIC SHIFTER

**BINARY SHIFTER:**



**NB:** to simplify the drawing, TG are reported as

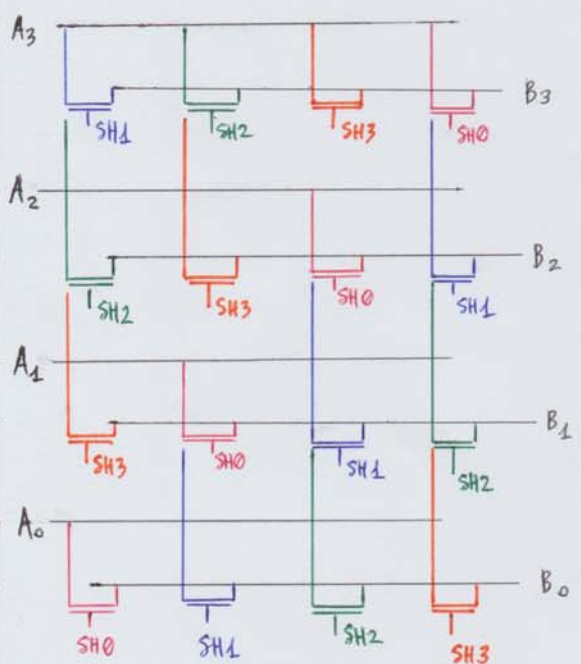
- IF  $\alpha = 1$  AND  $\beta = 0$  AND  $\gamma = 0 \Rightarrow Bi = Ai \Rightarrow$  NOP
- IF  $\alpha = 0$  AND  $\beta = 1$  AND  $\gamma = 0 \Rightarrow Bi = Ai-1 \Rightarrow$  SHL
- IF  $\alpha = 0$  AND  $\beta = 0$  AND  $\gamma = 1 \Rightarrow Bi = Ai+1 \Rightarrow$  SHR

**COMMENTS:** this is a very simple and compact shifter. It is possible to add more connections in order to perform shifts of more than one position in a single pass. More connections = more control signals.

**IMPORTANT:** always only once signal should be = 1

Usually is required a more configurable/programmable shifter ...

**BARREL SHIFTER:**



**NB:** to simplify the drawing, TG are reported as

- IF  $SH0 = 1$  AND all other signals = 0  $\Rightarrow$  SH0 = NOP
- IF  $SH1 = 1$  " " " "  $\Rightarrow$  SH1
- IF  $SH2 = 1$  " " " "  $\Rightarrow$  SH2
- IF  $SH3 = 1$  " " " "  $\Rightarrow$  SH3

**COMMENTS:** Regularity helps a lot in area compactness. Here in particular, is easy to build a transistor matrix, and then connect to input and output. Suppose a SHIFT REGISTER OF 8 BITS, you'll need 4x transistors and 4x area. To enable SHR, you'll need 2x transistors ...

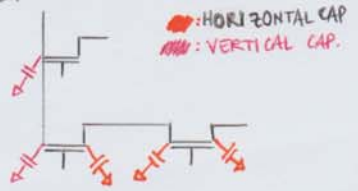
- NOTES:**
- 1) long lines: a line has a parasitic capacitance! longer the line, bigger the capacitance. Transistors connected to long lines should be correctly sized (BIGGER = FASTER). on the line
  - 2) parasitic capacitance of transistors: also those capacitance are summed with the line capacitance

**NB:** SH0, SH1, SH2, ... SHOULD DRIVE A LOT OF GATES  $\Rightarrow$  LOAD = SUM OF GATE CAPS.

$\downarrow$  DROP  $V$  CASE  $\Rightarrow$   $\downarrow$  PARASITIC RESIST.

PARASITIC CAPACITANCES

- HORIZONTAL LINES
- VERTICAL LINES
- DRAIN/SOURCE/GATE OF TRANSISTORS



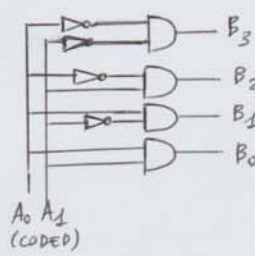
### DECODERS :

RECALL :



Combinational logic that converts binary information from the  $m$  coded inputs to a maximum of  $2^m$  unique outputs.

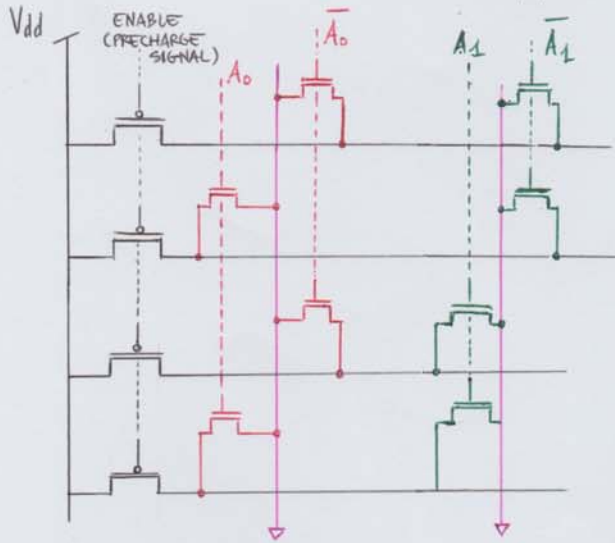
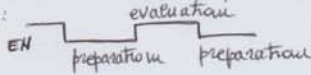
CMOS IMPLEMENTATION :



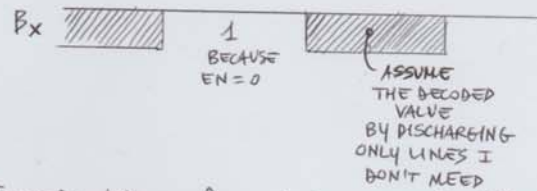
$$4 \text{ ANDS} = 24 \text{ TRANS} + 4 \text{ INV} = 8 \text{ TRANS} = 32 \text{ TRANS}$$

PASS TRANSISTORS IMPLEMENTATION :

NB: STATIC TOPOLOGY  $\equiv$  COMBINATIONAL  
 DYNAMIC TOPOLOGY  $\equiv$  COMBINATIONAL, BUT CHARACTERIZED BY A CONTROL SIGNAL. IT WORKS USUALLY IN 2 PHASES:



$A_0$	$\bar{A}_0$	$A_1$	$\bar{A}_1$	$B_3$	$B_2$	$B_1$	$B_0$
0	1	0	1	0	0	0	1
1	0	0	1	0	0	1	0
0	1	1	0	0	1	0	0
1	0	1	0	1	0	0	0

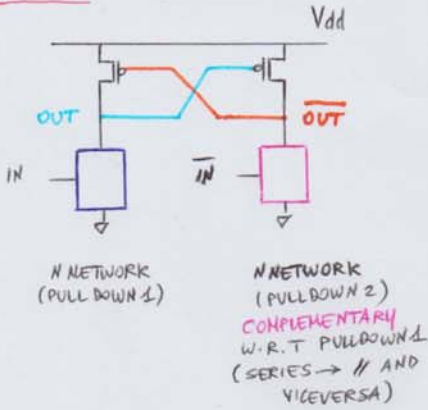


DRIVING CAPABILITIES: every signal should drive 2 transistors. Parasitic capacitance = caps of line + caps of transistor

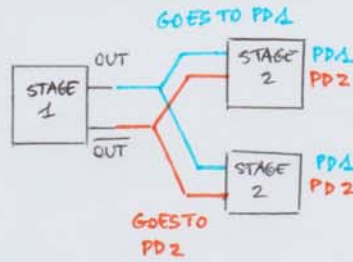
Typical problem of TG (out not perfect Vdd or GND) is avoided due to direct connection between Vdd or GND to output



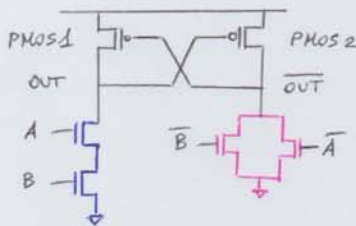
DC VSL : differential CMOS...



- NOTES:
- 1) TOPOLOGIES OF PD1 AND PD2 ARE COMPLEMENTARY
  - 2) INPUTS IN AND IN-bar ARE INDEPENDENT AND GENERATE FROM PREVIOUS STAGES: NO NEED FOR AN INVERTER!

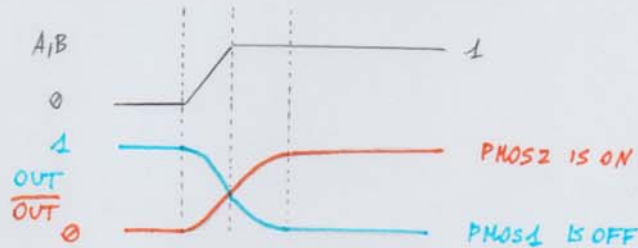


Example:



$$OUT = \overline{A \cdot B}, \quad \overline{OUT} = \overline{\overline{B + A}}$$

Suppose  $A=B=0$  at the beginning. PD2 is a short circuit to ground, so the gate of PMOS1 is connected to ground. It means PMOS1 is ON, and because PD1 is an open circuit, the gate of PMOS2 sees 1, and so it is OFF.



BECAUSE PMOS2 STARTS TO GO ON, OUT STARTS TO BECOME 1, SWITCHING OFF PMOS1... AT THE BEGINNING IT IS SLOW, BUT THE SPEED INCREASES OVER TIME

IT BEHAVE LIKE A POSITIVE REACTION!

ADVANTAGES:

1. POSITIVE REACTION SPEEDS-UP THE SWITCHING
2. SHORT CIRCUIT POWER (VDD-GND DIRECT CONNECTION DURING TRANSITION) IS REDUCED
3. AS IN NMOS-LIKE, GATE CAPACITANCE OF THE LOAD IS REDUCED; EACH OUT DRIVES ONLY 1 NMOS W.R.T CMOS -> 1 NMOS, 1 PMOS.
4. USED TO INCREASE S.N.R (RESILIENCE TO NOISE). IN FACT IT IS NOT USED A STANDALONE SIGNAL, BUT THE DIFFERENCE BETWEEN TWO SIGNALS.
5. INTRINSICALLY GENERATES OUT AND OUT-bar, NO INVERTER NEEDED

DISADVANTAGES:

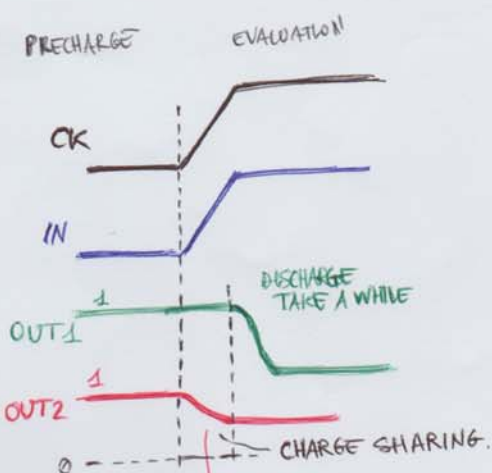
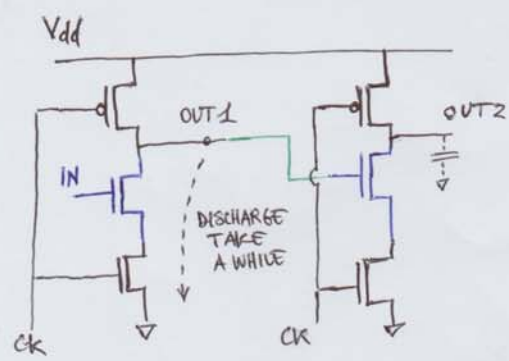
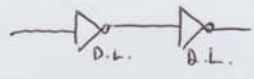
1. AREA: 2 PULLDOWN NETWORKS + 2 PMOS ARE USED (HOWEVER SAME NUMBER OF TRANSISTOR OF A CMOS IMPLEMENTATION + INV)

2. OPTIMAL SIZING IS A CRITICAL DECISION: CAPACITANCE OVER POSITIVE REACTION NETWORK SHOULD BE THE SAME

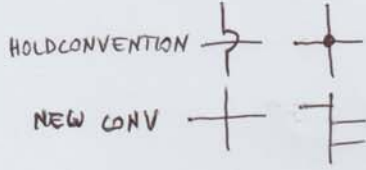


8/5/2018 PT1  
REG REG-  
HELECTR

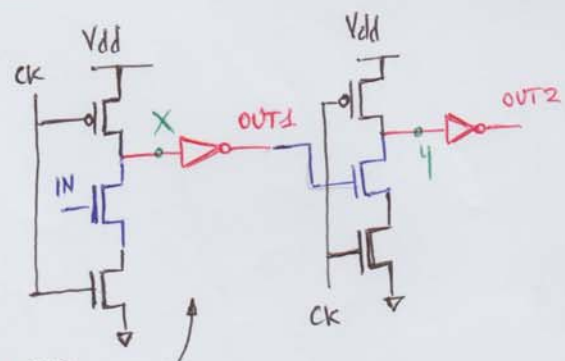
### CHARGE SHARING (IN DOMINO LOGIC)



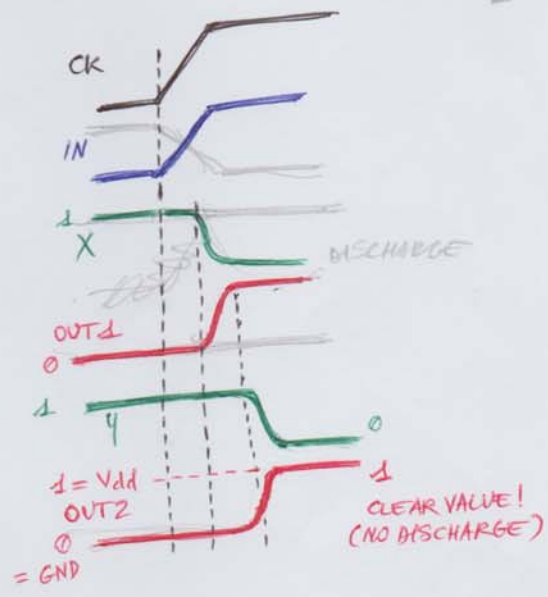
BECAUSE FOR A WHILE OUT1 REMAINS TO 1, BOTH NMOS NETWORK AND PMOS ARE ACTIVE AND THERE IS A PARTIAL DISCHARGE



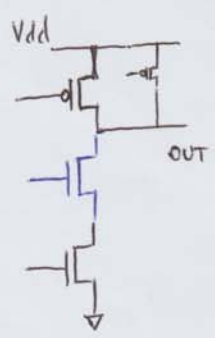
### SOLUTION



DOMINO GATE IS NOT-INVERTING (if inversion is needed -> CMOS INVERTER ADDED)



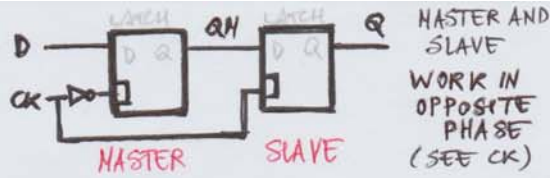
### SOLUTION 2: KEEP PER



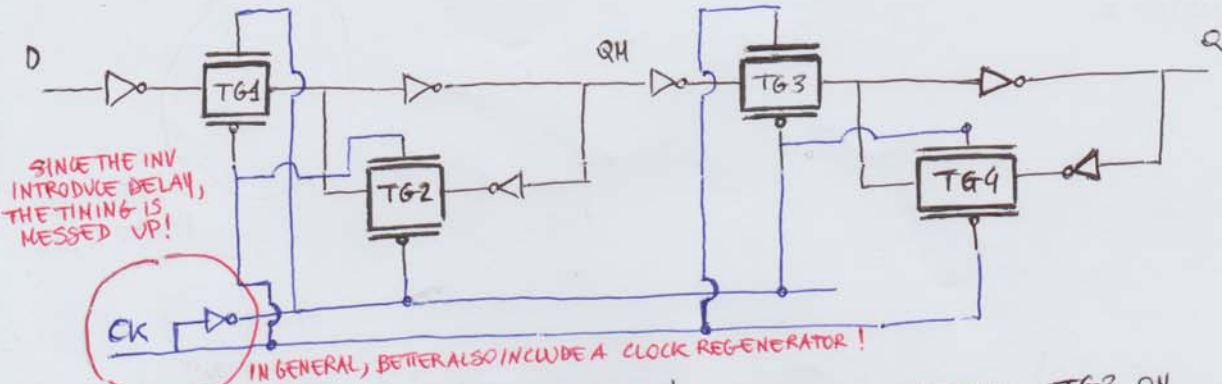
(1)



# FLIP-FLOP



8/05/18 PT2  
REG-REG-  
ELECTR.

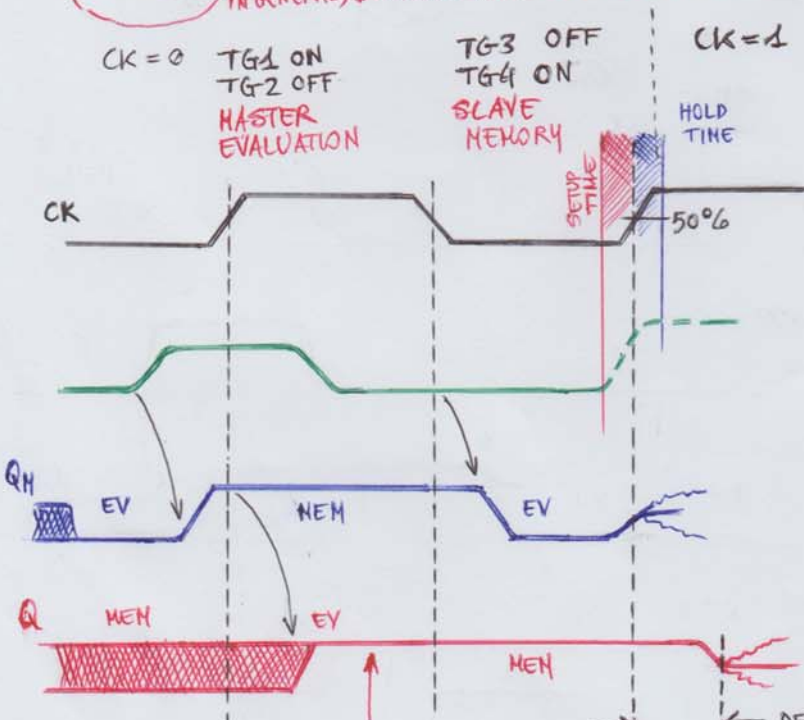


CK = 0  
TG1 ON  
TG2 OFF  
MASTER EVALUATION

TG3 OFF  
TG4 ON  
SLAVE MEMORY

CK = 1  
TG1 OFF  
TG2 ON  
MASTER MEMORY

TG3 ON  
TG4 OFF  
SLAVE EVALUATION



possible SETUP TIME violation of MASTER: for sure the problem is transfered to the slave (during slave evaluation)

**METASTABILITY!** It depends on many parameters

Q; EVEN IN EVALUATION, WILL REMAIN STABLE (AFTER THE EVALUATION) => NO RISK OF SETUP TIME VIOLATION

DELAY  $t_{CK \rightarrow Q} = t_{50\%Q} - t_{50\%CK}$   
IT DEPENDS ON

- SWITCH DELAY OF CK
- SWITCH DELAY OF 2 INV
- SWITCH DELAY OF TG3

**EDGE TRIGGERED**: Q CHANGE AFTER THE EDGE OF CLOCK (transition of CK) BECAUSE OF DELAYS!

**MITH BUSTING**: a flip flop isn't a magic box that sample at the clock edge and then is totally blind. Instead, inside it, a lot of things happen, without any stop!  
=> POWER CONSUMPTION  
=> NOISE

## CLOCK REGENERATOR:



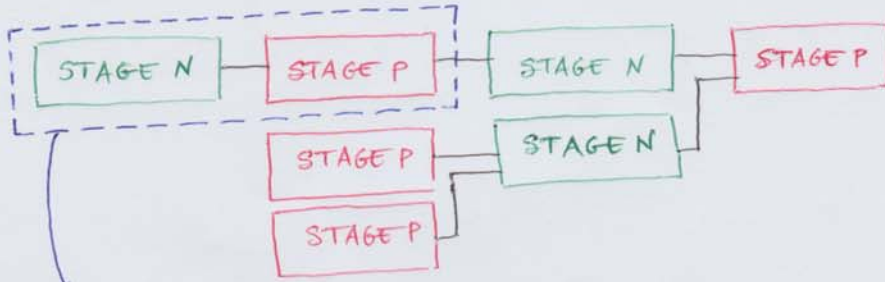
The system given clock should pilot a lot of logic inside our system. Should also travel long distances to reach our flipflop. A typical clock is a sine wave (V). This circuit introduce a little latency, but increase the clock slope!



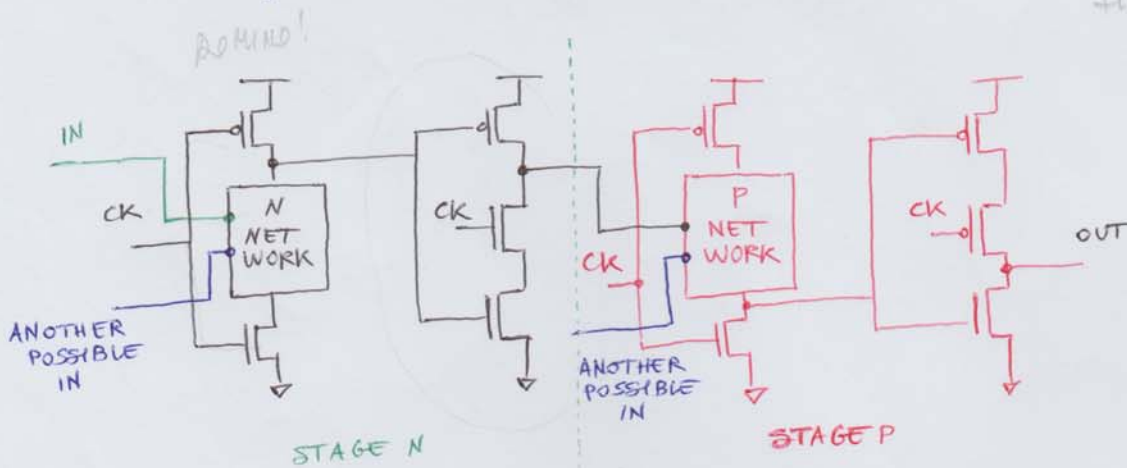
**DYNAMIC SEQUENTIAL LOGIC:**

TRUE SINGLE PHASE CLOCKING  $\Rightarrow$  TSPC

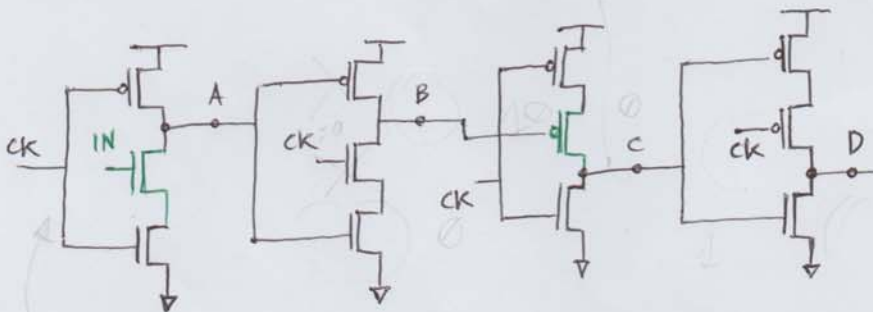
9/05/2018  
REG REG-



sim a MIN 6



Example: LOGIC FUNCTION 1 TRANSISTOR  $\rightarrow$  INVERTER



CK=0 PRECHARGE  
CK=1 EVALUATION

**MEMORY EVALUATION**

in precharge  $A=1$   
so PMOS is OFF.  
the NMOS connected to CK is OFF. It means neither the pull down or pull up is driving B, that maintain the old value!

**EVALUATION**  
PRE DISCHARGE

def the clock on the value of IN I can activate the pull up or maintain both disabled, maintaining zero

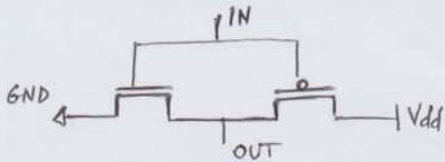
**EVALUATION**  
**MEMORY**

MIN 21-29

This logic allow to embed a logic function inside a latch. It's like computing on a double data rate, because we perform EXE+MEM in the same clock cycle.

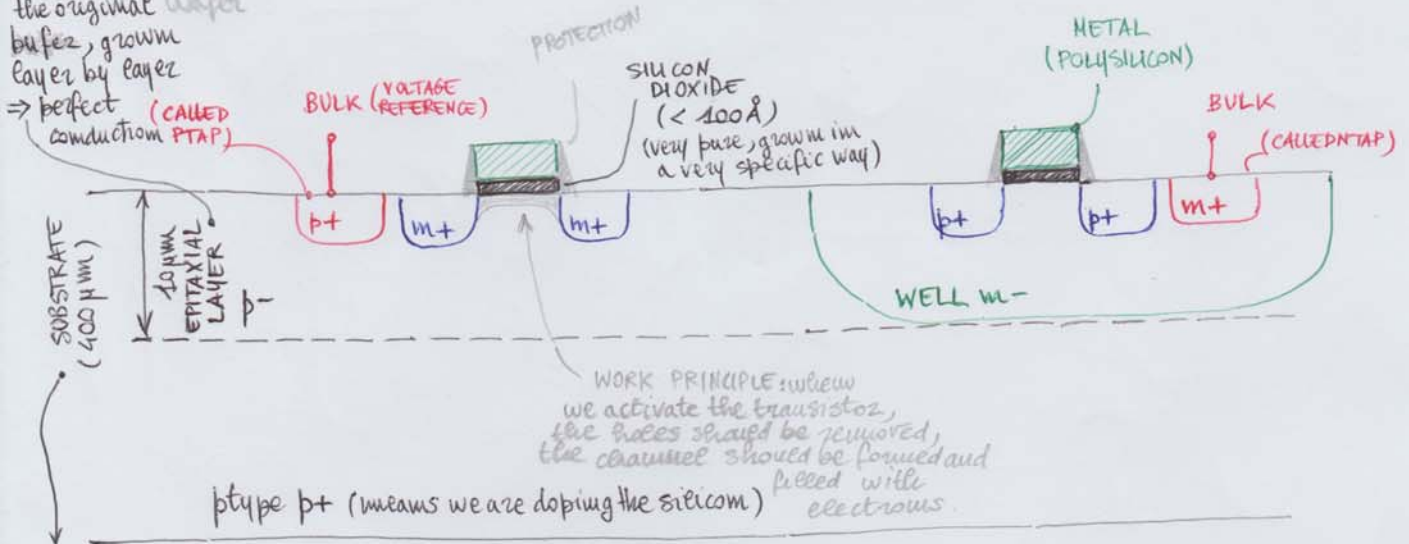


FROM SWITCH LEVEL TO TRANSISTOR LEVEL (TOWARDS GATE LIBRARY DESIGN & CHARACTERIZATION)



POLYSILICON: it is a silicon based conductor - Because we use silicon both for metal, oxide and semiconductor, the mobility, conduction and in general the compatibility is higher! (interface)

not part of the original wafer  
buffer, grown layer by layer  
=> perfect conduction (CALLED PTAP)



In order to get measurement we have to refer to a specific gate library => Below μm structure organization is different!

Why silicon is so widely used in the electronic industry, so that often is called "silicon industry"?

1. Easy to find on planet earth: for sure is true, and for sure it helped at the beginning of the electronic study and development, but is not the main reason.

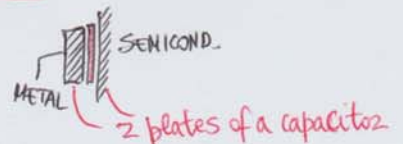
2. Silicon dioxide: in mos (metal oxide semiconductor) is really important the oxide layer.

Since silicon dioxide is a natural oxide of silicon: by exposing silicon to a certain oxygen (at a certain temperature and pressure) it will form naturally a even layer. => MASKS

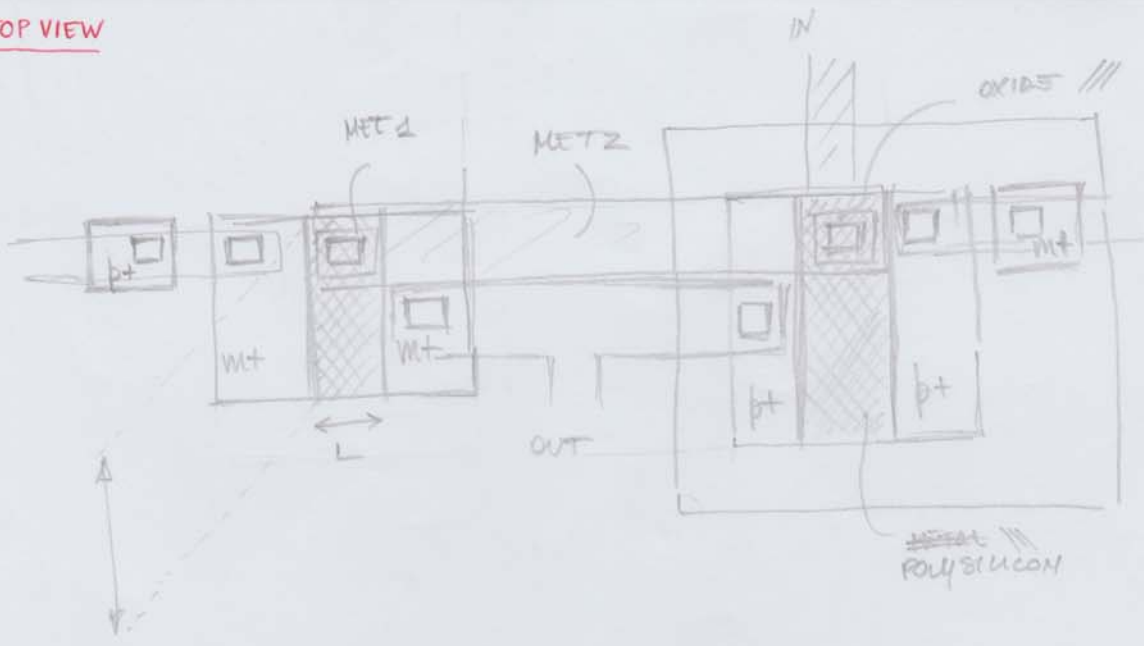
The layer between metal and semiconductor should be placed / builded in a different way, because it should be very defined. REMEMBER: between metal, oxide and semiconductor the channel should form if the transistor is active.

the thinner the oxide layer is, the better is for the overall transistor quality / behaviour:

1. Less gate capacitance
2. Threshold voltage
3. ...
4. ...

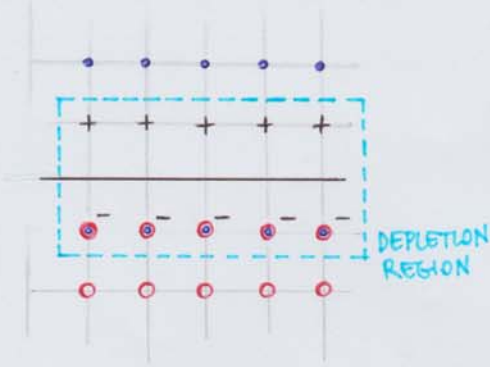


TOP VIEW



MASK  
VIEW

TRANSISTOR PMOS  
ARE BIGGER BECAUSE  
MOBILITY IS LOWER  
AND I WANT THE  
SAME DELAY OF NMOS

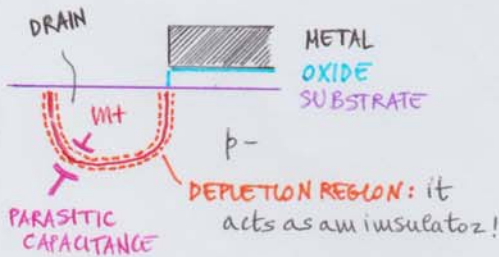


After the jump, the border of n-doped side is characterized by a negative charge, and the border of p-doped side is characterized by a positive charge.

This cause two things:

1. Electric field (called built-in electric field) is generated
2. A voltage difference is formed all across the border: a wall that is hard to pass by electrons. (it's around  $0.6V = V_g$  of a diode)

NIN 17-19



Another parasitic capacitance! Semiconductor (p-), insulator (depletion region), semiconductor (n+)  
= parallel plate capacitor

JUNCTION CAPACITANCE:

$$C_j = \frac{C_{j0}}{\sqrt{1 - \frac{V_D}{\phi_0}}} = A_D f(N_D, N_A, \dots)$$

OTHER PARAMETERS DEPENDS ON TECHNOLOGY

$= W \cdot X$   
DESIGN PARAMETER

BUILT-IN JUNCTION POTENTIAL

$$\phi_0 = \phi_T \ln\left(\frac{N_A N_D}{n_i^2}\right)$$

CONCENTRATION OF INTRINSIC SILICON

$V_p$  of diode that depends on  $V_T$  (threshold voltage) + current

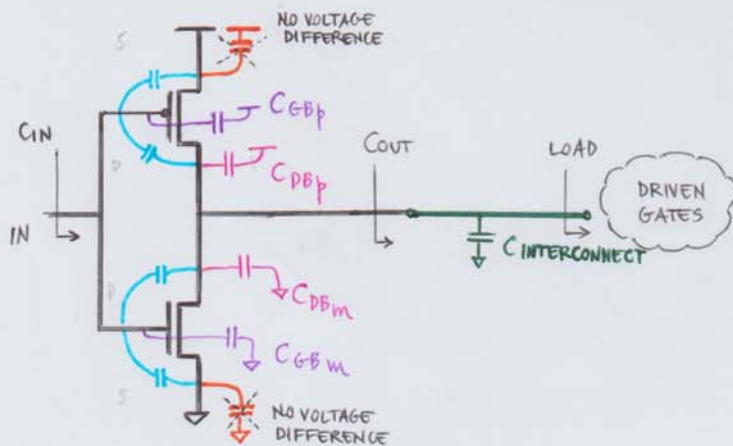
Recap example: inverters

pmos: BULK connected to Vdd

nmos: BULK connected to GND

$$i_D = I_S e^{V_D / mV_T}$$

$$V_D = mV_T \ln\left(\frac{i_D}{I_S}\right)$$



$$C_{IN} = C_{GBM} + C_{GBP} + C_{GDM} + C_{GDP} + C_{GSM} + C_{GSP}$$

$$C_{OUT} = C_{LOAD} + C_{INTERCONNECT} = C_{LOAD} + C_{DBM} + C_{DBP}$$

FEED THRU USH!

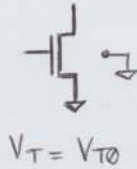


**BODY EFFECT**

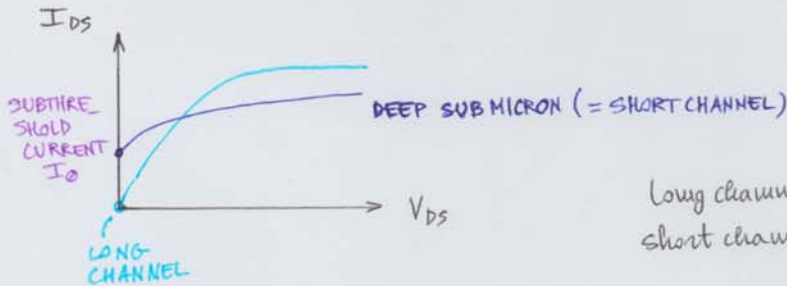
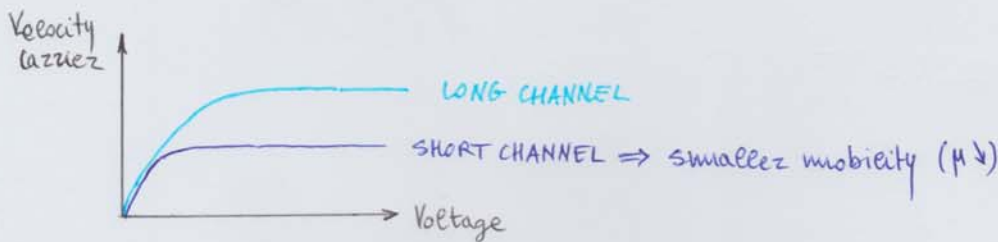
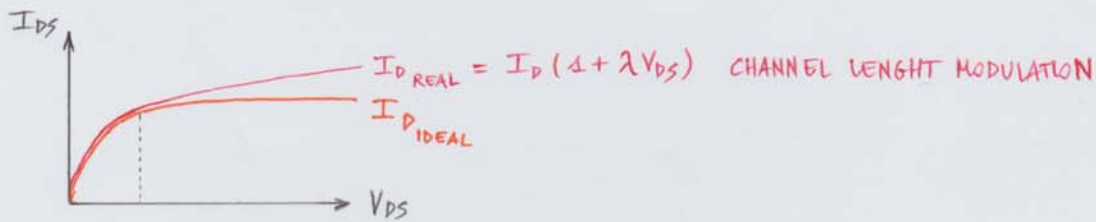
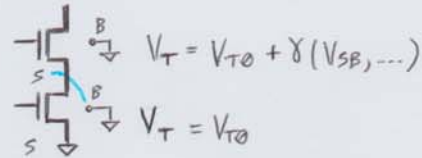
Even if we are able to build the exact same transistor multiple times (and it is not possible), every transistor can be characterized by a different threshold voltage  $V_T$ , resulting in a more difficult way to drive correctly the gate.

Example:

single transistor:



simple gate:



Long channel:  $V_{DS} = 0 \Rightarrow I_{DS} = 0$  no leakage  
 Short channel:  $V_{DS} = 0 \Rightarrow I_{DS} = I_0$  leakage!

HIGH  $V_T$  transistors have low  $I_0$ , but transition is slower (higher time to reach a bigger  $V_T$ )

DEPENDS ON TEMPERATURE QUADRATICALLY  
 DEPENDS ON  $V_T$

(A)  $V_{in} = 0$

Supposing a past history of transitions, in this condition  $V_{out} = V_{dd} = 1$   
(already powered on, not on first switch ...)

NMOS OFF  $\left\{ \begin{array}{l} V_{GS_m} = V_{in} - GND = V_{in} - 0 = 0 \\ V_{DS_m} = V_{out} - GND = V_{out} - 0 = V_{dd} \end{array} \right.$  NMOS OFF (in OFF region) because of NO CHANNEL  
 VERY HIGH: if  $V_{DS_m}$  is high, potentially we can move a lot of charges ... however, because  $V_{GS_m} = 0$  there are no charges (because there is no channel ...)

PMOS OFF  $\left\{ \begin{array}{l} V_{GS_p} = V_{in} - V_{dd} = 0 - V_{dd} = -V_{dd} \\ V_{DS_p} = V_{out} - V_{dd} = 0 \end{array} \right.$  VERY HIGH, BUT NEGATIVE: there are a lot of holes in the channel  
 THERE IS CHANNEL, BUT NOT ELECTRIC FIELD  
 there is not electric field in the channel: the charges are in the channel but are not moving  $\Rightarrow$  POFF



(B)  $V_{in} \uparrow$  ( $V_{in}$  is increasing)

$V_{GS_m} \uparrow$

$|V_{GS_p}| \downarrow$  but still is NEGATIVE

Are we going above the threshold or not? There is an empiric rule (however it depends on technology):  $V_T \approx V_{dd}/10 \Rightarrow$  the threshold voltage value is around  $\frac{1}{10} V_{dd}$  ...

Is very likely, because  $V_T$  is so small, that  $V_{GS_m}$  is around the threshold voltage  $\Rightarrow V_{GS_m} \approx V_T$   
 $V_{GS_p}$  remains above the threshold voltage  $\Rightarrow |V_{GS_p}| \gg V_T$

We are not putting numbers: however for sure the PMOS has (still) the channel formed and suppose that those conditions on  $V_{in}$  value are such that also NMOS is ON (formed a channel)

$\Rightarrow$  BOTH NMOS and PMOS HAVE CHANNEL FORMED  
 HAS FEW CHARGES ( $e^-$ ) BELOW THE CHANNEL  
 HAS A LOT OF CHARGES (HOLES) BELOW THE CHANNEL

$V_{DS_m} \downarrow$  remains high, remains a strong electric field across the channel  $\Rightarrow$  THERE IS A SMALL CHANNEL + BIG ELECTRIC FIELD  $\Rightarrow$  NMOS IS ON!

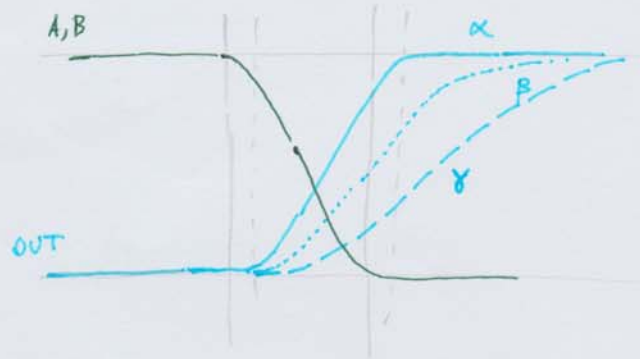
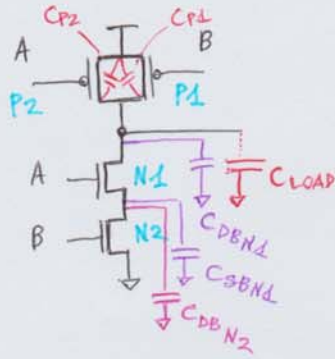
$V_{out} \downarrow$

By looking at the NMOS, probably it is in a SATURATION REGION because  $V_{DS_m} > V_{GS_m} - V_T$  ( $\approx 1 \gg 0$ )  
 PMOS SATURATION (?)  $V_{DS_p} > V_{GS_p} - V_T$  ( $0 \gg \text{NEG}$ )

$|V_{DS_p}|$  starting to increase, however it assume a NEGATIVE VALUE  $\Rightarrow$  a small electric field is formed  $\Rightarrow$  PMOS transistor is ON!



DEPENDENCY ON INPUT PATTERN



SITUATION  $\alpha$ :

$A_1 \rightarrow 0$

$B_1 \rightarrow 0$

$N_1, N_2$  ON  $\rightarrow$  OFF

$P_1, P_2$  OFF  $\rightarrow$  ON

Starting situation: for size

$V_{out} = 0$  and  $C_{load} + C_{dbn1} + C_{sbn1} + C_{dbn2}$  are all discharged completely

OUT  $\rightarrow 1$  supposing NMOS switch off immediately:

$P_1 + P_2$  CHARGING  $\left\{ \begin{array}{l} C_{load} \text{ CHARGED} \\ C_{dbn1} \text{ CHARGED} \\ C_{p1} \text{ CHARGED} \\ C_{p2} \text{ CHARGED} \end{array} \right.$

REASON WHY WHEN DON'T PUT TOO MANY INPUTS IN A SINGLE GATE = BETTER TO SPLIT !..

TRY A NOR

SITUATION  $\beta$ :

$A_1 \rightarrow 0$

$B_1 \rightarrow 1$  (remain at 1)

$N_1$  ON  $\rightarrow$  OFF

$N_2$  ON

$P_2$  OFF  $\rightarrow$  ON

$P_1$  OFF

OUT  $\rightarrow 1$

$P_2$  CHARGING

$\left\{ \begin{array}{l} C_{load} \\ C_{dbn1} \\ C_{p1}, C_{p2} \end{array} \right.$

even if one of the p is off. Only one p is doing the job  $\Rightarrow$  slower out: bigger propagation time increasing rise time

SITUATION  $\gamma$ : WORST CASE

$A_1 \rightarrow 1$

$B_1 \rightarrow 0$

$N_1$  ON

$N_2$  ON  $\rightarrow$  OFF

$P_2$  OFF

$P_1$  OFF  $\rightarrow$  ON

OUT  $\rightarrow 1$

$P_1$  CHARGING

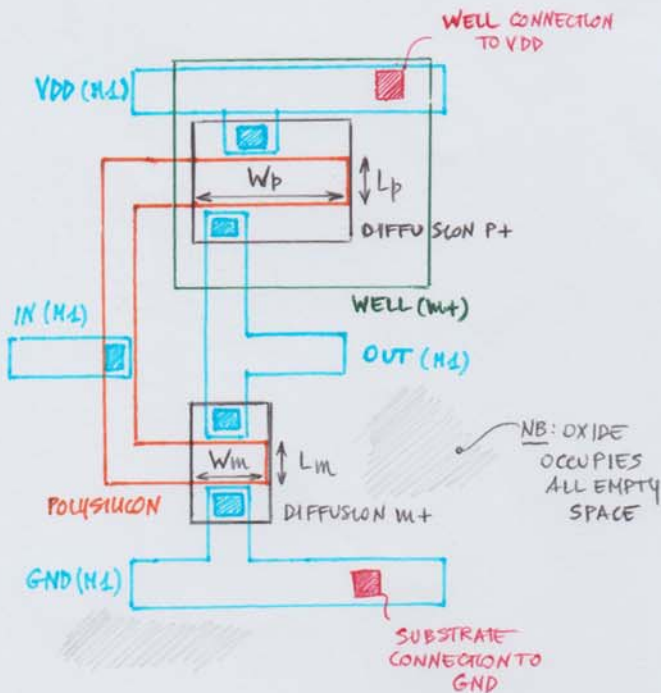
$\left\{ \begin{array}{l} C_{load} \\ C_{dbn1} \\ C_{p1}, C_{p2} \end{array} \right.$

$\left\{ \begin{array}{l} C_{sbn1} \\ C_{dbn2} \end{array} \right. \rightarrow$  Because  $N_1$  is ON!



③ LAYOUT OF EACH CELL

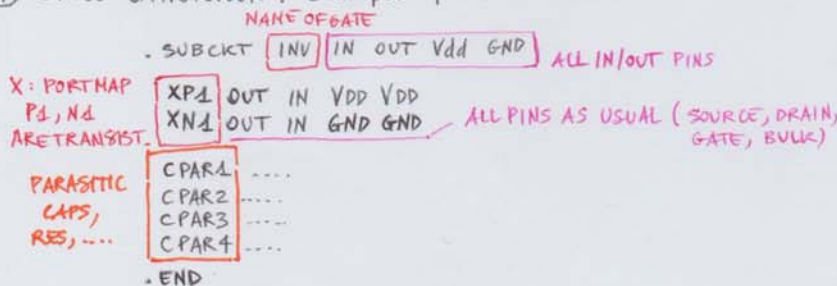
- ALL HEIGHTS H IDENTICALS
- WIDTH CHANGE DEPENDING ON CELL COMPLEXITY & TRANSISTORS SIZING



LAYOUT CHECKED VS DESIGN RULES  
SCHEMATIC (the figure represent an inverter!)

No matter how many simultaneous you run, you'll have always a finite number of values. Suppose that you z load is not exactly reported in the L.U.T.: just do an average! ⇒ FIND WHAT IS CALLED AN EFFECTIVE VALUE! (INTERPOLATION)  
Example:  $C_{L2} < C_{EFF} < C_{L3}$ ,  $t_{R2} < t_{RFF} < t_{R3}$   
is highlighted in the LUT down below...

④ SPICE EXTRACTION: example of an inverter



⑤ CHARACTERIZATION

Analysis of all delays, behaviours ... by simulating for a cot of different values of IN, etc...

⇒ SPICE SIMULATIONS FOR ALL CASES

⇒ MEASURE:  $t_{PROP, LH}$  (when OUT )

PROPAGATION TIME =  $t_{50\% OUT - t_{50\% IN}}$  (when OUT )

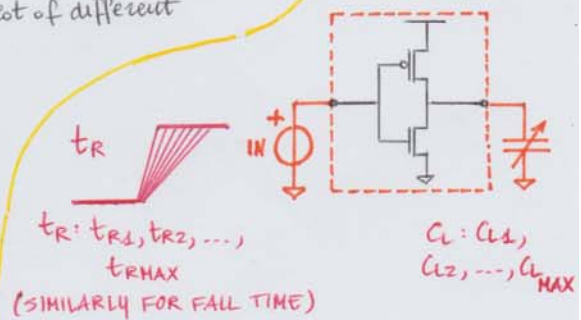
RISE AND FALL TIME MEASURED BETWEEN 10% - 90%  
 $t_{RISE}$  (10% to 90%)  
 $t_{FALL}$  (90% to 10%)  
+ ALL CURRENTS! (SWITCHING, SHORT CIRCUIT, LEAKAGE)

⇒ FOR COMPLEX GATES (more than one input) SIMULATIONS ARE REPEATED FOR ALL POSSIBLE INPUTS SWITCHING CASES!

⇒ EACH IN-OUT COMBINATION IS STORED IN A L.U.T.

⇒ ONE L.U.T FOR EACH MEASURE!!  
MANY L.U.T.S AT THE END

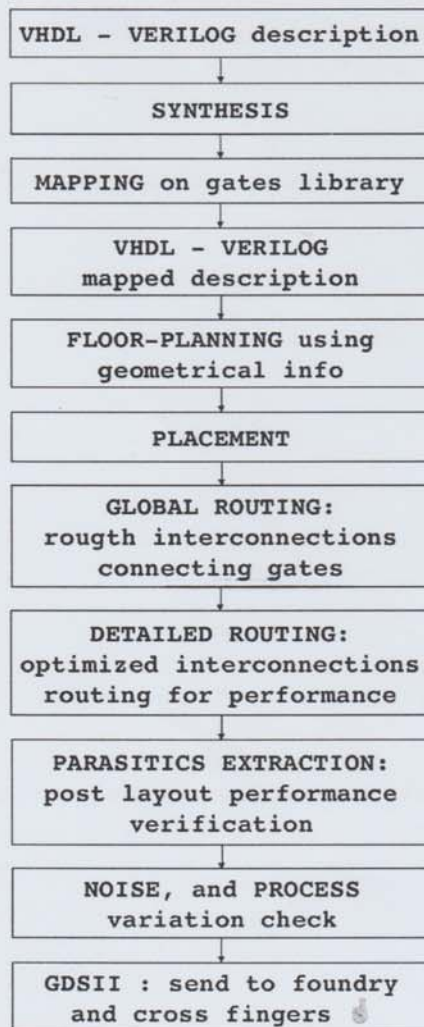
⇒ DESCRIBED IN LIBERTY FILES (.lib) AND USED FOR FURTHER SYNTHESIS PHASES



	CL1	CL2	CL3	CL4	...
tr1	tp11	tp21	tp31	tp41	...
tr2	tp12	tp22	tp32	tp42	...
tr3	tp13	tp23	tp33	tp43	...
...	...	...	...	tp44	...

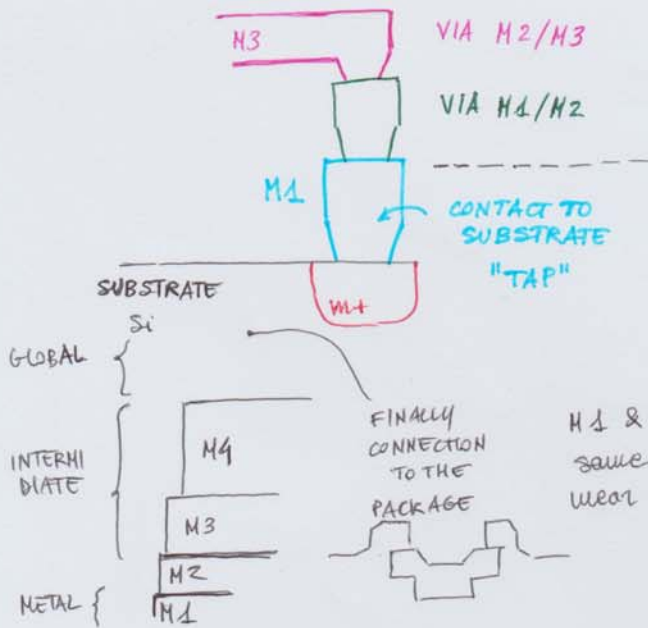
## ASIC design flow (top-bottom)

1. Based on a library of gates: a library contains all kinds of different gates (nand, and, or, inv, buffer, xor, mux, flip-flop, latches...), both elementary blocks and more complex ones.
2. Each gate is designed to drive multiple loads (X4, X8, ...). The driving capability is represent as the maximum number of elementary blocks the gate can drive without lose control of delay. Usually an elementary block correspond to an inverter (look documentation)
3. Each gate has different implementations, each one optimized for speed (HS) and/or for power (LL).
4. Each gate is described at
  - Transistor level (spice)
  - Layout level (LEF)
  - Timing-power-boolean (lib)
  - Behavioral (VHDL): VHDL is used as input of the synthesis part. However it can be used also after the synthesis analysis or even after the layout => the estimated delays can be used to run a realistic simulation (simulink). Maybe you need two change layout or architecture entirely because the circuit does not respect constraints.



29 MAY  
PT1

INTERCONNECTIONS



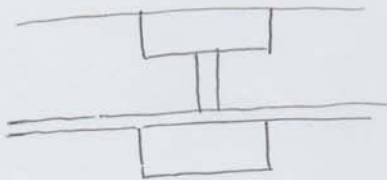
FIRST LAYER COULD BE DONE WITH A DIFFERENT TYPE OF METAL (TUNGSTEN) BECAUSE WE NEED A MECHANICAL STABILITY

M1 & M2 are used to connect points on the same cell, ... M3 ... to connect points of wear cells, ... MX ... to connect macroblocks

HIGHER LAYER = MORE METAL (HIGHER THICKNESS) BECAUSE OF ~~SMALLER~~ <sup>BIGGER</sup> SMALLER CURRENT DENSITY

GLOBAL DELIVER A VERY LARGE AMOUNT OF CURRENT

example: ST-ELECTRONICS

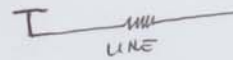


$k = 2.55 = \epsilon_r$

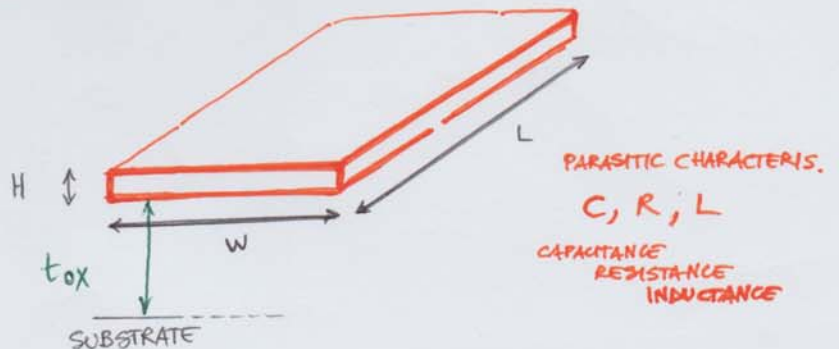
LOW DIELECTRIC CONSTANT TO REDUCE CAPACITANCE; HOWEVER THEY TEND TO BE WEAR  $\Rightarrow$  USE THIN MATERIAL IN THE MIDDLE (IN ORDER TO NOT CONTRIBUTE TOO MUCH ON CAP)

1. Heating

2. Voltage drop (especially on power supply lines)  $\equiv$  Noise



3. Reliability: protection from overcurrent. In a lot of systems there are fuse, that can deliver only a certain amount of current, otherwise it will brake... ELECTROMIGRATION





RESISTANCE DEPENDS ON FREQ. AND TEMP.

1. FREQUENCY R(f)

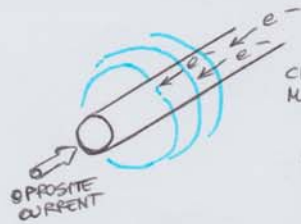


- SECTION OF A CONDUCTOR
- CURRENT IS EXPECTED TO DISTRIBUTE UNIFORMLY IN THE WHOLE SECTION

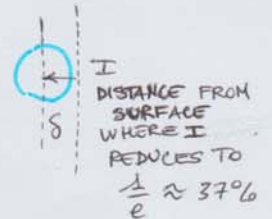


- IT MIGHT HAPPEN THAT CURRENT CONCENTRATE MORE IN THE PERIPHERY.

⇒ SKIN EFFECT



SKIN DEPTH:  $\delta = \sqrt{\frac{\rho}{\pi \mu f}}$   
 $\mu \parallel \nu$  PERMEABILITY



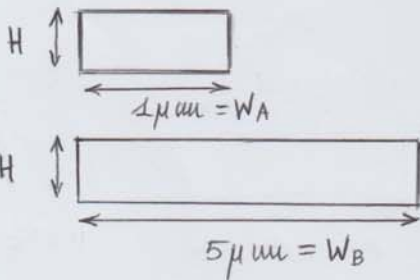
How FREQUENCY IS DEFINED FOR SQUARE WAVES?

ex: CLK



$T \rightarrow \frac{1}{T} = f + \text{other components}$

let's use only main frequency component:



f	δ (SKIN DEPTH)
100 MHz	11 μm = δ <sub>1</sub>
500 MHz	5.19 μm = δ <sub>2</sub>
1 GHz	3.67 μm = δ <sub>3</sub>
2 GHz	2.6 μm = δ <sub>4</sub>

$\delta_1, \delta_2, \delta_3, \delta_4 > W_A \Rightarrow$  NO SKIN EFFECT ON CONDUCTOR A

ESPECIALLY ON CASE 4:  $\frac{W_B}{2} \approx \delta_4$  SKIN EFFECT IS PRESENT!

skin effect has an impact: effective resistance increasing  
 it is like conductor B has a smaller section: smaller current flow  $\Rightarrow$  INCREASE RESISTANCE

TECHNOLOGY DEFINE A MAXIMUM WIDTH, IN ORDER TO AVOID THE SKIN EFFECT



$$I_{RMS} = \sqrt{\frac{1}{T} \int_0^T I^2(t) dt}$$

$f \uparrow \quad I_{RMS} \uparrow \quad P_m \uparrow \quad T_m \uparrow \quad R \uparrow$   
 $W \uparrow \quad \theta \downarrow \quad T_m \downarrow \quad R \downarrow$   
 $tox \uparrow \quad \theta \uparrow \quad T_m \uparrow \quad R \uparrow$  but HIGHER H-LAYER  $H \uparrow \quad R \downarrow$   
 HIGHER H-LAYER

HIN 34-37 UNTIL E.M.

**ELECTRO MIGRATION**

If electrons in a conductor have sufficient energy VOIDS (open circuits) or HILLOKS (short circuits) could be generated.



Electronic wind can accelerate electrons that can transfer some quantity of motion to metal atoms. If it happens continuously several atoms could move as an avalanche.

More probable if:

1. many electrons are present (METAL, POLY SILICON)
2. high current density ( $J = I/WH$ )
3. electrons free to move
  - BOUNDARIES
  - DEFECTS
  - GRAIN BOUNDARIES



**GRAIN BOUNDARIES:** free space. Here the atoms hit by electrons can be placed.

**GRAIN:** each grain has a specific orientation of the crystal.

BIG W reduce resistance but increases the risk of electro migration.

ACTIVATION ENERGY depends on the type of metal lattice

**BLACK'S LAW (empirical):**

$$H.T.F. \propto \frac{1}{J^2} e^{-\frac{Q}{kT_m}}$$

MEAN TIME TO FAILURE (circled)  
 CURRENT DENSITY (J)  
 ACTUAL METAL TEMPERATURE (T<sub>m</sub>)  
 BOLTZMAN CONSTANT (k)  
 ACTIVATION ENERGY (Q)

BETTER ROUTE POWER LINES ON HIGHER H. LAYERS IN ORDER TO REDUCE J (H ↑ H.T.F. ↑)

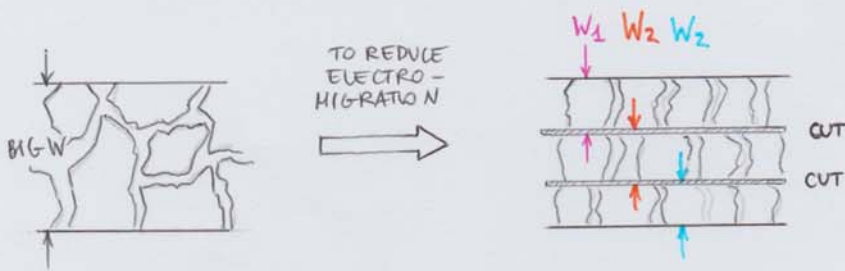
H.T.F. is usually measured in years (consumer electronics has 10 years of H.T.F.).  
 H.T.F. it is exponential dependent on the temperature!

30 MAY

RECAP: grain boundaries are a characteristic of polycrystalline materials, like the metal used for interconnections. They are defects in the crystal structure, and tend to decrease both electrical and thermal conductivity of the material. For our analysis, it is important to remember that atoms of metal hit by electrons (current = flow/motion of electrons) ~~can~~ can move into the grain boundaries.

By reducing the order of magnitude of  $W$  (width) of the wire, grain boundaries has a perpendicular ( $\perp$ ) direction with respect to the electrons directions.

This physical layout reduce the probability that an atom of metal, hiten by electrons, falls into the grain boundaries.



Several wts creates the equivalent of several small lines in parallel.  
 $W_1 + W_2 + W_3 \approx W$   
 Same current driving capabilities.  
 $W_1, W_2, W_3 \approx$  SIZE OF GRAIN

When a lot of current has to flow, a big wire is required: BIG CURRENT = BIG WIRE

**CHEESED INTERCONNECT.**  
 Especially used in power and ground lines

Bamboo effect is favoured  $\Rightarrow$  Electro-migration risk reduced.

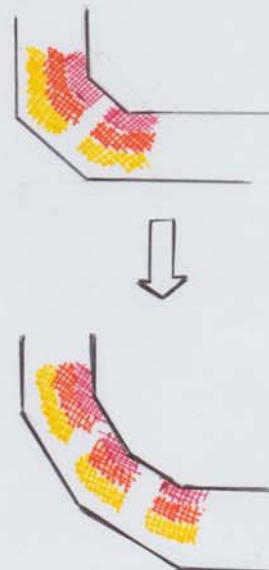
CORNERS:

Usually in interconnection layouts, lines are straight, or characterized by  $90^\circ$  corners. However:

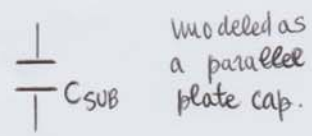
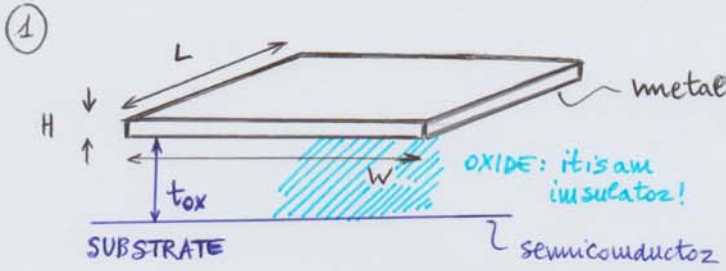


Explanation: high speed electrons choose always the minimum resistance path. This means higher temperature on the corner inside  $\Rightarrow$  Increase the risk of electro migration

SOLUTION: FROM MANHATTAN STYLE TO A MORE BROKEN-LINE LAYOUT (ROUNDED)



**CAPACITANCE** 2 main contributions → // PLATE CAP.  
 ↳ FRINGE CAP.



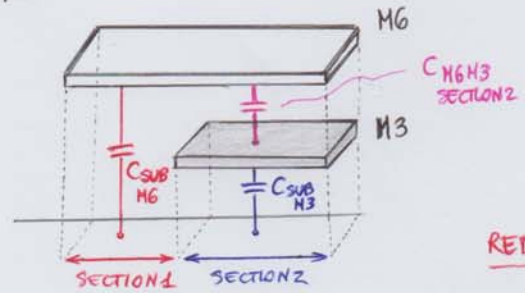
$$C_{SUB} = \epsilon_0 \epsilon_r \frac{WL}{t_{ox}}$$

SILICON DIOXIDE (SiO<sub>2</sub>)  
 DIELECTRIC CONSTANT = 3.9

W ↑ C ↑ R ↓  
 HIGH LEVEL H.L. tox ↑ C ↓ Tm ↑ R ↑

capacitance opposite w.r.t. resistance ☹️

Example:

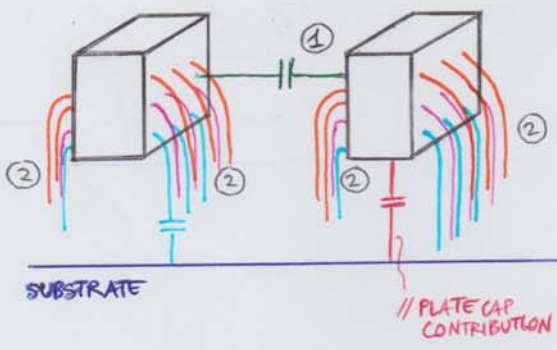


supposing W<sub>M3</sub> = W<sub>M6</sub> ...

$$C_{M6M3} = \epsilon_0 \epsilon_r \frac{W_{M3} L_{SECTION2}}{t_{ox M6-M3}}$$

**REDUCE CAP:** a way is to reduce  $\epsilon_r$  by using a different material, instead of SiO<sub>2</sub>. There are some organic materials with lower  $\epsilon_r$ , but they are more flexible ... interleave them with some thin stronger material, that does not introduce too much cap.  $\epsilon_r \downarrow K_{ox} \uparrow$  TERNAL CONDUCTIVITY

**(2) FRINGE CAP**

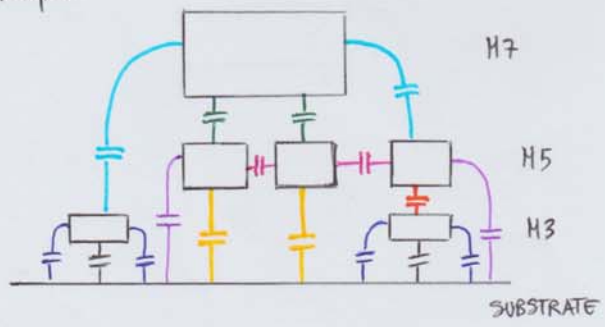


- Due to electric field lateral component:
1. contribution from the side of the wire (conductor) to the side of nearby conductor (wire, usually)
  2. contribution from the side of the wire to the substrate.

This happens a lot in buses, because they are long parallel wires (# depends on parallelism)

FAT WIRES = BIG HEIGHT = HIGH H.L.

Example:

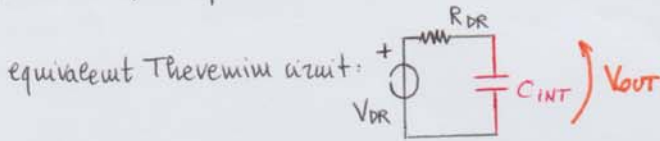
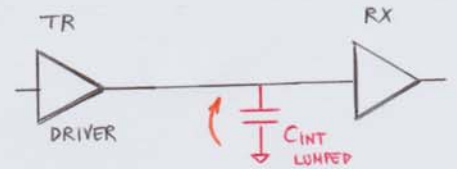




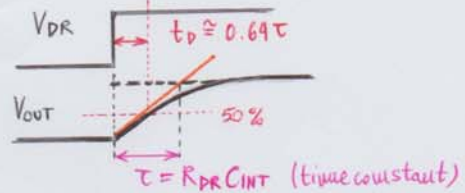
5/09/18

① LUMPED-RC

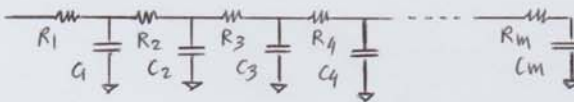
The capacitance is distributed across the line; however the model consider it as a concentrated capacitance. It neglects completely the resistance of the line



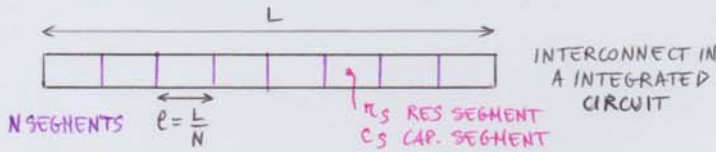
$$V_{OUT} = V_{DR} (1 - e^{-t/\tau})$$



② ELMORE



We can use this model to describe interconnects:



50% DELAY

$$t_D \approx C_1 R_1 + C_2 (R_1 + R_2) + \dots + C_m (R_1 + R_2 + \dots + R_m) = R_1 (C_1 + C_2 + \dots + C_m) + R_2 (C_2 + \dots + C_m) + \dots + R_m C_m$$

Divide the interconnection in N segment of length l. Each one will be characterized by a certain resistance and capacitance:

$$r_s = r l = r L / N \quad c_s = c l = c L / N$$

RES per unit length      CAP per unit length

We can apply now the elmore model:

$$t_D \approx C_s r_s + C_s (r_s + z_s) + \dots + C_s (r_s + \dots + r_s)$$

$$= C_s r_s + C_s (2r_s) + \dots + C_s (N r_s)$$

$$= r_s C_s \frac{L^2}{N^2} (1 + 2 + 3 + \dots + N) = r_s C_s \frac{L^2}{N^2} \frac{N(N+1)}{2}$$

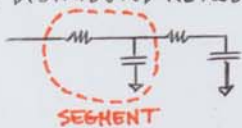
if N is big  $\Rightarrow N+1 \approx N$

$$= \frac{r_s C_s}{2} L^2 \quad \text{DEPENDENCY ON TECH PARAMETERS}$$



recall example SQUARE DEPENDENCY ON L

③ DISTRIBUTED RC MODEL



$l \rightarrow \Delta l$

$$r_s \rightarrow r \Delta l$$

$$c_s \rightarrow c \Delta l$$

solve differential equations:

$$t_D \approx 0.38 RC \quad \text{TOTAL CAP}$$

$L$  TOTAL RES

④ DISTRIBUTED RCL MODEL



$$l \rightarrow \Delta l \quad r_s \rightarrow r \Delta l$$

$$c_s \rightarrow c \Delta l$$

$$l_s \rightarrow l \Delta l$$



⑤ TRANSMISSION LINE

Up to MODEL 4:



Model 5:



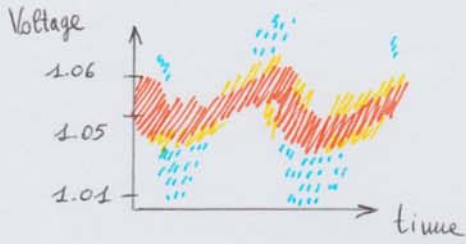
$$V_2 = V_1 \Gamma_{RX}$$

$$\Gamma_{RX} = \frac{Z_{RX} - Z_{LINE}}{Z_{RX} + Z_{LINE}}$$



EXAMPLE: measured supply noise from 90nm Itanium Microprocessor

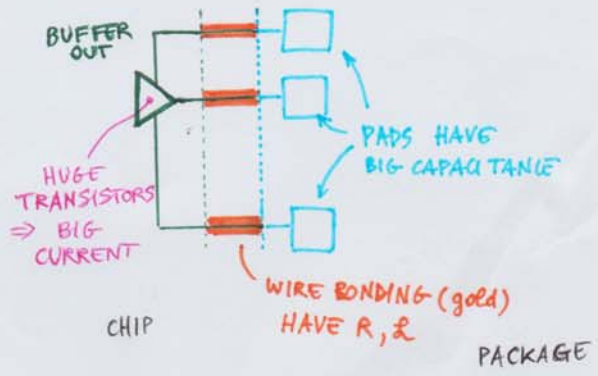
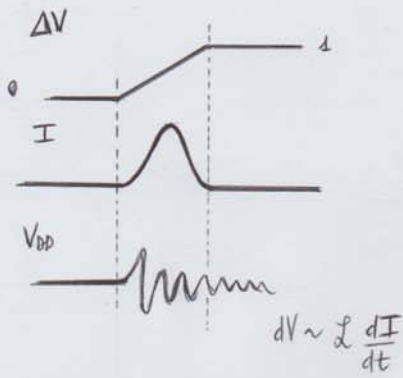
PT2  
@ 1.30  
(-1 REG ONLY)



The Itanium  $\mu$ processor works with a 1.1 Volts power supply ---  
However the average voltage is around 1.054 V and in general power supply has a sinusoidal behaviour / progression...  
There are also some sporadic variations on even lower voltage ---

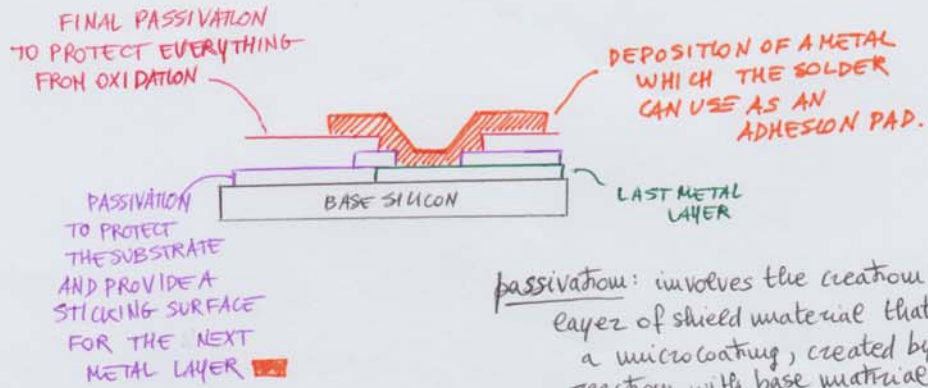
FINO A 1.37

SSN: SIMULTANEOUSLY SWITCHING NOISE



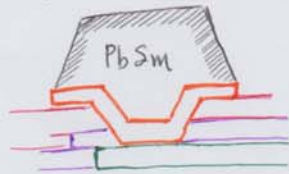
C4 process steps :

① Metallization of contact pads :



passivation: involves the creation of an outer layer of shield material that is applied as a microcoating, created by chemical reaction with base material, or allowed to build from spontaneous oxidation in the air.

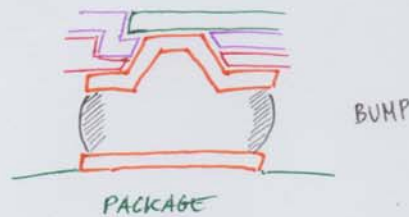
② solder deposition



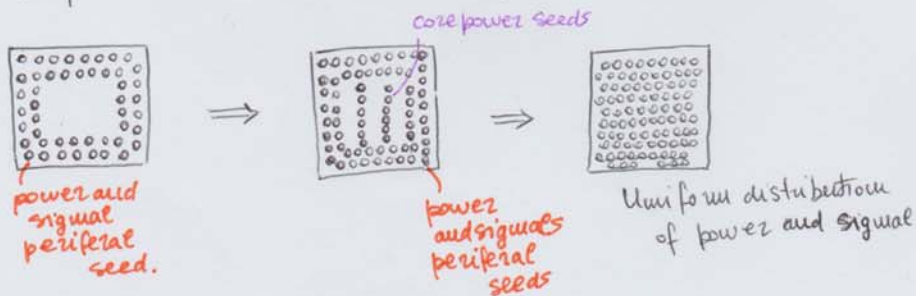
③ Reflow process: with heat, the solder melts into a ball



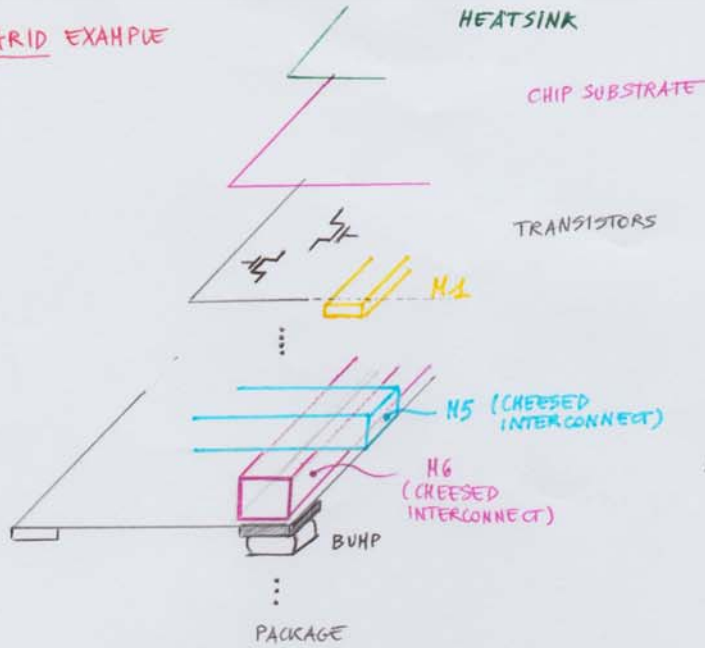
④ Bonding process: the chip is flipped, and connected to the package - (Alignment!!)



Historically:



POWER GRID EXAMPLE

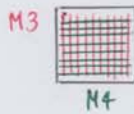


Usually the two highest metal layers have  $\perp$  directions and just carry Vdd and GND.

- 3 metal layer approach: at the beginning the metal 3 was used for power supply / clock routing. The organization was only one direction:

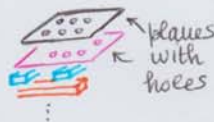


- 4 metal layer approach: in order to reach in a more precise way all parts of the circuit, a metal layer M4 is introduced,  $\perp$  with respect to M3. It also give more spots to put power supply pads. Still M3 and M4 mostly used for power supply / clock. Usually more than 50% of total metal is used for Vdd / Gnd.



POWER GRID

- 6 metal layer approach (planes). Instead of two Metal layers dedicated to Vdd and Gnd, two solid planes are used. This reduce a lot the resistance! Also, because there are less straight thin metal line, also chip inductance is reduced. This approach however is very expensive.



- Lowcost, low frequency applications: routed power distribution on two stacked layers of metal (one for VDD, one for GND).

