

Campus

A5

7mm×24行 40页

System on Chip Techniques

# Campus

A5 点线本 | 7mm×24行 | 40页

KOKUYO

Campus 无线装订本  
A5 40页

WCN-CNB3430



6 937748 309494  
MADE IN CHINA

国誉商业(上海)有限公司

<http://www.kokuyo.cn/st/>

上海市奉贤区人杰路128号

TEL : 400-820-0798 FAX : 021-3255-8508

产地:上海市 QB/T1438-2007合格

A5 210×148mm

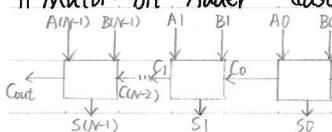
● 采用日本进口纸张,  
纸质细滑,牢固不掉页。

Half Adder: No Cin  
Full Adder: With Cin

### # Binary Adder

$$\begin{array}{ccccccc}
 \text{Cin} & A & B & S & C & S = \text{Cin}'A' \cdot B + \text{Cin}'A \cdot B' + \text{Cin} \cdot A' \cdot B' + \text{Cin} \cdot A \cdot B & \left. \begin{array}{l} 2X \cdot \text{XOR} \\ + \end{array} \right. \\
 0 & 0 & 0 & 0 & 0 & = \text{Cin}'(A'B + AB') + \text{Cin}(A'B' + AB) & \\
 0 & 0 & 1 & 0 & 1 & = \text{Cin}'(A'B + AB') + \text{Cin}(A'B + AB')' = A \oplus B \oplus \text{Cin} & \left. \begin{array}{l} 2X \cdot \text{OR} \\ + \end{array} \right. \\
 0 & 1 & 0 & 0 & 0 & \\
 0 & 1 & 0 & 1 & 1 & C = \text{Cin}'AB + \text{Cin}A + \text{Cin}B = \text{Cin}(A+B) + AB & \left. \begin{array}{l} 2X \cdot \text{AND} \\ + \end{array} \right. \\
 1 & 0 & 1 & 0 & 0 & [\text{If do not use logic minimisation}] & \\
 0 & 0 & 1 & 0 & 0 & \hookrightarrow C = \text{Cin}'AB + \text{Cin}A'B + \text{Cin}AB' + \text{Cin}AB & \left. \begin{array}{l} 2X \cdot \text{XOR} \\ + \end{array} \right. \\
 1 & 0 & 1 & 0 & 1 & = \text{Cin}(A \oplus B) + AB & \left. \begin{array}{l} 2X \cdot \text{AND} \\ + \end{array} \right. \\
 1 & 1 & 0 & 0 & 1 & S = A \oplus B \oplus \text{Cin} & \left. \begin{array}{l} 1X \cdot \text{OR} \\ + \end{array} \right.
 \end{array}$$

### # Multi-bit Adder (cascade N number of full adders)



Linear increase of delay with number of bits

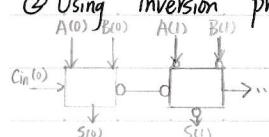
$$T_d(\text{worst}) = (N-1)T_{\text{carry}} + T_{\text{sum}}$$

### Carry chain optimisation

① Express  $S$  in terms of the signal  $C$ :  $S = ABC\text{cin} + (A+B+C\text{cin}) \cdot C'$

Thus for each adder the  $C$  signal is coming earlier than the sum signal.

### ② Using inversion property



Reduce one CMOS invert per 2 adders / Significant decrease time  
将内部(cadder)输出端的CMOS反向器去掉, 第N个adders输出助C位反向, 利用反演属性, 将输入端  
 $p > N, N > p$ , 使第N个adders输出原C位值(第N个adders不去掉CMOS反向器).

### # Subtractor

$$\text{2's complement: } 5_{10} = 1001_2 \quad -5_{10} = 00001001_2 \rightarrow 11110110_2 \rightarrow 11110111_2$$

$$A - B = \overline{A} \oplus B$$

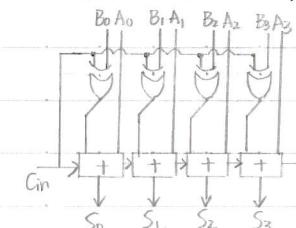
$$A \quad B \quad Z$$

$$0 \quad 0 \quad 0$$

$$0 \quad 1 \quad 1$$

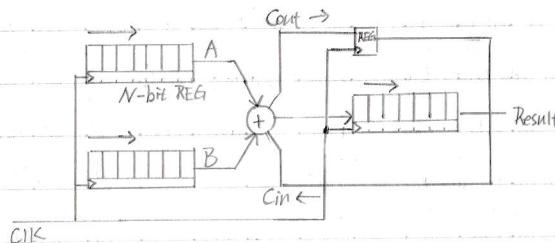
$$1 \quad 0 \quad 1$$

$$1 \quad 1 \quad 0$$



↑ Carry bit  
0: adder  
1: subtractor

## # Bit Serial Adder



$$\text{carry } (t) = A(t) \cdot B(t) + C(t-1) [A(t) + B(t)]$$

$$\text{sum } (t) = \text{carry } (t) [A(t) + B(t) + C(t)] + A(t) \cdot B(t)$$

- ① Equal sum and carry delay
- ②  $N$  number of cycles required
- ③ Good for technology in the range 2-5  $\mu\text{m}$
- ④ Nibble (4 bits) and byte (8 bits) adders are frequently used
- ⑤ Carry-save-addition

# Carry-lookahead Adder (CLA) try to predict  $C_k$  earlier than  $T_c^* K$ 

Look at the 0th stage carry and predict all the carry corresponding to different stages.

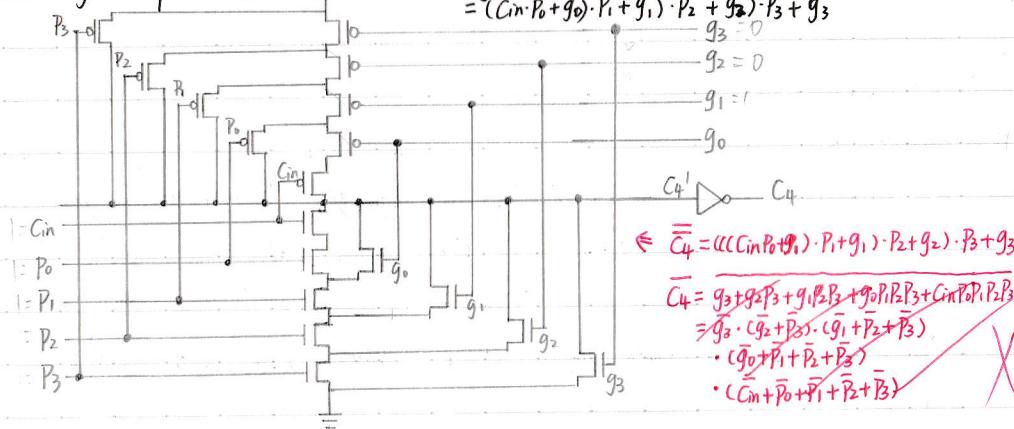
## [Carry propagation]

A	B	Cin	C	$\{ g = AB \}$
0	0	X	0	Kill
0	1	C	C	propagate
1	0	C	C	propagate
1	1	X	1	generate

$$C_4 = \text{OR} \begin{cases} g_3 \\ g_2 \cdot P_3 \\ g_1 \cdot P_2 \cdot P_3 \\ g_0 \cdot P_1 \cdot P_2 \cdot P_3 \\ \text{Cin} \cdot P_0 \cdot P_1 \cdot P_2 \cdot P_3 \end{cases}$$

generate at bit position 3  
... 2, propagated 3  
... 1, ... 2, 3  
... 0, ... 1, 2, 3  
... input, ... 0, 1, 2, 3

## [Static logic implementation]



$$C_4 = g_3 + g_2 \cdot P_3 + g_1 \cdot P_2 \cdot P_3 + g_0 \cdot P_1 \cdot P_2 \cdot P_3 + \text{Cin} \cdot P_0 \cdot P_1 \cdot P_2 \cdot P_3$$

$$= (\text{Cin} \cdot P_0 + g_0) \cdot P_1 + g_1) \cdot P_2 + g_2) \cdot P_3 + g_3$$

$$\bar{C}_4 = \overline{(C_0 \cdot P_0 \cdot P_1 \cdot P_2 \cdot P_3)} = \overline{(g_0 \cdot P_0 + \bar{P}_0) \cdot P_1 + g_1) \cdot P_2 + g_2) \cdot P_3 + g_3}$$

$$\bar{C}_4 = \overline{g_3} + \overline{g_2 \cdot P_3} + \overline{g_1 \cdot P_2 \cdot P_3} + \overline{g_0 \cdot P_1 \cdot P_2 \cdot P_3} + \overline{\text{Cin} \cdot P_0 \cdot P_1 \cdot P_2 \cdot P_3}$$

$$= \overline{g_3} + \overline{g_2} + \overline{P_3} \cdot \overline{g_1} + \overline{P_2} \cdot \overline{P_3} \cdot \overline{g_0} + \overline{P_1} \cdot \overline{P_2} \cdot \overline{P_3} \cdot \overline{g_0}$$

P-MOS passes 11 very efficiently  
N-MOS passes 10 very efficiently

$$\downarrow f(A, B)$$

$$C_4 = g_3 + g_2 \cdot P_3 + g_1 \cdot P_2 \cdot P_3 + g_0 \cdot P_1 \cdot P_2 \cdot P_3 + \text{Cin} \cdot P_0 \cdot P_1 \cdot P_2 \cdot P_3$$

$$= (\text{Cin} \cdot P_0 + g_0) \cdot P_1 + g_1) \cdot P_2 + g_2) \cdot P_3 + g_3$$

$$f(\bar{A}, \bar{B})$$

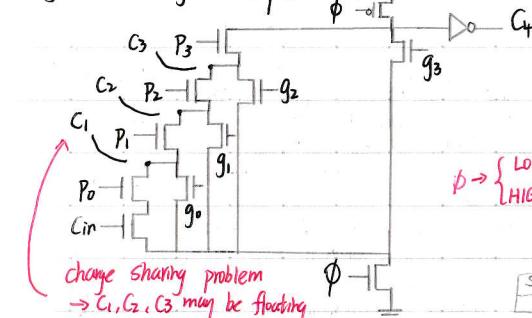
$$C_4 = \overline{g_3} \cdot \overline{g_2} \cdot P_3 + \overline{g_1} \cdot P_2 \cdot P_3 + \overline{g_0} \cdot P_1 \cdot P_2 \cdot P_3 + \text{Cin} \cdot P_0 \cdot P_1 \cdot P_2 \cdot P_3$$

$$= P_3 (g_3 + P_2 (g_2 + P_1 (g_1 + P_0 (g_0 + \text{Cin})))$$

$$\bar{C}_4 = \bar{P}_3 + \bar{g}_3 \cdot (\bar{P}_2 + \bar{g}_2 \cdot (\bar{P}_1 + \bar{g}_1 \cdot (\bar{P}_0 + \bar{g}_0 \cdot \bar{\text{Cin}})))$$

$$\left\{ \begin{array}{l} C_4 = (\text{Cin} \cdot P_0 + g_0) \cdot P_1 + g_1) \cdot P_2 + g_2) \cdot P_3 + g_3 \\ \bar{C}_4 = (\text{Cin} \cdot \bar{P}_0 + \bar{g}_0) \cdot \bar{g}_1 + \bar{P}_1) \cdot \bar{g}_2 + \bar{P}_2) \cdot \bar{g}_3 + \bar{P}_3 \end{array} \right.$$

## [Dynamic Logic Implementation]



$$C_4 = g_3 + P_3 (g_2 + P_2 (g_1 + P_1 (g_0 + \text{Cin})))$$

↳ { Low : Pre-charge  
High : Set  $C_4$  }

S D G ↳ charge the MOSFET capacitance

## [CLA - Summary]

CLA compared to ripple-carry adder : { Faster  
Larger Area }

Limitation: cannot go beyond 4 bits of look-ahead

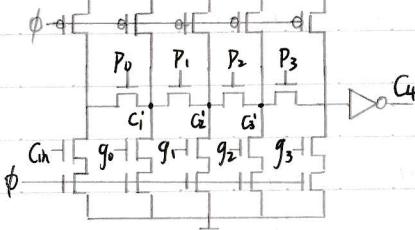
↳ large  $P, g$  fan-out slows down carry generation

Way out: Manchester carry chains

↳ Tries to reuse logic by pre-charging each carry position

## # Manchester Carry Chain

Precharge internal nodes to avoid charge-sharing problem



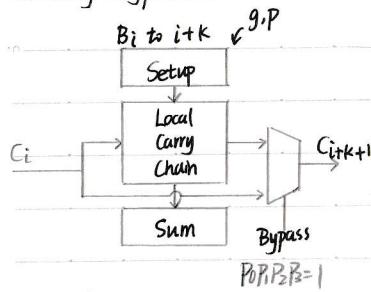
Compared to CLA:

- ① Smaller area
- ② Cin close to the output

Carry chain can be any length



## # Carry-bypass adder



## [Pass mode]

Local carry vector computes intermediate carries (possibly incorrect). At the same time, mux selection set to pass. When input arrives, intermediate carries might be recomputed. Meanwhile, input carry is sent to Cont.



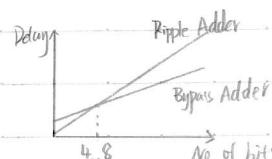
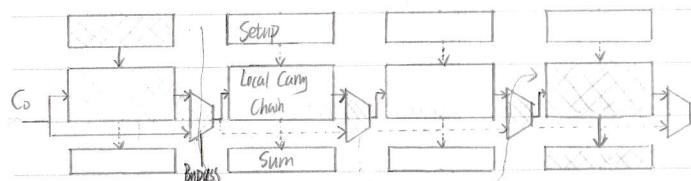
## [If not pass mode]

Local carry vector computes intermediate carries (bits 10, 11 correct). At the same time, mux selection set to local. Meanwhile, output carry is sent to Cont correctly. When input carry arrives, intermediate carries C8 and C9 will be recomputed correctly.



## [Timing]

每个Cont同时产生



$$\text{Delay} = t_{\text{setup}} + \max\{t_{\text{select}}, 4t_{\text{FA}}\} + 3 \times t_{\text{mux-pass}} + 4 \times t_{\text{FA}} + t_{\text{sum}}$$

Faster than ripple adder but still linear

## [Exam]

b 8 0 2

$$A_{\text{xor}} = 1 \mu\text{m}^2$$

$$A_{\text{and}} = 0.5 \mu\text{m}^2$$

$$A_{\text{fa}} = 2 \mu\text{m}^2$$

$$A_{2:1 \text{ multi}} = 0.2 \mu\text{m}^2$$

$$A_{\text{setup}} + A_{\text{acc}} + A_{\text{multi}} = 14.2 \mu\text{m}^2$$

$$A_{\text{setup}} = 4 \cdot (A_{\text{and}} + A_{\text{xor}}) = 4 \times 1.5 = 6 \mu\text{m}^2$$

$$A_{\text{acc}} = 4 \times 2 = 8 \mu\text{m}^2$$

$$A_{2:1 \text{ multi}} = 0.2 \mu\text{m}^2$$

$$A_{\text{all}} = 3 \times 14.2 \mu\text{m} + 14$$

$$t_{\text{and}} = 1$$

$$t_{\text{xor}} = 2$$

$$t_{\text{fa/bit}} = 4$$

$$t_{\text{sum}} = 2$$

$$t_{\text{2:1 mux}} = 1$$

$$t_{\text{setup}} = 2$$

$$t_{\text{FA}} = 4 \times 4 = 16$$

$$t_{\text{sel}} = P_0 P_1 P_2 P_3 = 2$$

$$t_{\text{mux}} = 3 \times 1 = 3$$

$$2 + 16 + 3 + (6 + ?)$$

$\nearrow$

## Multiplier

$$\begin{array}{r} 1101 \\ \times 0101 \\ \hline 1101 \\ 0000 \\ \hline 01000001 \end{array}$$

← Multiplicand 被乘数  
← Multiplier 乘数  
← Product

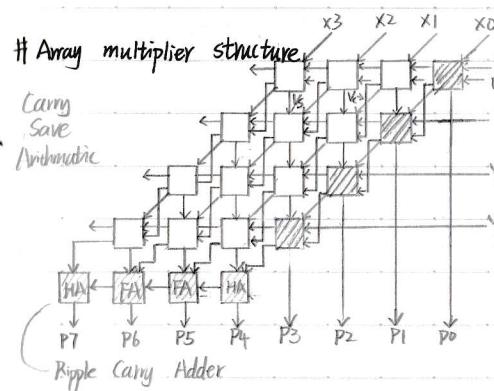
# Unsigned Digital Multiplication

If length of the multiplicand is  $N$ -bit  
and length of the multiplier is  $R$ -bit, then  
length of the product  $P = N+R$ .

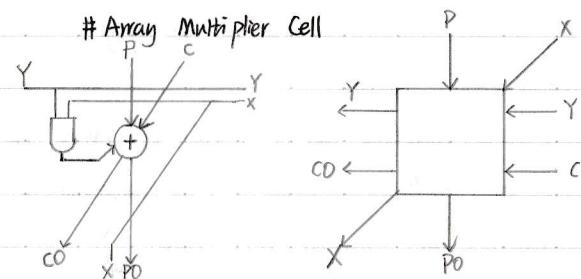
# Parallel Array Multiplier 并行阵列

$$\begin{array}{ccccccc} x_3 & x_2 & x_1 & x_0 & & & \\ y_3 & y_2 & y_1 & y_0 & & & \\ \hline x_3 y_0 & x_2 y_0 & x_1 y_0 & x_0 y_0 & & & \\ x_3 y_1 & x_2 y_1 & x_1 y_1 & x_0 y_1 & & & \\ x_3 y_2 & x_2 y_2 & x_1 y_2 & x_0 y_2 & & & \\ \hline x_3 y_3 & x_2 y_3 & x_1 y_3 & x_0 y_3 & & & \\ \hline p_7 & p_6 & p_5 & p_4 & p_3 & p_2 & p_1 & p_0 \end{array}$$

## # Array Multiplier structure



## # Array Multiplier Cell



In each stage, LSB drops out to form part of partial product. The accumulating product is shifted right one place by wiring between each stage.

Total

 $6 \cdot T_{FA} + T_{RA}$  for each line

## # Signed Binary Multiplication

$$\begin{array}{r} 1010 \\ \times 0101 \\ \hline 1111010 \\ \text{Sign extension} \\ \rightarrow 0000000 \\ \hline 111010 \\ 000000 \\ \hline 11100010 \end{array}$$

← Signed  
Multiplicand  
← 2's complement of 30

$$\begin{array}{r} 0110 \\ \times 1011 \\ \hline 00000110 \\ 0000110 \\ 000000 \\ \hline 11010 \\ 11100010 \end{array}$$

← Signed  
Multiplier

## # Booth's algorithm

For  $i$ th bit of the multiplier

$r_i$	$r_{i-1}$	$r_{i-1} - r_i$	Action	$r_i$	$r_{i-1}$	$r_{i-1} - r_i$	Action
0	0	0	Add zero (do nothing)	*	1	0	$\leftarrow 2^i$ complement of 6
1	0	-1	Subtract multiplicand	1	1	0	$\leftarrow 1-1=0$ Add zero
0	1	1	Add multiplicand	0	0	1	$\leftarrow 1-0=1$ Add 6
1	1	0	Add zero (do nothing)	1	0	0	$\leftarrow 0-1=-1$ $2^i$ complement of 6

$$\begin{array}{r}
 0110 \leftarrow 6 \\
 \times 1011 \leftarrow -5 \\
 \hline
 11111010 \leftarrow 0-1=-1 \text{ } 2^i \text{ complement of 6} \\
 00000000 \leftarrow 1-1=0 \text{ Add zero} \\
 000110 \leftarrow 1-0=1 \text{ Add 6} \\
 \hline
 11010 \leftarrow 0-1=-1 \text{ } 2^i \text{ complement of 6}
 \end{array}$$

Not necessary to know when we come to the MSB and switch to a subtraction

## # Modified Booth Algorithm

A faster multiplier can be realised by looking at two bit

For  $i$ th and  $(i+1)$ th bit of multiplier, the algorithm is:

$r_{i+1}, r_i$	$r_{i-1}$	$\{r_{i+1}, r_i\} + r_{i-1}$	Action
(0) 0 0	0	0	Do nothing
(1) 0 1	0	1	Add 1 <sup>st</sup> Multiplicand
(-2) 1 0	0	-2	Subtract 2 <sup>nd</sup> Multiplicand
(-1) 1 1	0	-1	Subtract 1 <sup>st</sup> Multiplicand
(0) 0 0	1	1	Add 1 <sup>st</sup> Multiplicand
(1) 0 1	1	2	Add 2 <sup>nd</sup> Multiplicand
(-2) 1 0	1	-1	Subtract 1 <sup>st</sup> Multiplicand
(-1) 1 1	1	0	Do nothing

$$\begin{array}{r}
 0110 \\
 \times 1001 \\
 \hline
 11111010 \leftarrow 2^i \text{ complement of 6} \\
 111010 \leftarrow \text{Subtract } 6 \text{ (进两位)}
 \end{array}$$

Grouping multiplier bits into pairs  $\Rightarrow$  Reduces number of partial products to half

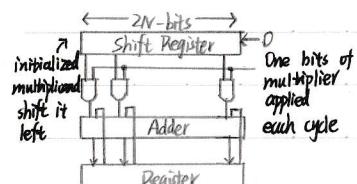
Use of higher radix reduces number of intermediate addition

- Radix-8 should implement  $*3, *-3, *4, *-4$

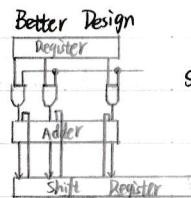
- Complicated, trade-off has to be made

Depends on architecture

## # Sequential Multiplier



{ 2N-bit adder  
2N-bit register }



Shift result register to right  
{ N-bit And gate  
N-bit Adder }

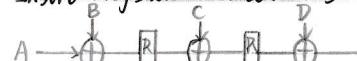
## # Shift and Rotate Operation

Difference between { Shifter : shifts any bits

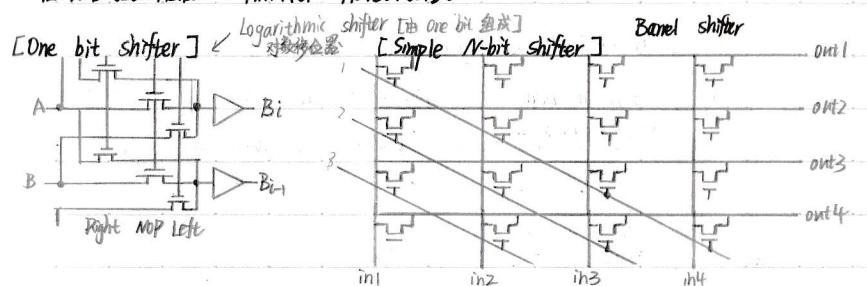
Shift register : shifts only one bit

## # Pipelined multiplier

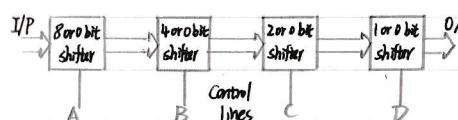
Insert registers between rows



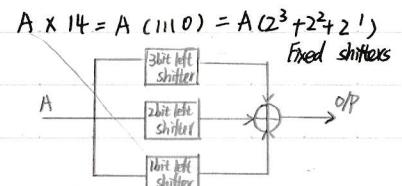
T0: A0 B0 C0 D0 A0+B0  
T1: A1 B1 C1 D1 A1+B1  
T2: A2 B2 C2 D2 A2+B2  
A0B0C0D0  
A0B0C0  
A0B0C0D0  
A0B0C0D0



## # Programmable shifter 16-bit



## # Multiplication using shifter



## # Summary

Trade-offs between area and delay

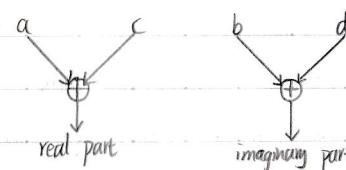
- Barrel shifter: fastest O(1),  $N^2$  transistors
- Logarithmic shifter:  $O(\log N)$ ,  $N \log N$  transistors
- One bit shifter:  $O(N)$ ,  $N$  transistors

Barrel shifter: wire dominated circuit

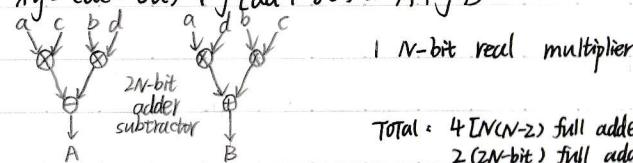
### # Complex numbers and addition

$$x = a + b \cdot j \quad \begin{cases} a: \text{real part} \\ b: \text{imaginary part} \end{cases}$$

$$(a+j \cdot b) + (c+j \cdot d) = (a+c) + j(c+d)$$



$$xy = (ac - bd) + j(ad + bc) = A + j \cdot B$$



$\left\{ \begin{array}{l} N(N-2) \text{ full adders} \\ N \text{ half adders} \\ N^2 \text{ AND gates} \end{array} \right.$

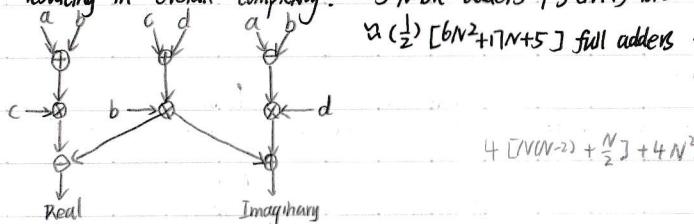
$$\begin{aligned} \text{TOTAL: } & 4[N(N-2) \text{ full adders} + (N/2) \text{ full adders} + N^2 \text{ AND}] + \\ & 2(2N \text{-bit}) \text{ full adders} \\ & = 2N(2N-1) \text{ full adders} + N^2 \text{ AND gates} \end{aligned}$$

Reduce complexity

$$\Rightarrow xy = (ac - bd) + j(ad + bc) = [(ca+b) - b(d+c)] + j[d(a-b) + b(c+d)]$$

Resulting in overall complexity:  $3N$

$$\approx \frac{1}{2} [6N^2 + 17N + 5] \text{ full adders} + 3(N+1)^2 \text{ AND gates.}$$



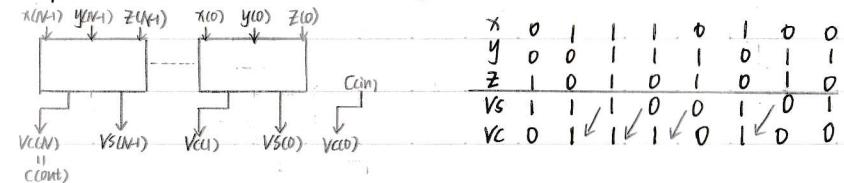
$$4[V(N-2) + \frac{N}{2}] + 4N^2$$

### Redundant Arithmetic

#### # Carry-Save addition (CSA)

Perform an addition of three vectors (or inputs) using an array of 1-bit adders without propagating the carries.

The sum  $r$  of 3 numbers  $x, y, z$  is represented as two vectors - carry vector ( $vc$ ) and pseudo-sum vector ( $vs$ ).  $x + y + z = vc + vs = r$

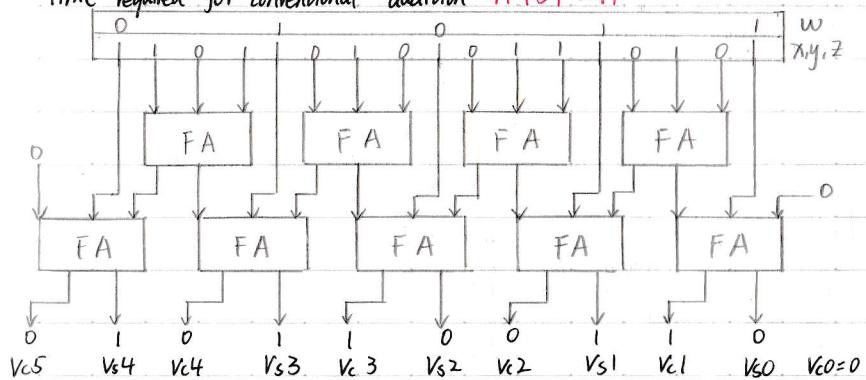


The final result of CSA is in carry-save format. Just add  $vs$  and  $vc$ .

#### # Multiple operand addition

$$\begin{cases} x = 1010 \\ y = 0111 \\ z = 1000 \end{cases} \Rightarrow \begin{cases} vs = 00101 \\ vc = 10100 \end{cases} \Rightarrow 00101 + 10100 = 11001$$

Considering one full adder delay =  $T$   
Time required for carry-save addition  $T + 5T = 6T$  (Need a 5-bit adder for converting)  
Time required for conventional addition  $4T + 5T = 9T$



4:2 adder

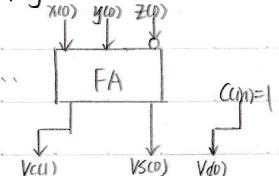
Conventional addition delay =  $4T + 5T + 5T = 14T$   
Carry-save delay (include conversion) =  $2T + 5T = 7T$

Actual carry-save delay is independent of the wordlength until you do conversion

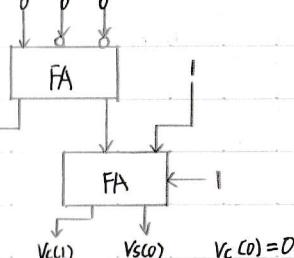
# Minus - Minus - Plus (MAP)

# Subtraction

$x + y - z$  [Plus - Plus - Minus operation]



$x - y - z$  [Minus - Minus - Plus]



[Determine sign of CSA number]

In conventional 2's complement sign can be determined by checking the MSB

Not applicable for CSA number

Sign can be determined by adding  $v_s$  and  $v_c$ : speed advantage gone

Use generate-propagate formulation.

Take  $v_{ci}$  and  $v_{ci+1}$  as a bit pair and compute propagate  $p_{ci}$  starting from MSB position until you encounter 11 or 00 condition.

- ① If  $p_{(N-1)} \cdot p_{(N-2)} \dots p_{(k+1)} = 1$  and  $g_{(k)} = 1 \Rightarrow$  Positive Number
- ② If  $p_{(N-1)} \cdot p_{(N-2)} \dots p_{(k+1)} = 1$  and  $g_{(k)} > 0$  and  $p_{(k)} = 0 \Rightarrow$  Kill condition = Negative number
- ③ If  $g_{(N-1)} = 1 \Rightarrow$  Negative Number

[The Logic Function]

0 as output will signify a +ve number and 1 will signify a -ve number

Considering detecting a -ve number, i.e., output logic function sign should be 1

Adding two 2's complement number will be -ve when

①  $g_3 = 1$ : both numbers are -ve

②  $p_3k_2 = 1$ : bit 3 propagate, bit 2 kill (00 condition  $k_{(i)} = a_{(i)} \cdot b_{(i)}$ )

③  $p_3p_2k_1 = 1$ : bit 3 and bit 2 propagate, bit 1 kill

④  $p_3p_2p_1k_0 = 1$ : bit 3, 2 and 1 propagate, bit 0 kill

⑤  $p_3p_2p_1p_0 = 1$ : all bits propagate

Therefore  $\text{sgn} = p_3p_2p_1p_0 + p_3p_2p_1k_0 + p_3p_2k_1 + p_3k_2 + g_3$

If  $\text{sgn} = 1$ , the number is

[EXAMPLE]

$g_3 = 1$ : addition of two -ve numbers is -ve number

$p_3k_2 = 1$ :  $1011 + 0011 = -5 + 3 = -2$

$p_3p_2k_1 = 1$ :  $1001 + 0101 = -7 + 5 = -2$

$p_3p_2p_1k_0 = 1$ :  $1010 + 0100 = -6 + 4 = -2$

$p_3p_2p_1p_0 = 1$ :  $1010 + 0101 = -6 + 5 = -1$

$1011 \quad 1001 \quad 1010 \quad 1010$

$0011 \quad 0101 \quad 0100 \quad 0101$

$p_3k_2 \quad p_3p_2k_1 \quad p_3p_2p_1k_0 \quad p_3p_2p_1p_0$

$1 \quad 1 \quad 1 \quad 1$

Therefore, in the worst case we need to block all the bit pairs - a propagation chain exists from the LSB to MSB, slow operation particular for longer wordlength  
 Way out: partition the wordlength and apply the criteria in parallel and finally combine  
 In this case in the sgn expression use  $g_{(i-1)}$  only for the final block  
 For the intermediate blocks replace generate function by kill function

Consider block of 4 the sgn expression for 16-bit wordlength is

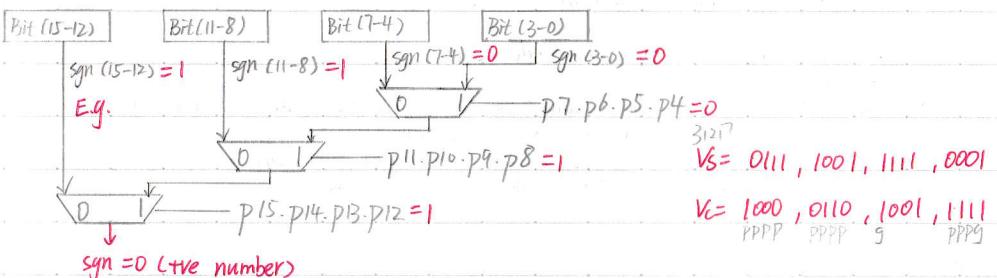
$$\text{sgn}_{(15-12)} = p_{15} \cdot p_{14} \cdot p_{13} \cdot p_{12} + p_{15} \cdot p_{14} \cdot p_{13} \cdot k_{12} + p_{15} \cdot p_{14} \cdot k_{13} + p_{15} \cdot k_{14} + g_{15}$$

$$\text{sgn}_{(11-8)} = p_{11} \cdot p_{10} \cdot p_9 \cdot p_8 + p_{11} \cdot p_{10} \cdot p_9 \cdot k_8 + p_{11} \cdot p_{10} \cdot k_9 + p_{11} \cdot k_{10} + k_{11}$$

$$\text{sgn}_{(7-4)} = p_7 \cdot p_6 \cdot p_5 \cdot p_4 + p_7 \cdot p_6 \cdot p_5 \cdot k_4 + p_7 \cdot p_6 \cdot k_5 + p_7 \cdot k_6 + k_7$$

$$\text{sgn}_{(3-0)} = p_3 \cdot p_2 \cdot p_1 \cdot p_0 + p_3 \cdot p_2 \cdot p_1 \cdot k_0 + p_3 \cdot p_2 \cdot k_1 + p_3 \cdot k_2 + k_3$$

[Sign extractor circuit - 16 bit CSA number]



# Summary

CSA can give us the advantage of carry-free or limited carry propagation - very fast addition

Addition operation is independent of wordlength

Higher radix redundant arithmetic (in particular radix-4) has been successfully used in high-speed datapath design

Finding sign of a CSA number is difficult - the sign is given by the most significant non-zero digit

More than one bit (more wires) are needed for representing a digit

Conversion to conventional 2's complement requires traditional adders

Not good for non-arithmetic purpose

## ⇒ Low Power Design

### # Short Circuit Power Dissipation

The input of a real life logic gate (e.g., an inverter) is actually driven by input waveforms that has finite rise and fall time

- Both the NMOS and PMOS transistors may conduct simultaneously for a very short time during switching, forming a direct current path between the supply and ground.
- The NMOS starts conducting when rising input voltage exceeds the threshold voltage of the NMOS transistor
- The PMOS remains "on" until the input reaches the voltage level ( $V_{DD} - I \cdot R_D$ ), creating a time window both when both transistors are on.

The current drawn from the supply when the PMOS and NMOS transistor both conduct does not contribute to the capacitance charging

- Hence the name short circuit current.

Maybe reduced by making the output voltage transition time longer and/or by making input voltage transition time smaller : should be balanced carefully

### # Leakage Power dissipation

NMOS and PMOS transistors generally have nonzero reverse leakage and subthreshold currents.

- so even when the gate is idle, it still draws a very small amount of current.
- Mainly controlled by the process parameters
- For a chip having millions of transistors, this leakage power dissipation may contribute significantly.

Two main components : biggest concern in deep submicron (近々)

- Reverse gate leakage current
  - Tunnelling of electrons through gate oxide ( $O\text{gA}$ )
- Subthreshold current
  - A weak current that flows due to carrier diffusion between the source and drain of a transistor in weak inversion region ( $V_{GS}$  is less than the threshold voltage but not exactly zero)
  - If  $V_{GS}$  is smaller than but very close to the threshold voltage then the subthreshold power dissipation may become comparable in magnitude to the switching power dissipation

Process level:

- channel level
  - Scaling the vertical dimension of the MOSFET like oxide thickness, junction depth
  - Substrate doping : to control the depletion width
- Well engineering
  - change of doping profile in the channel : <sup>RAKE</sup>retrograde doping (1-D characteristics of well profile), hub doping (localized 2-D dopant distribution near S/D extension)

Circuit level : Controlling voltages at different device terminals.

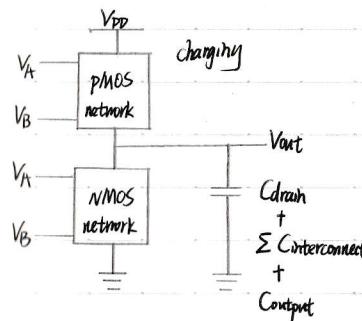
- Standby leakage control using transistor stack
- Multiple threshold design
- Sleep transistor insertion

## Low Power Design

### # Why

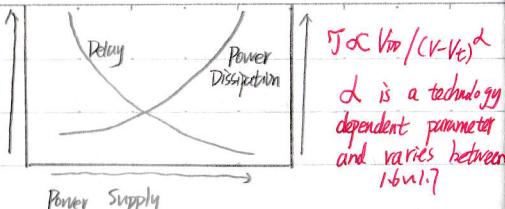
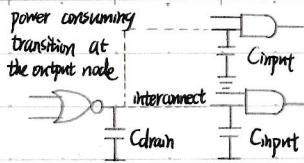
- Increasing prominence of portable and mobile systems
  - Requirement of low power dissipation to maximise battery life
- High-performance microprocessors and DSPs require complex chips with high clock frequency
  - Power dissipation and hence temperature increase linearly with frequency and complexity
- Reliability
  - Close correlation exists between the peak power dissipation of digital circuits and reliability
- General-purpose microprocessors and DSPs are not a good solution
  - Each instruction execution typically requires a number of clock cycles and thus, the power dissipation is at least 1 or 2 orders of magnitude higher than the hardware approach
- Technology level
  - Selection of right device geometry, interconnect and process
- Circuit level
  - proper choice of circuit design styles, reduction of voltage swing, appropriate clocking strategy
- Architecture level
  - Smart power management strategy for various blocks, utilization of pipelining and parallelism 流水线 并行性
- Algorithm level
  - Proper selection of the data processing algorithms, specifically to minimise the number of switching events for a given task
- Dynamic Power
  - Power dissipation due to switching of different nodes of a CMOS circuit
  - Predominant power dissipation for CMOS circuits up to  $0.13\mu\text{m}$
- Short-circuit power dissipation
  - Power dissipation due to non-ideal switching characteristics of CMOS logic gates
- Leakage power
  - Power dissipation due to non-zero reverse leakage and subthreshold currents in the NMOS and PMOS transistors
  - Becoming nearly same as the dynamic power dissipation for deep sub-micron technology

### # Dynamic (Switching) Power Dissipation



#### Dynamic Power

- Dynamic power is dissipated when energy is drawn from power supply to charge up the output node capacitance
- During the charge-up phase, the output voltage typically makes a full rail transition from 0 to  $V_{DD}$ , and the energy drawn is relatively independent of the function performed by the circuit
- In a complex structure the load capacitance that has to be switched typically consists of
  - The output capacitance of the driving gate itself
  - Total interconnect capacitance
  - The input characteristic  $\text{dr}^+$  capacitances of the driven gates.



## # Dynamic Power Dissipation Expression

Average dynamic power can be calculated by considering a step input where the rise and fall times are much less than the pulse repetition period.  
Let frequency  $f = 1/t_p$ ,  $t_p$  being the time period

$$P_d = \frac{1}{t_p} \int_0^{t_p} i_{in}(t) V_{out} dt + \frac{1}{t_p} \int_{t_p}^{2t_p} i_{in}(t) (V_{DD} - V_{out}) dt$$

$$i_{in} = C_L \frac{dV_{out}}{dt} \text{ and } i_p = C_L \frac{d(V_{DD} - V_{out})}{dt}$$

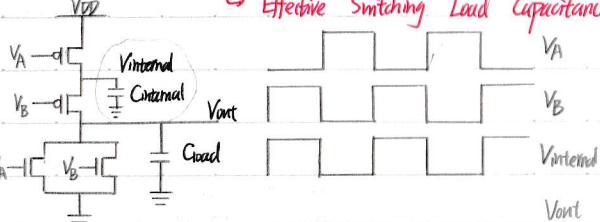
Substituting these values we get

$$P_d = \frac{C_L}{t_p} \int_{V_{DD}}^{V_{out}} V_{out} dV_{out} + \frac{C_L}{t_p} \int_{V_{DD}}^{V_{out}} (V_{DD} - V_{out}) d(V_{DD} - V_{out}) = \frac{C_L V_{DD}^2}{t_p} = C_L \cdot V_{DD}^2 \cdot f$$

So the average dynamic power is proportional to the load capacitance and frequency and quadratically proportional to supply voltage

$$\boxed{P_d = \sum_i G_i \cdot C_i \cdot V_{DD}^2 \cdot f}$$

Effective Switching Load Capacitance



Internal nodes are making power consuming transitions but output node is not showing the effect effect

$$\begin{aligned} I=1: & P_d = \frac{1}{t_p} \int_0^{t_p} V_{out} \cdot i_{in} dt \\ & \left\{ \begin{array}{l} i_{in} = C_L \frac{dV_{out}}{dt} \\ i_p = C_L \frac{d(V_{DD} - V_{out})}{dt} \end{array} \right. \\ I=0: & P_d = \frac{1}{t_p} \int_{t_p}^{2t_p} (V_{DD} - V_{out}) \cdot i_p dt \\ & P_d = \frac{1}{t_p} \int_0^{t_p} V_{out} \cdot C_L \frac{dV_{out}}{dt} dt + \frac{1}{t_p} \int_{t_p}^{2t_p} (V_{DD} - V_{out}) \cdot C_L \frac{d(V_{DD} - V_{out})}{dt} dt \\ & = \frac{C_L}{t_p} \left[ \frac{1}{2} V_{out}^2 \right]_0^{t_p} + \frac{C_L}{t_p} \int_{t_p}^{2t_p} \left[ \frac{1}{2} (V_{DD} - V_{out})^2 \right] dt \\ & = \left( \frac{1}{2} V_{DD}^2 + \frac{1}{2} V_{DD}^2 + 0 \right) \cdot \frac{C_L}{t_p} = V_{DD}^2 \cdot \frac{C_L}{t_p} = C_L \cdot V_{DD}^2 \cdot f \end{aligned}$$

$$\frac{1}{4 \times 10^{-3}} = \frac{1}{4} \times 250$$

## # Generalised average dynamic power dissipation expression

$$P_{avg} = \left( \sum_i d_i \cdot T_i \cdot C_i \cdot V_i \right) \cdot V_{DD} \cdot f_{clk}$$

$d_i$  is the switching probability of the  $i$ th node capacitance  $C_i$ .  
(The assumption that all the nodes switches at every clock cycle is not correct and thus, the number of power consuming transistors for the nodes is a probabilistic function and depends on the input state.)  
 $V_i$  is the voltage swing of the  $i$ th internal node  
In general,  $V_i = V_{DD}$

## # Low power by Reduction of Switching Activity

Switching activity ( $d_i$ ) gives the probabilistic measurement of number of transitions of a particular node

- Number of times we need to charge the capacitance at that node, i.e. power consumption
- How to calculate switching activity (activity factor)?
- ① Consider that at the primary input logic 0 and 1 have equal chance of occurrence and they are uncorrelated
- ② Find out the probability of occurrence of a logic 1 at the output node ( $P_{O1}$ )
- ③ Find out the probability of occurrence of a logic 0 at the output node ( $P_{O0}$ )
- ④ The power consuming transition (switching activity) can be derived as  $P_{O1} = P_{O0} = P_{O1} \cdot P_{O0}$

$$P_{avg} = \sum_i d_i \cdot C_i \cdot V_{DD}^2 \cdot f$$

#  $d_i = P_{O1} \cdot P_{O0}$   
Cont = 1pf  
 $C_m = 1\text{pf}$

$$\begin{array}{ll} A & B \quad X \Rightarrow P(O1) = \frac{1}{4} \quad (\text{probability}) \\ D & D \quad 0 \quad 0 \quad P(O0) = \frac{3}{4} \\ Z & C \quad 0 \quad 1 \quad P_{O1} = P(O1) \cdot P(O0) = \frac{3}{16} \\ C & D \quad 1 \quad 0 \quad P_{O0} = P(O1) \cdot P(O0) = \frac{3}{16} \\ D & Y \quad 1 \quad 1 \end{array}$$

$$P_{avg} = [d(X) C_x + d(Y) C_y + d(Z) C_z] V_{DD}^2 \cdot f$$

# XOR  
 $P(X) = \frac{1}{4} \quad P(O1) = \frac{3}{4} \Rightarrow P_{O1}(X) = \frac{3}{4} \cdot \frac{1}{4} = \frac{3}{16} \quad d(X)$   
 $P(Y) = \frac{3}{4} \quad P(O2) = \frac{1}{4} \Rightarrow P_{O2}(Y) = \frac{3}{4} \cdot \frac{1}{4} = \frac{3}{16} \quad d(Y)$

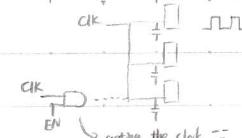
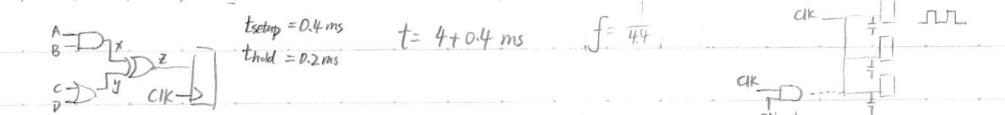
$$\begin{array}{ll} Z = \bar{X} \cdot Y + X \cdot \bar{Y} & \Rightarrow P_{O1}(Z) = P(O1) \cdot P(O2) + P(O1) \cdot P(O2) \\ & = \frac{3}{4} \cdot \frac{3}{4} + \frac{1}{4} \cdot \frac{1}{4} = \frac{10}{16} = \frac{5}{8} \end{array}$$

$$P_{avg} = \frac{5}{8} \cdot P_{O1} \cdot P_{O2} = \frac{15}{64} \quad \boxed{d(Z)}$$

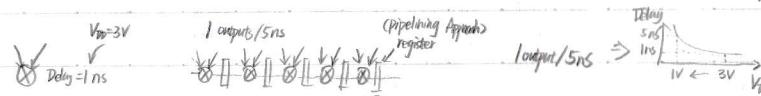
$$\begin{array}{ll} C_A = C_m + C_{out} = 2\text{pf} & P_X(1) = P_{O1}(1) \cdot P_{O2}(1) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} \\ C_Y = C_m + C_{out} = 2\text{pf} & P_X(0) = 1 - P_X(1) = \frac{3}{4} \\ C_Z = C_{out} = 1\text{pf} & P_{O1}(X) = P(O1) \cdot P(O2) = \frac{3}{16} \\ & P_{O2}(Y) = P(O1) \cdot P(O2) = \frac{3}{16} \\ & P_{O1}(Z) = 1 - P_{O1}(1) = \frac{3}{4} \\ & P_{O2}(Y) = P(O1) \cdot P(O2) = \frac{3}{16} \end{array}$$

$$\begin{array}{ll} t_{X+dday} = 1\text{ms} & t = 4\text{ms} \Rightarrow f = 250\text{Hz} \\ t_1 = 2\text{m} & t_2 = 2\text{m} \\ t_3 = 2\text{m} & \\ & V_{DD} = 3.3V \\ & \dots \end{array}$$

$$P_{avg} = \left[ \frac{3}{16} \cdot 2 \times 10^{-6} + \frac{3}{16} \times 2 \times 10^{-6} + \frac{15}{64} \cdot 1 \times 10^{-6} \right] \cdot 3.3V \cdot 250\text{Hz}$$



huge current is required to charge all the flip-flops when EN goes to 1 from 0

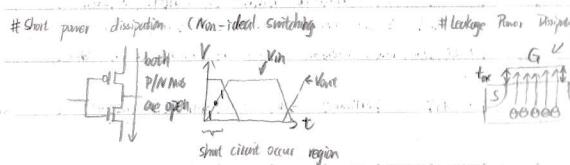


not anything new  $\gg N-1$

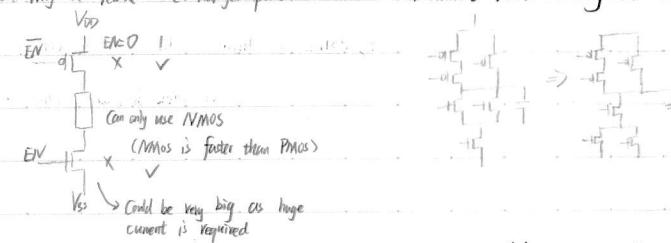
Using parallel to reduce the specification constraint and then reduce V<sub>D</sub> to get the suitable delay value

become faster

↳ Increase the propagation delay



Best way to reduce the leakage power



Logically equivalent circuits could have different power consumptions

The one of the right is better because of low capacitance on output node

## # Power Reduction by Clock Gating

Power consumption breakdown in MCU and ASIC

Clock gives to all the registers in the chip.

Even if some part of the chip is not working, still it is dissipating power because of clock activity at the input of the register. Total wastage

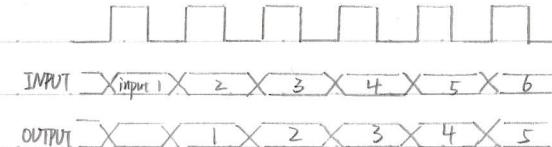
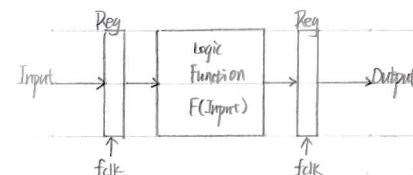
Reduce signal toggling when unnecessary

Pros:

Increases validation complexity

Timing budget may not allow it

## # Low Power Design Using Pipelining Approach



Total capacitance ( $C_{\text{total}}$ ) consists of

- The capacitances switched in the input register array
- The capacitances switched to implement the logic function
- The capacitances switched in the output register array

The expression for power dissipation

$$\text{Preference} = C_{\text{total}} \cdot V_{DD}^2 \cdot f_{\text{clk}}$$

Create an  $N$ -register pipeline for the circuit

- Logic function  $F(\text{Import})$  has to be partitioned into  $N$  successive stage each having approximately equal delay
- A total of  $(N-1)$  register arrays have to be introduced in addition to the original input and output registers
- All registers are clocked at  $f_{\text{clk}}$  (through it is possible to increase clock speed to  $N \cdot f_{\text{clk}}$ ) and thus, delay for each logic block can be increased by a factor of  $N$  without sacrificing throughput

Physically, we are running the logic segments below their actual speed capacity

- This gives us a chance to lower the supply voltage while maintaining the clock rate of  $f_{\text{clk}}$

$$\text{Pipeline} = [C_{\text{total}} + (N-1) C_{\text{reg}}] \cdot V_{DD, \text{new}}^2 \cdot f_{\text{clk}}$$

$$\text{Total power reduction factor: } \frac{\text{Pipeline}}{\text{Preference}} = \left[ 1 + \frac{C_{\text{reg}}}{C_{\text{total}}} (N-1) \right] \cdot \frac{V_{DD, \text{new}}^2}{V_{DD}^2}$$

$C_{\text{reg}}$  represents the capacitance switched by each pipeline reg

## [Example]

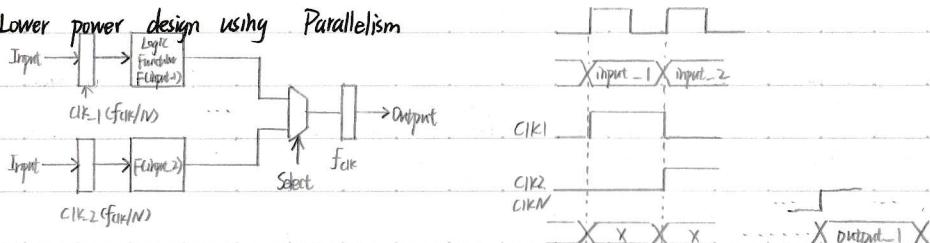
Consider replacing a single-stage logic block ( $V_{DD}=5V$ ,  $f_{\text{clk}}=20\text{MHz}$ ) with a four-stage pipeline, running same  $f_{\text{clk}}$

- The propagation delay of each pipeline stage can be increased by a factor of 4 without sacrificing data throughput
- Assume the magnitude of threshold voltage for each transistor is  $0.8V$
- The designed desired speed reduction (4 times) can be achieved by reducing the supply voltage from  $5V$  to  $\approx 2V$
- Assume  $C_{\text{reg}}/C_{\text{total}}=0.1$

$$\text{Then the overall power reduction factor} = \frac{1}{5} \cdot (1+0.1 \cdot 3) \cdot \frac{4}{25} = \frac{5.2}{25}$$

Thus a dynamic power saving of about 80% is achieved while maintaining the same throughput

## # Lower power design using Parallelism



Assume that the consecutive input data arrive at the same rate as in the single-stage case

The input data are routed to all registers in  $N$  processing blocks

- External clock signals with clock period of  $N$  Tclk are used to load each register every  $N$  clock cycles.
- The individual clock signals are skewed by Tclk, such that each one of the  $N$  consecutive data is loaded into a different input register
- Each register is thus clocked at a lower frequency  $f_{clk}/N$ , the time allowed to compute the function for each input data is increased by a factor of  $N$

The supply voltage can be reduced until the path delay equals the new clock period of  $N$  Tclk

$$\text{Parallel} = [1 + \frac{\text{Cyc}}{\text{Clock}}] \cdot \text{Total: } V_{DD,\text{new,filt}}$$

$$\text{Parallel} = \frac{V_{DD,\text{new}}^2}{V_{DD}} \cdot [1 + \frac{\text{Cyc}}{\text{Clock}}]$$

A judicious balancing between the pipelining and parallelism depending on specific circuit circumstances is needed to make the architecture efficient in terms of power dissipation

## # Penalty of Pipelining and Parallelism

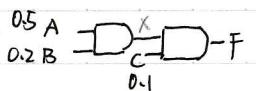
Pipeline: ① Increasing area

② Increase of latency from 1 cycle to  $N$  cycles

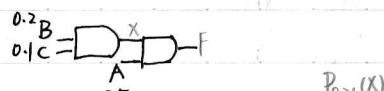
Parallelism: ① Increase area and latency

② Silicon area will grow faster than the number of processor because of signal routing and the overhead circuitry

Input Reordering to reduce Switching Activity



$$(1-0.5 \times 0.2) \times (0.5 \times 0.2) = 0.09$$



$$(1-0.2 \times 0.1) \times (0.2 \times 0.1) = 0.196$$

## Principles of Timing in Synchronous Systems

Explains the difference between Synchronous and Asynchronous systems

### # Introduction

Circuit design style can be classified into two major categories, which differ in how they indicate when the signals are in a stable 0 or 1 state (and are thus ready to be sampled reliably):

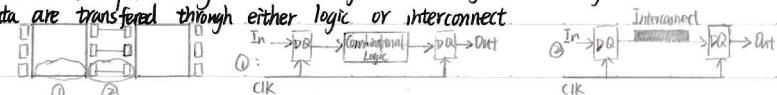
- Synchronous methodologies globally distribute a timing signal (a clock) to all parts of the circuit. Transaction Transitions (rising and/or falling depending on the design) on this clock line indicate moments at which the data signals are stable
- Asynchronous methodologies are fundamentally different; there is no common and discrete time. Instead the circuits use handshaking between their components in order to perform the necessary synchronization, communication, and sequencing of operations

Asynchronous systems don't have clock generator

### # Synchronous Systems (Logic vs. Interconnect)

In synchronous system all memory elements in the system are simultaneously updated using a globally distributed periodic synchronization signal (i.e., a global clock signal)

Data are transferred through either logic or interconnect



This functionality is ensured by strict timing constraints applied to the design. These constraints depend on the type of

Explains the timing constraints for synchronous systems (setup time, hold time)

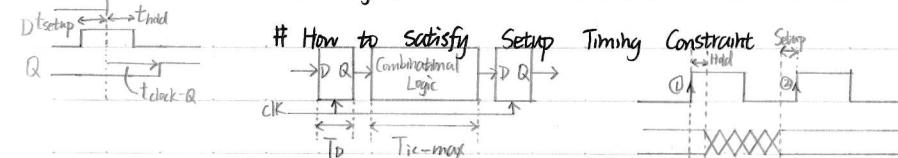
### # Flip-flop Based Synchronous Design

$Q$  becomes  $D$  after clock edge

Timing Constraint: The data applied to the  $D$  flip-flop should be stable for at least  $T_{\text{setup}}$  and should stay stable for at least  $T_{\text{hold}}$

Setup Time ( $T_{\text{setup}}$ ): is the minimum amount of time the data should be stable before the application of the clock signal

Hold Time ( $T_{\text{hold}}$ ): is the minimum amount of time the data should be stable after the application of the clock signal



Setup Time Constraint (for ideal clock):

$$T_{\text{clk-min}} \geq T_{\text{D}} + T_{\text{c-max}} + T_{\text{setup}}$$

Holding Constraint

$$T_{\text{clk-min}} \leq T_{\text{D}} + T_{\text{c-min}}$$

$T_{\text{clk-min}}$ : the minimum clock period

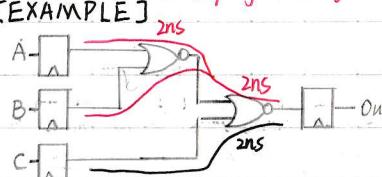
$T_{\text{D}}$ : the delay of the flip-flop

$T_{\text{c-max}}$ : the critical path delay

$T_{\text{clk-min}}$ : the minimum IC delay

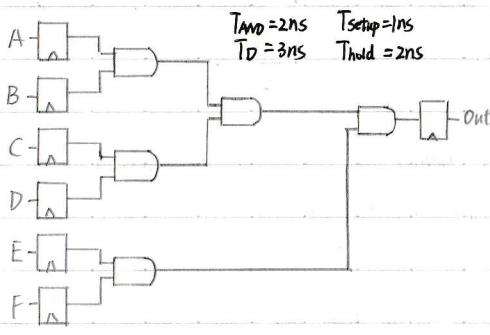
(contamination delay)

\* Calculate the maximum frequency of a digital synchronous circuit



$$\begin{aligned} \textcircled{1} \quad & T_{\text{clk}} \geq T_p + T_{\text{ic-max}} + T_{\text{setup}} \\ & T_{\text{clk}} \geq 3 + 4\text{ns} + 1\text{ns} \\ & T_{\text{clk}} \geq 8\text{ns} \\ & f_{\text{max}} = 1/T_{\text{clk}} = 125\text{MHz} \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad & T_{\text{hold}} \leq T_p + T_{\text{ic-min}} \\ & 2\text{ns} \leq 3\text{ns} + 2\text{ns} \\ & 2\text{ns} \leq 5\text{ns} \end{aligned}$$



$$\textcircled{1} \quad T_{\text{ic-max}} = 2\text{ns} + 2\text{ns} + 2\text{ns} = 6\text{ns}$$

$$\textcircled{2} \quad T_{\text{ic-min}} = 2\text{ns} + 2\text{ns} = 4\text{ns}$$

$$\textcircled{3} \quad T_{\text{clk}} \leq T_p + T_{\text{ic-max}} + T_{\text{setup}} = 3\text{ns} + 6\text{ns} + 1\text{ns} = 10\text{ns}$$

$$f_{\text{clk}} = 1/T_{\text{clk}} = 100\text{MHz}$$

$$\textcircled{4} \quad T_{\text{hold}} \leq T_p + T_{\text{ic-min}} = 3\text{ns} + 4\text{ns} = 7\text{ns}$$

$$\textcircled{5} \quad T_{\text{and}} = 0.5\text{ns} \quad T_p = 0.5\text{ns}$$

$$T_{\text{ic-min}} = 0.5\text{ns} + 0.5\text{ns} = 1\text{ns}$$

$$T_{\text{hold}} = 2\text{ns} \leq T_p + T_{\text{ic-min}} = 1.5\text{ns} \quad X$$

\* Perform timing analysis that includes interconnect effects

## # The Role of Interconnect

Buses are implemented as long metal lines on a silicon wafer transmitting data using electromagnetic waves (finite speed limit)

In addition to data transfer Metal lines (interconnect) are also used power distribution networks and for clock distribution network in synchronous design.

Interconnect networks have become increasingly complicated with greater Integration

## # Computation Vs. communication



### Operation

32 ALU Operation

32 Register Read

Read 32b from 8K RAM

Transfer 32b across chip (10mm)

Transfer 32b across chip (20mm)

### Delay (0.43nm)

650ps

125ps

32780

1400

2200

2800

4600

### Delay (0.15km)

250ps

500ps

1000

2000

3500

6000

2:1 global on-chip comm to operation delay

## # The impact of interconnect in Modern Technologies

**Limit Performance:** Signals cannot travel across the entire die within a global clock cycle.

For example, 6-to clock cycles are needed to transfer data on a bus at 50 nm technology  
Unpredictability in signal propagation time has serious consequences for performance and correct functioning of synchronous digital circuits

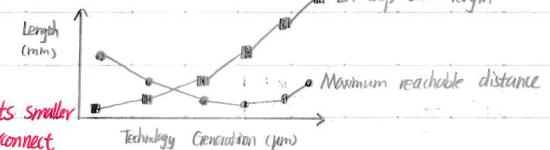
Consume up to 50% of chip power

Affect system reliability (crosstalk errors)

⇒ The figure shows that

Chip size increases while the technologies gets smaller

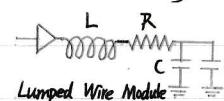
This leads to increased delays of global interconnect



## # Parasitic Modelling



## # Wire Modelling



Accurate and efficient models are essential to estimate the impact of interconnect on design budget (delay, power, reliability) at early stage of the design.  
Most design tools use RC models.

## # Distributed Models vs. Lumped Models

Thin on-chip are often modeled by distributed RC lines rather than lumped model

Three standard approximations are L-model, T-model, and Π-model, so named because of their shape.

The Π model is the most accurate, three segments are sufficient around for 97% accuracy level.



Lumped Model



Distributed Model



L-model



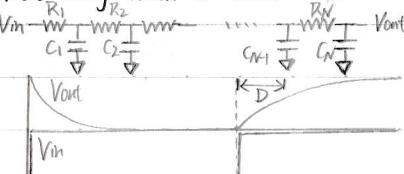
Π model



T-model

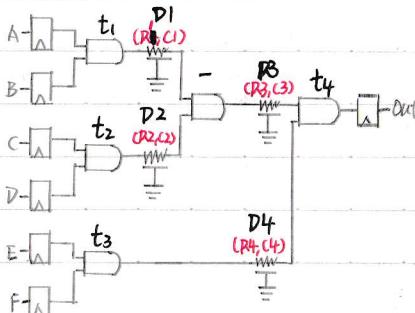
## # Wire delay modelling

Elmore analyzed the distributed model and came up with figures for delay in 1948 before VLSI  
Elmore delay which is taken to be the time between the 50% crossings of the input and output



$$D_{\text{wire}} = 0.4 * R_{\text{wire}} * C_{\text{wire}}$$

## [Exercise]



The circuit in figure 1 is synchronous and uses positive edge triggered flip-flops.

The delay of each logic gate is 0.1ns.  
The delay of each flip-flop is 0.3ns.  
 $T_{\text{setup}} = 0.1\text{ns}$     $T_{\text{hold}} = 0.1\text{ns}$

$$f_f = 10^{-15} \text{ F}$$

$$R_1 = R_2 = 500\Omega, R_3 = 400\Omega, R_4 = 1500\Omega$$

$$C_1 = C_2 = 100\text{fF}, C_3 = 80\text{fF}, C_4 = 250\text{fF}$$

Calculate the maximum achievable clock frequency  $f_{\text{max}}$

$$D1 = D2 = 0.4 * 500 * 100 * 10^{-15} = 0.02\text{ns}$$

$$D3 = 0.4 * 400 * 80 * 10^{-15} = 0.0128\text{ns}$$

$$D4 = 0.4 * 1500 * 250 * 10^{-15} = 0.15\text{ns}$$

!!!  $D4 > T_{\text{hold}}$

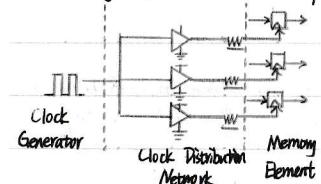
$$T_{\text{IC-max}} = t_3 + D4 + t_3 = 0.35\text{ns}$$

$$T_{\text{CLK-min}} \geq T_0 + T_{\text{IC-max}} + T_{\text{setup}} = 0.3 + 0.35 + 0.1 = 0.75\text{ns}$$

$$f_{\text{max}} = 1/T_{\text{CLK-min}} = 1.33\text{ GHz}$$

## ★ Explain the principles of clock generation.

## # Clocking System Main Components



**Clock Generation:** Usually a system on a chip receives one or more primary clock signals from an external clock chip and, in turn, generates necessary derivatives for its internal use

**Clock Distribution Network:** The primary design goal in clock distribution networks is to ensure that a clock signal arrives at every register within the entire synchronous system at precisely the same time (i.e. reduce clock skew and jitter)

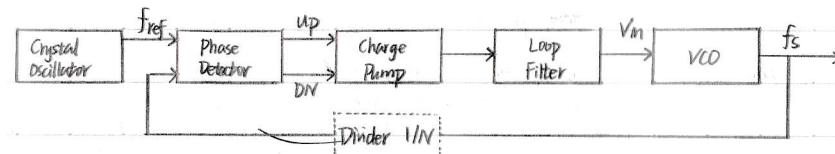
**Memory Elements:** Store the data values of the output of combinational logic blocks

## # Clock Generation : PLL Basic Operation

Phase-Locked Loop is a control system that generates a signal that has a fixed relation to the phase of a "reference" signal. A phase-locked loop circuit responds to both the frequency and the phase of the input signals, automatically raising or lowering the frequency of a controlled oscillator until it is matched to the reference in both frequency and phase

Crystal  $f_{\text{ref}}$   $\rightarrow$  PLL  $\rightarrow f_s$

## # Clock Generation : PLL Basic Operation



Phase-locked loop (PLL) is typically used to generate the high frequency required by SoC. A PLL takes an external low-frequency reference crystal frequency signal and multiplies its frequency by a rational number  $N$ .

The reference clock is typically generated off-chip from an accurate crystal reference. The reference clock is typically compared to a divided version of the system clock (i.e. the divided clock) using a phase detector that compares the phase difference between the signals and produces an Up or Down signal when the local clock lags or leads the reference signal. It detects which of the two input signals arrives earlier and produces an appropriate output signal.

Next, the Up and Down signals are fed into a charge pump, which translates the digital encoded control information into an analogue voltage. An Up signal increases the value of the control voltage and speeds up the VCO, which causes the local signal to catch up with the reference clock. A Down signal, on the other hand, slows down the oscillator and eliminates the phase lead of the local clock.

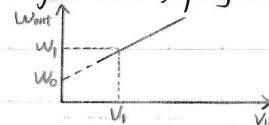
The Loop filter removes the high frequency components from the VCO control voltage and smooths out its response

## # Voltage Controlled Oscillator

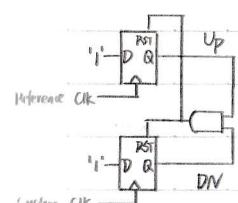
An ideal VCO generates a periodic signal whose frequency is a linear function of its input voltage

$$W = W_0 + k_{\text{vco}} \cdot V_{\text{in}}$$

↗ gain of the vco  
 ↘ fixed frequency offset

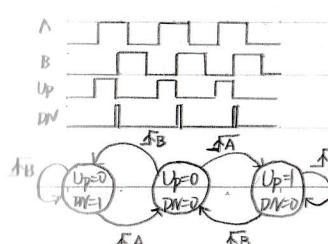


## # Phase Frequency Detector



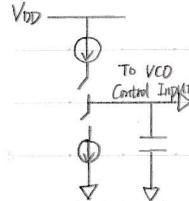
**Function:** It determines the relative phase and frequency difference between two incoming signals and outputs a signal that is proportional to this difference

**Operation:** When input A leads B, the Up output is asserted on the rising edge of input A. The Up signal remains in this state until a low-to-high transition occurs on input B. At that time, the DN output is asserted, causing both flip-flops to reset through the asynchronous reset signal. The duration of the small pulse on the DN output is equal to the delay through the AND gate and register reset delay. The pulse width of the Up pulse is equal to the phase error between the two signals. The roles are reversed for the case when input B leads A. The circuit also acts as a frequency error detector. For the case when A is at a higher frequency than B, the PFD generates a lot more Up pulses while the DN pulses average close to zero. The opposite is true for the case when B has a frequency larger than A — many more pulses are generated on the DN output than the Up output.



## # Charge Pump

**Function:** It converts the UP/DN signal pulses into an Analog Voltage that controls the VCO.



**Operation:** A pulse on the Up signal adds a charge packet to the storage capacitor proportional to the size of the Up pulse, and a pulse on the DN signal removes a charge packet proportional to the DN pulse. If the width of the Up pulse is larger than the DN pulse, then there is a net increase in the control voltage. This effectively increases the frequency of the VCO.

## # Loop Filter



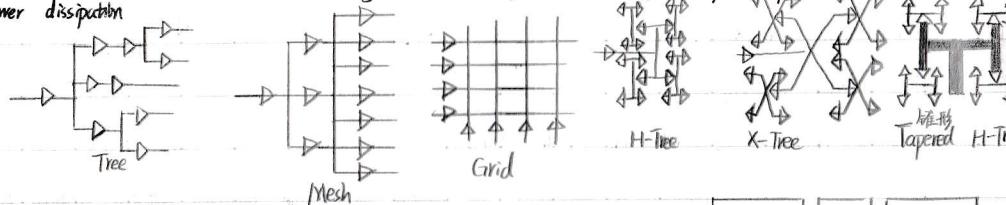
**Function:** Low-pass filter that removes the high-frequency components from the VCO control voltage and smooths out its response, which results in a reduction of the jitter (jitter is the temporal variation in the signal period). It is usually made of passive components.

## # Clock Distribution Networks (CDN)

**Design Requirement:** To distribute a tightly controlled clock signal to each synchronous register on a large non-redundant hierarchically structured integrated circuit

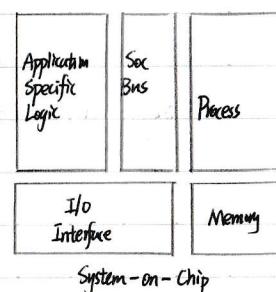
**Solutions:** Include the use of asymmetric structures such as buffered tree or symmetric structure such as H-trees

**Main issues to consider when making a CDN design choice** is system speed, physical die area, and power dissipation



**# Symmetric CDN structure** such as H-tree and X-tree are difficult to implement in irregular VLSI-based systems due to routing constraints and different fan-out requirements.

In these types of systems, buffered tree topologies integrated with structured custom design methodologies should be used



## # Clock Tree Synthesis for ASICs

Modern back-end tools include clock tree synthesis

- Creates balanced RC-trees
- Uses special clock buffer standard cells
- Can add clock shielding
- Can exploit useful clock skew

Consider a zero-skew clock routing network that was constructed based on estimated routing parasitic

• The primary design goal in clock distribution networks has been to ensure that a clock signal arrives at every register within the entire synchronous system at precisely the same time (i.e. reduce clock skew and jitter)

■ Reduce Power: In many high-speed process processors, a majority of the dynamic power is dissipated in the clock network. One option to reduce dynamic power is to use clock gating (shutting down disabling the clk) units

■ Reduce Area

■ Clock Skew Rate: To maintain signal integrity minimum skew rates are required. However, aggressive reduction of skew may lead to increase power consumption and signal reflections.

\* Optimize the power consumption for the clocking systems

## # Power reduction techniques

Dynamic Power Consumption of a clocking scheme can be given as:

$$P_{clk} = 8 \cdot f \cdot C \cdot V_s^2$$

$f$ : Switching activity of the clock distribution network

$f$ : Clock Frequency

$V_s$ : The swing voltage of the clock signal

$$C = C_{CDN} + C_{ME}$$

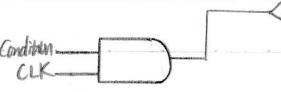
$C_{CDN}$ : the total capacitance of the clock distribution interconnect including all buffers

$C_{ME}$ : the sum of the capacitances of all memory elements used in the clocking system

(e.g. flip-flops or latches)

## # Power reduction techniques : clock gating

Principle: CG reduces switching activity by **Not clocking** some parts of the system if their outputs are not needed. Most clock-gating is done at the Register Transfer Level.



Combinational Clock Gating: reduces power by disabling the clock on registers when the output is not changing

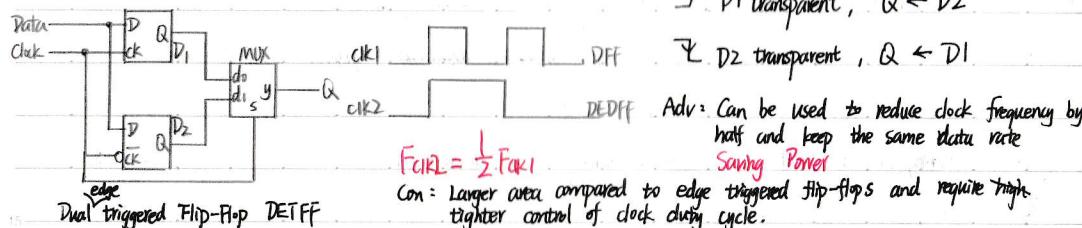
System-level clock-gating stops the clock for an entire block, effectively disabling all functionality. On the contrary, combinational and sequential selectively suspend clocking while the block continues to produce an output. Condition for disabling a certain clock is identified by the designer.

## # Power reduction techniques : Dual Edge Triggered Flip-Flop

Principles: Use the special types of storage elements such that the clock frequency can be reduced by half by and keep the same data rate

$$\text{FDE} = \frac{1}{2} \text{FDT}$$

P1 transparent, Q  $\leftarrow$  D2  
P2 transparent, Q  $\leftarrow$  D1

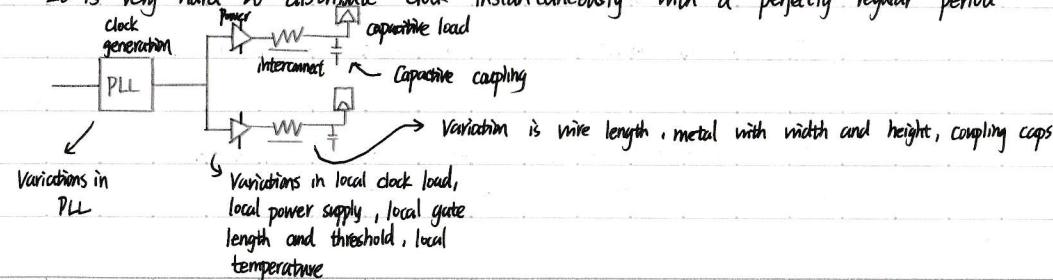


Perform timing analysis for Synchronous Systems with skew and jitter

## # FF based Synchronous Design

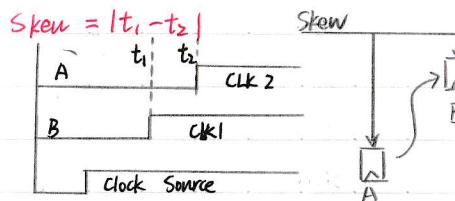
### # Clock Signal is not ideal

It is very hard to distribute clock instantaneously with a perfectly regular period



## # Clock Skew : Spatial Clock Variation

Clock skew: static (constant over time) point-to-point variation of clock arrival time to FFs. Clock skew between two flip-flops is computed as the difference in arrival times of clock signal



## # Causes of Clock Skew

Designed Variations: mismatch in buffer load sizes, interconnect lengths

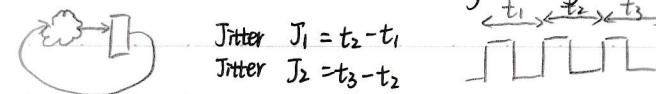
Process Variations: process speed spread across die yielding different L<sub>off</sub>, T<sub>ox</sub>, etc. values

Temperature Gradients: change MOSFET performance across die

IR voltage drop in power supply: changes MOSFET performance across die.

## # Clock Jitter : Temporal Clock Variation

Clock Jitter is a variation in clock edge timing between clock cycle i.e. Cycle-to-Cycle variation



## # Sources of Jitter

Variations in PLL oscillation frequency: This is an analogue circuit and highly sensitive to intrinsic device noise and power supply variations. A major problem is the coupling from the surrounding noisy digital circuitry through the substrate.

High Frequency Power Supply Variations: This is caused by instantaneous IR drops along the power signal grid due to fluctuation (leakage) in switching activity can affect the delay of buffers in the clock distribution network which lead to temporal variation in clock cycle.

Variation In Capacitive Load: Due to coupling between the clock lines and adjacent signal wires. The exact coupling to the clock network is not fixed from cycle-to-cycle.

## # Skew vs Jitter

① Uncertainty in the clock generation circuit

Jitter

② Process variation in devices

Skew

③ Interconnect variation ✓ design variation

Skew

④ Power supply noise

Jitter

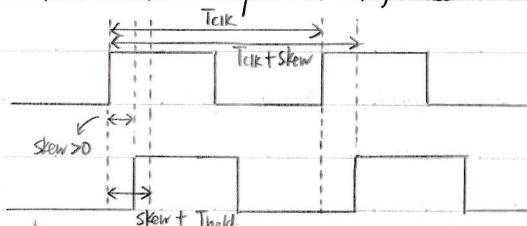
⑤ Data Dependent Load Capacitance

Jitter

⑥ Static Temperature Gradient

Skew

## # Positive Skew Impact on Performance



$$T_{clk-min} + \text{skew} \geq T_D + T_{IC-max} + T_{setup}$$

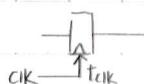
$$T_{hold} \leq T_D + T_{IC-min} - \text{skew}$$

Skew > 0 : Improves performance, but makes Thold harder to meet

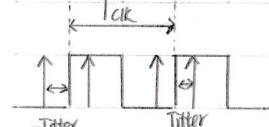
If Thold is not met (race conditions), the circuit malfunctions independent of the clock speed.

## # Jitter Impact on Performance

Jitter causes T to vary on



a cycle-by cycle basis



$$T_{clk} \geq T_D + T_{IC-max} + T_{setup} + 2 \cdot \text{Jitter}$$

$$T_{hold} \leq T_D + T_{IC-min} - 2 \cdot \text{Jitter}$$

Jitter directly reduces the performance of a sequential circuit.

## # Combined Impact

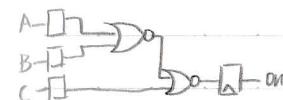
$$T_{clk} \geq T_D + T_{IC-max} + T_{setup} - \text{skew} + 2 \cdot \text{Jitter}$$

$$T_{hold} \leq T_{hold} + T_{IC-min} - \text{skew} - 2 \cdot \text{Jitter}$$

Skew > 0 with jitter : Degrades performance, and makes Thold even harder to meet

Campus

## [EXAMPLE]



$$\begin{aligned} T_D &= 3 \text{ ns} \\ T_{HOLD} &= 2 \text{ ns} \\ T_{SETUP} &= 1 \text{ ns} \\ T_{HOLD} &= 2 \text{ ns} \\ \text{(c) skew} &= 1 \text{ ns} \\ \text{Jitter} &= 1 \text{ ns} \end{aligned}$$

$$T_{clk} \geq T_{D-min} + T_{setup} + 2 \cdot \text{Jitter} - \text{skew}$$

$$= 3 + 4 + 1 + 2 - 1 = 9$$

$$f_{max} = \frac{1}{T_{clk}} = 111.11 \text{ MHz}$$

$$T_{hold} = 2 \leq T_D + T_{IC-min} - 2 \cdot \text{Jitter} - \text{skew} = 3 + 2 - 2 - 1 = 2$$

## \* Differentiate between Critical path and False Path

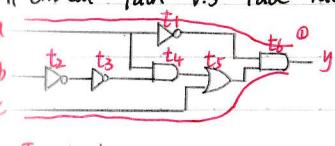
# Timing Analysis (There are two main types of timing verification)

**Dynamic Timing Analysis :** Circuit timing is verified by means of exhaustive simulation of the circuit using actual gate delay. The test vectors to the inputs of a circuit, then the changes in all gates and followed through to the outputs in order to find the worst case time. It can take an extraordinary amount of time to run.

**Static Timing Analysis (STA) :** It does not need simulations, the delay of all paths are computed using gate delay information to check for any setup or hold time violations. This method runs quickly due to the use of simplified delay models, and because of its limited consideration of the effects of logical interactions between signals. It is used extensively in all EDA tools.

STA suffer from false path problem which may lead to 保守的 conservative estimation of circuit performance

## # Critical Path v.s False Path



$$\textcircled{a} \quad T_{IC} = t_1 + t_6$$

$$\textcircled{a} \quad T_{IC-min} = t_5 + t_6$$

The clock speed is determined by the slowest feasible path between registers in the design referred to as **Critical path**

A **false path** is a path which cannot be exercised due to Boolean gate conditions or it is not relevant to the circuit path. False paths cause pessimistic timing constraints 悲观的

Because of these effects, static timing analysis might be overly conservative and predict a delay that is greater than you will experience in practice.

$$a=0 \Rightarrow b=x$$

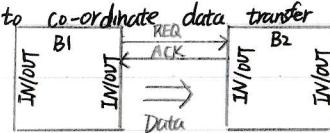
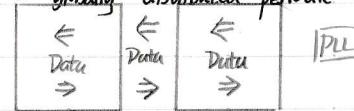
$$a=1 \Rightarrow b,c,x$$

## Basic Principle of Asynchronous Circuit Design

\* To explain the basic operation principles of asynchronous circuits

### II Introduction (Synchronous)

Synchronous systems uses a globally distributed periodic signal to co-ordinate data transfer among system's parts.



Asynchronous Systems use handshaking protocols to co-ordinate data transfer among system's parts

### II Key Design Difference

#### Synchronous logic design:

- Proceeds without taking timing correctness (hazards, signal ack-ing etc.) into account
- Combinational logic and memory latches (registers) are built separately
- Static timing analysis of combinational logic is sufficient to determine the Max Delay (clock period)

#### Asynchronous logic design

- Must ensure hazard-freedom, signal ack-ing, local timing constraints
- Combinational logic and memory latches (registers) are often mixed in "complex gates"
- Dynamic timing analysis of logic is needed to determine relative delays between paths

To avoid complex issues, circuits may be built as Delay-insensitive and/or Speed-independent

### II Verification and Testing Difference

#### Synchronous logic verification and testing:

- Only functional correctness aspect is verified and tested
- Testing can be done with stand ATE and at low speed

#### Asynchronous logic verification and testing:

- In addition to functional correctness, temporal aspect is crucial: e.g. causality and order
- Testing must cover faults in complex gates (logic + memory) and must proceed at normal operation rate

deadlock - freedom  
死锁自由

### II Motivation for Asynchronous

#### Pro

##### Low power consumption

Less emission of electromagnetic noise (No peak current around clock edges)  
Robustness (健強性) toward variation in supply voltage, temperature and fabrication process  
Modularity for system-on-chip design (Plug and play interconnectivity)

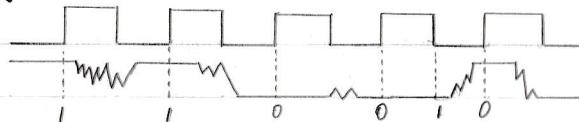
##### Average-case performance

No clock distribution and clock skew problems

#### Con

Difficult to design (Hazards, synchronization)  
Complex timing analysis (Difficult to estimate performance)  
Difficult to test (No way to stop the clock)

## # Synchronous communication

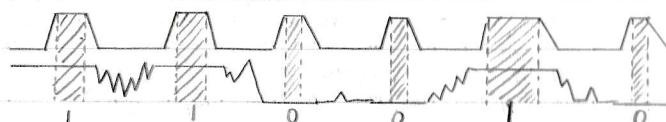


Clock edges determine the time instants where data must be sampled.

Data wires may glitch between clock edges (set-up/hold times must be satisfied)

Data are transmitted at a fixed rate (clock frequency)

## # Bundled data communication



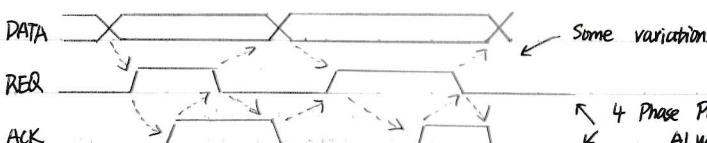
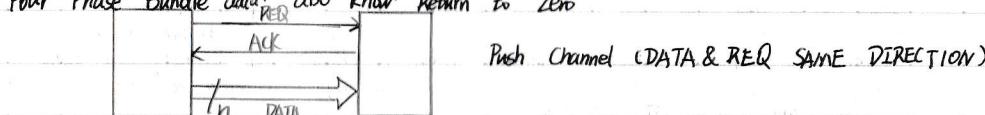
Validity signal • Similar to an aperiodic local clock.

$n$ -bit data communication requires  $n+1$  wires

Data wires may glitch when no valid

## # Bundled Data Protocol

Four Phase Bundle data (also known as Return to Zero)



## # Timing Assumption of Bundled

Each data line is a single wire • "Bundled data also known as single rail"

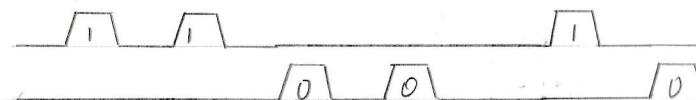
On sender side, time (DATA) < time (REQ)

This order should be preserved on receiving end: valid (DATA) < REQ [data should be valid precedes REQ]

Non-trivial: inter-line skew must be taken care of and hidden

- Placement and routing
- Safety margins at sending -data end
- Buffer insertion

## # Dual rail communication



Two wires with L (low) and H (high) per bit

• "LL" = "spacer", "LH" = "0", "HL" = "1".

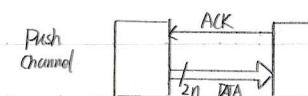
$n$ -bit data communication requires  $2n$  wires

Each bit is soft-timed

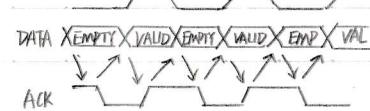
Other delay-insensitive codes exist (e.g. k-of-n) and event-based signalling (choice criteria: pin and power efficiency)

## # 4-phase dual rail protocol

Each data bit encoded into 2 wires



VALUE	dt	df
EMPTY	0	0
VALID"0"	0	1
VALID"1"	1	0
Not used	1	1



Delay Insensitive (DI): Each bit can propagate at own speed.

4 phase at higher level (than signals):

1. Sender sends valid word (V)
2. Receiver sets ACK  $\uparrow$
3. Sender sends empty word (E) (removes the data)
4. Receiver sets ACK  $\downarrow$

Each acknowledge is acknowledged / indicated

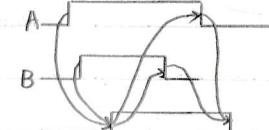
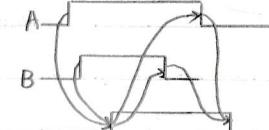
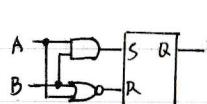
Problems: Glitches, Hazards.

\* To describe the structure and the functionality of the Muller C-element

### # Event Logic - The Muller - C Element

A	B	F <sub>M1</sub>
0	0	0
0	1	F <sub>N</sub>
1	0	F <sub>N</sub>
1	1	1

Schematic



### # Acknowledgement / Indication

A gate / circuit acknowledges its input if, for every input change, there is an output change (e.g. wave)

Non-indicating example : AND gate { Acknowledge :  $\{01, 10, 00\} \rightarrow 11$  }  
 AND  
 Does not acknowledge :  $00 \rightarrow \{01, 10\}$

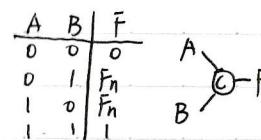
### # Acknowledgement / Indication

#### C-element

Acknowledges all ones :  $\{00, 10, 01\} \rightarrow 11$

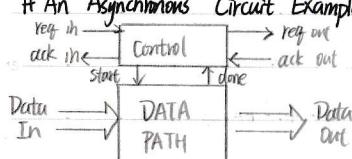
Acknowledges all zeros :  $\{10, 01, 11\} \rightarrow 00$

A	B	F
0	0	0
0	1	F <sub>N</sub>
1	0	F <sub>N</sub>
1	1	1



\* To design asynchronous circuit for simple combinational functions based on dual rail protocol

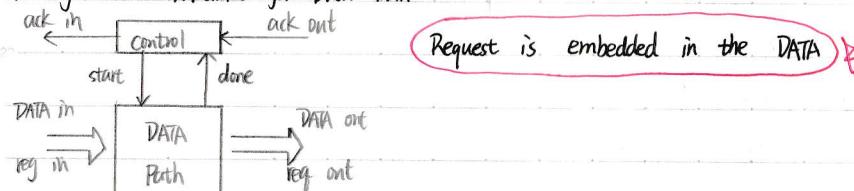
### # An Asynchronous Circuit Example



Two essential components of asynchronous circuits:

- ① Control which generates the handshake signal
- ② Data path computes data and generates the completion detection signal

### # Asynchronous Modules for Dual Rail



Request is embedded in the DATA

### # Dual Rail Circuit : Datapath Design

The datapath performs the computation and generates the completion detection signal (i.e. done). The datapath also generates the request signals. There are many techniques to design asynchronous dual rail logic. We will study one method called Delay insensitive Min-term Synthesis

#### Dual Rail Logic Design

##### Delay Insensitive Min-term Synthesis (DIMTS)

This technique resembles the traditional sum of products approach, but there are also few important differences :

① The minterms are formed using C-elements (instead of AND gates)

② Reduction of the Boolean Equations by combining minterms into simpler terms is (in general) not allowed. Together, these assumes requirements assure that our combinational circuits do not produce any valid output signals until all input signals are valid, and that none of the output signals change back to the empty value until all inputs are empty

#### DIMTS Design Example

How to build a DIMTS circuit for a logic function

① Write a truth table for the function

② Replace each binary value with its dual rail equivalence

At Af Bt Bf OUTt OUTf

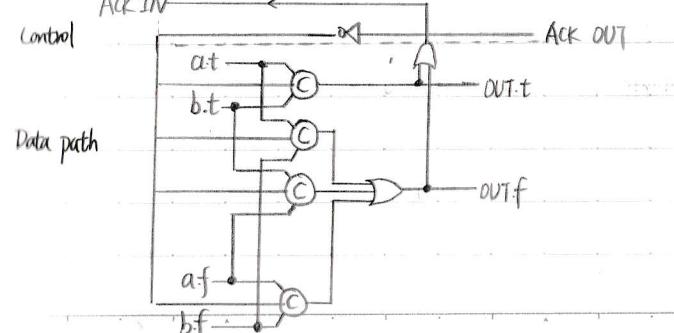
0	1	0	1	0	1
0	1	1	0	0	1
1	0	0	1	0	1
1	1	1	0	1	0

OR gates

③ Build your functions based on the minterms using C-elements (instead of AND-gates) and

OUTt = At.Bt

OUTf = At.Bf + Af.Bt + Af.Bf



## # Dual Rail Logic Design : DIMS Design Example : Adder

## ① Truth Table

Cin	A	B	S	Cont
0	0	0	0	0
0	0	1	1	0
0	1	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1
1	1	0	1	1
1	1	1	1	0

## ② Dual Rail Equivalence

Cin-t	Cin-f	A-t	A-f	B-t	B-f	S-t	S-f	Cout-t	Cout-f
0	1	0	1	0	1	0	1	0	1
0	0	0	0	1	1	0	1	0	1
0	1	1	0	0	0	1	0	0	1
0	1	0	1	1	0	1	0	1	0
1	0	0	1	0	0	1	0	0	1
1	0	1	0	1	0	1	0	1	0
1	1	0	0	1	0	0	1	0	1
1	1	1	1	1	0	1	0	1	0

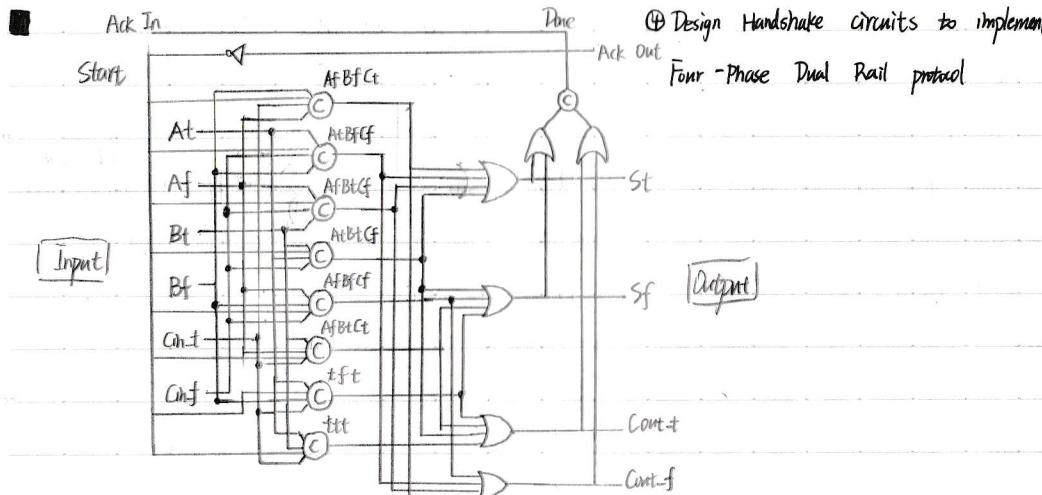
## ③ Build functions

$$St = Af \cdot Bt \cdot Cin-f + At \cdot Bf \cdot Cin-f + Af \cdot Bf \cdot Cin-t + At \cdot Bt \cdot Cin-t$$

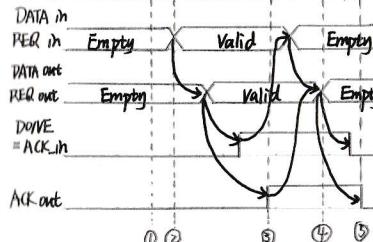
$$Sf = Af \cdot Bf \cdot Cin-f + At \cdot Bt \cdot Cin-f + Af \cdot Bt \cdot Cin-t + At \cdot Bf \cdot Cin-t$$

$$Cont-t = At \cdot Bt \cdot Cin-f + At \cdot Bf \cdot Cin-t + Af \cdot Bt \cdot Cin-t + Af \cdot Bf \cdot Cin-t$$

$$Cont-f = Af \cdot Bf \cdot Cin-f + At \cdot Bf \cdot Cin-f + Af \cdot Bt \cdot Cin-f + Af \cdot Bf \cdot Cin-t$$



## # Four Phase, Dual Rail Protocol



## ① At the Beginning

Both inputs and outputs are in the empty state

$$\begin{cases} \text{Input} = \text{Empty} & \text{Ack\_out} = 0 \\ \text{Start} = 1 & \\ \text{Output} = \text{Empty} & \text{Done} = 0 \\ \text{Ack\_in} = 0 & \text{Ack\_in} = 0 \end{cases}$$

## ② Inputs become valid

④ When the next stage of the pipeline receives the valid data, it asserts Ack\_out , Ack\_out goes high and start goes low

The circuit will wait until all the inputs transition to the empty state

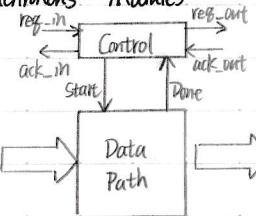
After all the input transitioning to the empty state and Ack\_out is asserted, the outputs will transition to the empty state

## ⑤ After all outputs transition to an empty state Ack\_in will be de-asserted

The next stage of the pipeline will de-assert Ack\_out when it receives the empty data. This takes us back to the initial state.

## Logic Design for Asynchronous Control Circuit and Interfaces

### # Asynchronous Modules



Two essential components of asynchronous circuits :

1. Control which generates the handshake signal
2. Data path computes data and generates the completion detection signal.

### # Delay Modules

① Bounded delays (BD) : Rel realistic for gates and wires

- Technology mapping is easy, verification is difficult

② Speed independent (SI) : Unbounded (pessimistic) delays for gates and "negligible" (optimistic) delays for wires

- Technology mapping is more difficult, verification is easy

③ Delay insensitive (DI) : Unbounded (pessimistic) delays for gates and wires.

- DI class (built out of basic gates) is almost empty

### # Asynchronous Control Circuits

There are numerous methods, tools and different assumptions (e.g. delay models)

In this course :

Assumption : Speed independent control circuits

Method : Signal transition graph (A special type of Petri net)

Tool : Workcraft (mostly developed at Newcastle University)

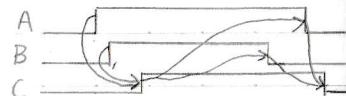
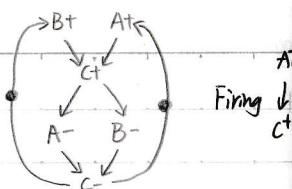
# Transfer timing specifications into a signal transition graph

## # Signal Transition Graph

Vertices represent signal transitions

Arcs represent caused relations between transitions

Marking is the assignment of tokens to places and represents the state of the circuit.



## # Token preservation

Tokens do not disappear

Tokens do not appear (from nowhere)

One token does not overtake another

A signal transition graph with  $n$  inputs can generate  $m$  outputs

- wait for  $n$  tokens on inputs
- generate  $m$  tokens on outputs



## # Signal Transition Graph: Concurrency & Choice

### Concurrency

$a^+$  and  $b^-$  are concurrent with each other.  
Either  $a$ - fires before  $b$ - or  $b$ - fires before  $a$ - or  $a$ - and  $b$ - fires simultaneously

### Choice



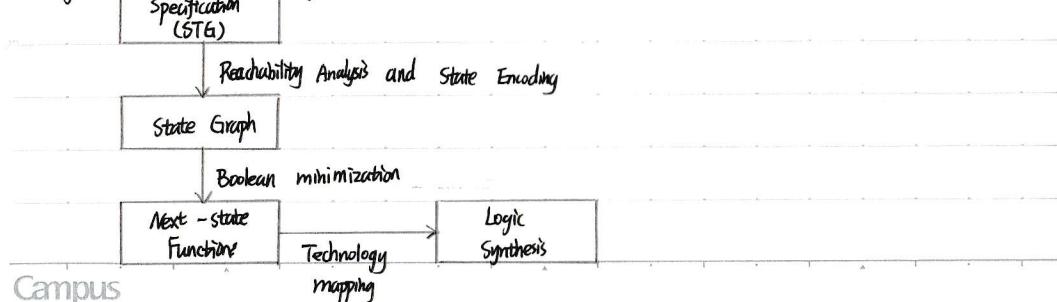
Either  $a$ - fires

or  $b$ - fire  
but not both

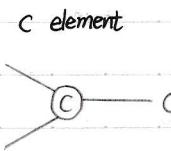
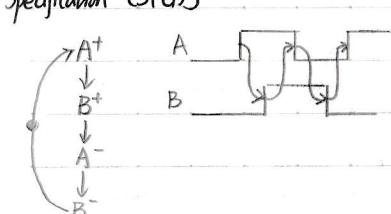
choice made by environment

# Digital logic circuits for asynchronous controllers

## # Asynchronous Circuits' Design Flow



## [Specification (STG)]



Buffer

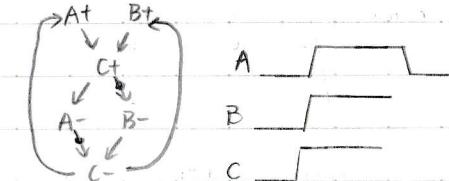
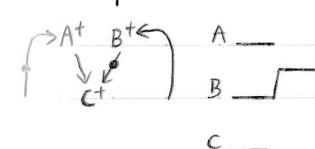
C element

## [State Graph]

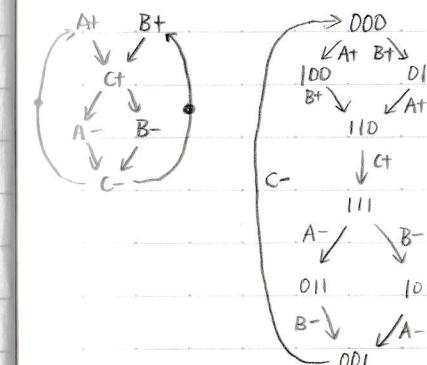
### Reachability Analysis

→ The aim of this analysis is to find all reachable states for a given STG by simulating token flow

### State Graph : Transition

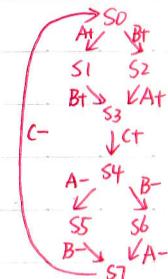


### State Graph : State Coding



### State Graph : Reachability Analysis

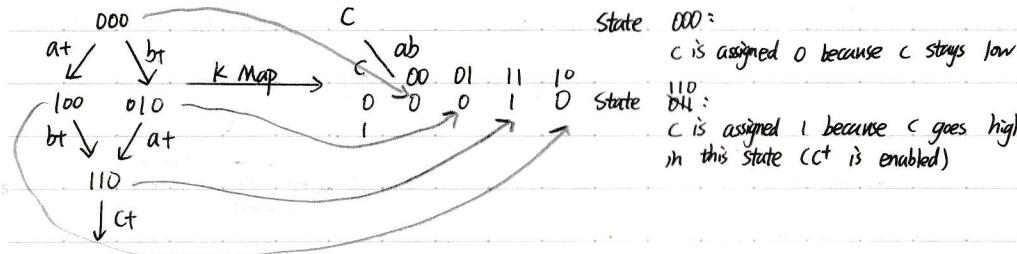
For more complex examples reachability analysis is performed before binary state coding



State Graph

Reachability Graph

## [Next State Functions]

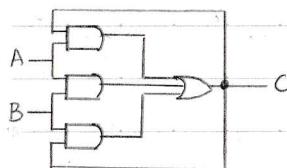


## [Logic Synthesis]

Using Complex gates

	ab	00	01	11	10	R: Rising	F: Falling
0	0	0	0	1	0	1	0
1	F	1	1	0	1	0	1

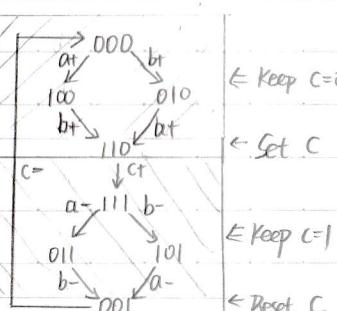
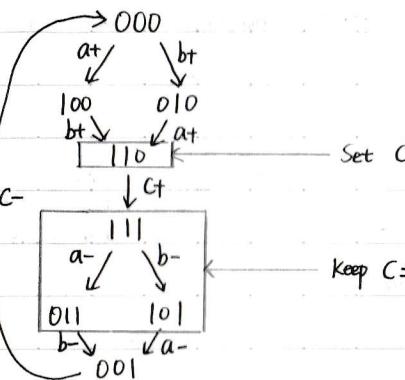
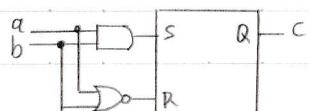
$$C = ab + ac + bc$$



HAZARDOUS Material

Using SR latches

	ab	00	01	11	10	R	S
0	0	0	0	1	0	1	1
1	F	1	1	0	1	1	1



## II Logic Synthesis : using SR Latch

The set function (cover) for a variable  $C$ :

1. Must contain all states in the SET regions
2. May contain states from the KEEP  $C=1$  regions
3. May contain states not reachable by the circuits

The reset function (cover) for a variable  $C$ :

1. Must contain all states in the RESET regions
2. May contain states from the KEEP  $C=0$  regions
3. May contain states not reachable by the circuits

- ① Needs (mutually exclusive) SET,RESET regions
- ② May use  $KEEP\ 0$ ,  $KEEP\ 1$  regions
- ③ May use ~~unreach~~ unreachable states
- ④ Area/performance/power may be more or less than gate circuits
- ⑤ May use C elements instead of SR latches
- ⑥ SR latches enable implementation with standard libraries

## Explain synthesis conditions for speed independent circuit's specifications

### Synthesis Conditions

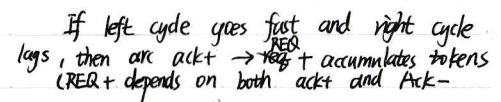
In a speed independent circuit, the delay of the gate does not affect the functionality of the circuit.

To guarantee a hazard-free speed independent implementation for an STG, the following properties must be satisfied:

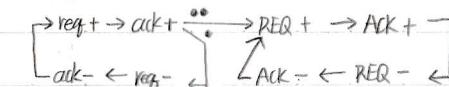
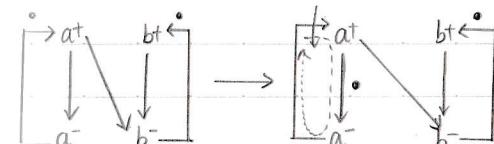
1. Boundedness 有界性
2. Consistency 一致性
3. Complete State Decoding
4. Output Persistence 持久性

**Safety**: • STG is bounded (safe) if no place or arc can ever contain more than one token  
• Unbounded is always caused by one-side dependency

EXAMPLE: the STG below is unbounded (not safe)



### Boundedness



Boundedness  $\leftrightarrow$  Finite reachability graph

左边比右边快，右边会同时有1/个token在运行  
无限累积

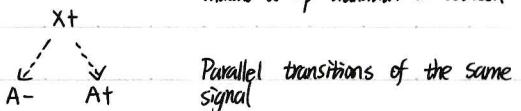
## # Consistent State Assignment (CSA)

The subset of an STG makes no sense

### INCONSISTENT

$A+ \dashrightarrow A+$  Two subsequent up-transitions without a down-transition in between

$A- \dashrightarrow A-$  Two subsequent down-transitions without an up-transition in between



Parallel transitions of the same signal



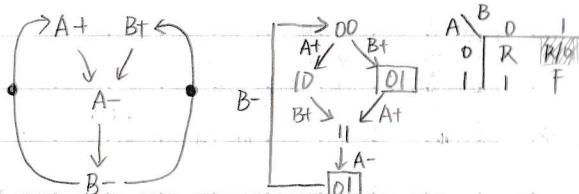
### CONSISTENT

$A+ \dashrightarrow A- \dashrightarrow A+$

Two up-transitions with a down-transition in between (and vice versa)

Boundedness + CSA  $\Leftrightarrow$  Finite Binary State Graph

## # Complete State Coding

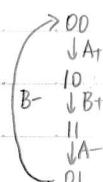


### CSC - conflict:

Two states have the same binary codes, and different set of enabled non-input signals

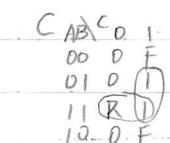
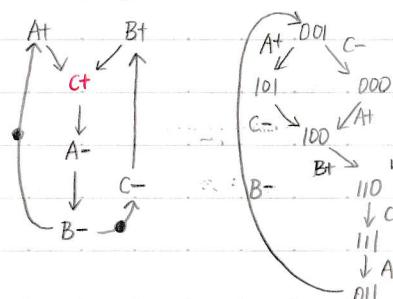
## Solving CSC by concurrency reduction

$A+ \rightarrow B+$   
↑  
Add ARC



No equivalent transformation - concurrency is reduced!

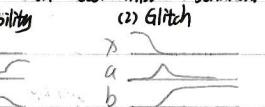
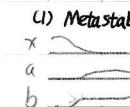
## Solving CSC by Signal Insertion



$$A = \bar{B} + \bar{C}$$

$$C = B(C+A)$$

Non-deterministic behaviour

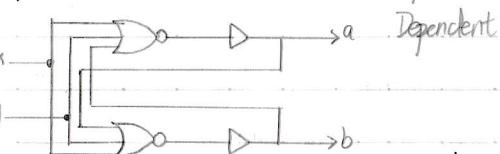
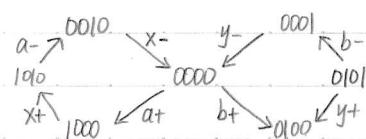


## # Output Persistency

Output disables output

Assume:  $a, b$  are outputs and  $x, y$  inputs

Depend on which buffer is stronger Speed



A state graph is said to be output-persistent if no event can be disabled by another event (unless they are both inputs).

This property determines whether or not the circuit is speed independent

速度

## # Violations of Implementability

### Violations

Boundedness

Consistency

CSC (reducible)

Persistency

### Meaning

Infinite Reachability  
No binary state graph

Change Spec

Change Spec

No Logic

Hazards

### How to correct

Change Spec

Change Spec

Add Signals or Change Spec

Change Spec

### Type of Implementation

N/A

N/A

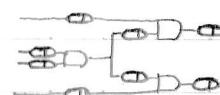
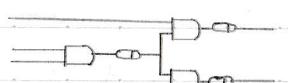
Same inputs and outputs

No speed independent

\* Verify functional correctness under different delay models.

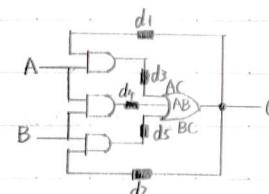
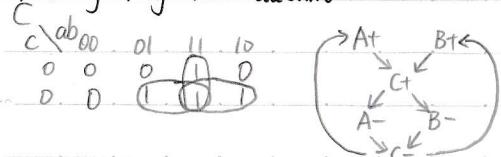
### # Gate vs wire delay models

Gate delay model: delays in gates, no delay in wires



Wire delay model: delays in gates and wires.

### # Taking Delay into account



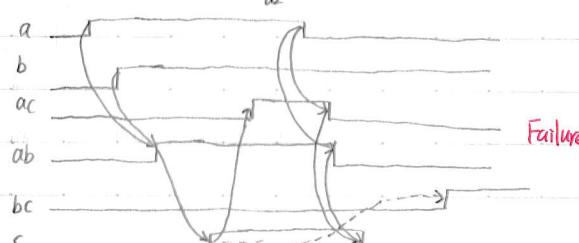
Delay assumption:

$$\begin{aligned}d_1 &= d_2 = 2 \text{ time units} \\d_3 &= d_4 = 1 \text{ time units} \\d_5 &= 10 \text{ time units}\end{aligned}$$

Gate: 1 time unit

(The logic inside ↗)

This design functions correctly for unbounded (complex) gate delay models  
This design is incorrect for unbounded wire delay models  
For bounded delay model, correctness depends on delay constraints



### # Summary

This synthesis method produces speed independent implementations that function correctly under the assumption of unbounded gate delay models.

Wire delay are ignored, in real system this can be achieved by dividing the system into smaller modules within which wire delay are ignored. These modules communicate using delay insensitive interfaces or interfaces satisfying certain timing constraints.

It is your duties as a designer to make sure that these assumptions are correct.

### # Asynchronous Design Tools

[1] Workcraft <http://www.workcraft.org/>

Developed by Newcastle University

Used for synthesis of asynchronous circuits from Petri net and Transition graph specification

[2] The Balsa Asynchronous Synthesis System

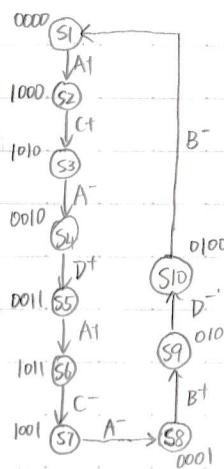
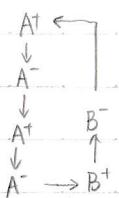
Developed over a number of years at the University of Manchester.

Balsa is built around the Handshake Circuits methodology and can generate gate level netlist from high-level descriptions in the Balsa language.

Both dual-rail (QDI) and signal-rail (bundled data) circuits can be generated.

Latest release supports Xilinx FPGA technology targets

<http://apt.cs.man.ac.uk/projects/tools/balsa/>



"[ ]" indicate a feedback to input "out"

Inorder = AB CD

Outorder: [B] [C] [D]

[B] = A' C' D

[C] = A' C + AD'

[D] = B' D + A' C

$$\begin{array}{c} \text{AB} \\ \diagdown \quad \diagup \\ \text{00} \quad \text{01} \quad \text{11} \quad \text{10} \\ \text{00} \quad \text{0} \quad \boxed{1} \quad \text{0} \quad \text{0} \\ \text{01} \quad \text{0} \quad \text{0} \quad \boxed{1} \quad \text{0} \\ \text{11} \quad \text{0} \quad \text{0} \quad \text{0} \quad \text{0} \\ \text{10} \quad \text{0} \quad \text{0} \quad \text{0} \quad \text{0} \end{array}$$

$$B = A' C' D$$

$$\begin{array}{c} \text{AB} \\ \diagdown \quad \diagup \\ \text{00} \quad \text{01} \quad \text{11} \quad \text{10} \\ \text{00} \quad \text{0} \quad \text{0} \quad \boxed{1} \quad \text{0} \\ \text{01} \quad \text{0} \quad \text{0} \quad \text{0} \quad \boxed{1} \\ \text{11} \quad \text{0} \quad \text{0} \quad \text{0} \quad \text{0} \\ \text{10} \quad \text{0} \quad \text{0} \quad \text{0} \quad \boxed{1} \end{array}$$

$$\checkmark C = A C + A D'$$

$$C = A' B' C + A B' D'$$

$$\begin{array}{c} \text{AB} \\ \diagdown \quad \diagup \\ \text{00} \quad \text{01} \quad \text{11} \quad \text{10} \\ \text{00} \quad \text{0} \quad \boxed{1} \quad \boxed{1} \quad \text{0} \\ \text{01} \quad \text{0} \quad \text{0} \quad \text{0} \quad \text{X} \\ \text{11} \quad \text{0} \quad \text{1} \quad \text{1} \quad \text{0} \end{array}$$

$$\checkmark D = B' D + A' C$$

$$D = B' D + A B' C$$

## Metastability and Synchronizers

# Describe the need for synchronization in modern system.

# Why Synchronization : 1. Sampling Asynchronous Signals

Before an asynchronous signal can be input to a synchronous module, it must first be synchronized with the local clock to ensure that it makes transitions only during allowed portions of the clock cycle. Such synchronization is required, for example, when sampling real-world input devices whose output can change at any time asynchronously. If such signals that were input directly into a synchronous system, they might cause a synchronization failure.

# Why Synchronization = 2. Crossing Clock Domains

Increasing complexity of modern SoC requires the use of many different IPs with different clocking requirements

Example of a smartphone: ① v2o IPs are getting usual

- ② Physical interface peripherals  $\sim$  200MHz  
 ③ Modern DSP processing  $\sim$  400MHz  
 ④ MAC processing  $\sim$  600MHz  
 ⑤ GP Host  $\sim$  600MHz  $\rightarrow$  1GHz  
 ⑥ Graphical & Video accelerators  $\sim$  500MHz  
 ⑦ Interconnect  $\sim$  400MHz

The use of a single clock for such systems is becoming less feasible, therefore modern systems have multiple clock domains.

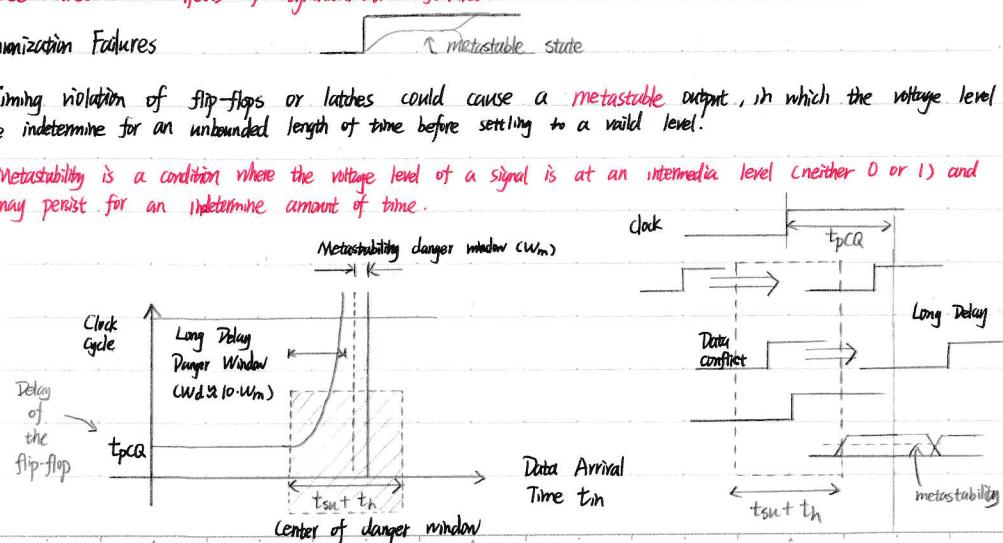
When transferring data between different clock domains, it must first be synchronized with the local clock of the destination to ensure that it makes transitions only during allowed portions of the clock cycle otherwise this may cause a synchronization failure.

Describe causes and effects of synchronization failures

## # Synchronization Failures

Timing violation of flip-flops or latches could cause a **metastable** output, in which the voltage level may be indeterminate for an unbounded length of time before settling to a valid level.

Metastability is a condition where the voltage level of a signal is at an intermediate level (neither 0 or 1) and which may persist for an indefinite amount of time.



## # Long Delay / Metastability windows

Both Long Delay and Metastability windows are bad.

We will consider the combined window in our calculation:  $W = W_d + W_m$  (either long delay or metastability)

If data arrives during the danger window  $W$  then we will get a synchronous failure

$W$  is Experimentally, less than 4 gate delays (F04)

## # How often does metastability happen?

Assuming that : ① Data arrives or changes uniformly over clock cycle  $T$

②  $W$  is the width danger zone

The probability of entering Metastability assuming a clock period is given by:

$$P = \frac{W}{T} = W \times F_c$$

Multiplying the probability of becoming metastable during a cycle that has a data event by the rate of data events yields the rate of entering into metastability (MR)

$$MR = W \times F_c \times F_d$$

## [EXAMPLE]

Calculate the rate of Metastability that is caused by moving one data item from IP1 to IP2, assuming the clock frequency of IP1 is 800MHz, the clock frequency of IP2 is 400MHz

The Data transfer rate is 200KHz  
The width of the danger zone for all flip-flops in the system is 132ps  $10^{-12}$

$$F_c = 400\text{MHz} \quad F_d = 200\text{kHz} \quad W = 132\text{ps}$$

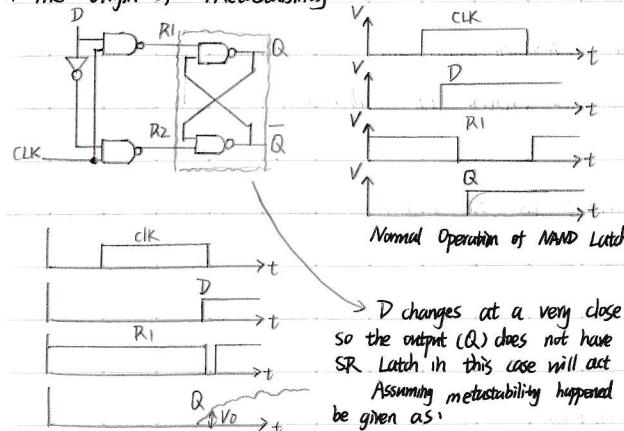
$$MR = W \cdot F_c \cdot F_d = 400 \times 10^9 \times 200 \times 10^3 \times 132 \times 10^{-12} = 10560$$

$$80000 \times 132 \times 10^{-3}$$

$$\frac{132}{10560} \times 80$$

10,000 synchronization failure occur per second

## # The Origin of Metastability

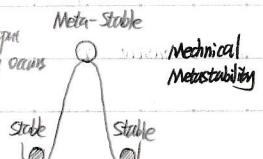


D changes at a very close time to the falling edge of the clock, so the output (Q) does not have sufficient time to make transition from 0 to 1. The SR Latch in this case will act as a cross coupled sense amplifier.

Assuming metastability happened at  $(t=0)$ , the output voltage of (Q) at time  $t$  will be given as,

$$V(t) = V_0 \cdot e^{\frac{t}{\tau}}$$

Value of the output when metastability occurs



## # What happens after metastability

Latch goes 1 or 0 randomly

Noise can have effect on the final output

Metastability does not usually last very long, but it can cause system failure

## # How long does the metastability last?

Metastability ends approximately when the output voltage of the latch reaches the threshold voltage of the NAND gate, by substituting  $V$  with  $V_{th}$  in the previous equation we get,

$$V_{th} = V_0 \cdot e^{\frac{t}{\tau}}$$

Therefore Metastability will approximately lasts :  $t_m = T \cdot \ln \frac{V_{th}}{V_0}$

Probabilistic Analysis shows that given a metastability event at  $t=0$ , the probability of metastability at  $t>0$  ( $P_m$ ) can be calculated as :

$$P_m(t) = e^{\frac{t}{\tau}}$$

## # Mean Time Between Failures (MTBF)

Failure means that flip-flop became metastable after the clock's sampling edge, and that it is still metastable  $S$  time later. The two events are independent, so we can multiply their probabilities.

The rate of failure ( $FR$ ) is calculated as the product of the metastability rate and the probability of staying in metastable state for  $S$  time.

$$FR = MP \cdot P_m(S)$$

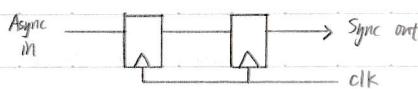
$$FR = W \times F_c \times F_d \times e^{-\frac{S}{T}}$$

$$MTBF = \frac{1}{Rate(failure)} = \frac{e^{\frac{S}{T}}}{W \times F_c \times F_d}$$

\* To design synchronous synchronizers for multiple clock domain SoCs

## # How to avoid Metastability

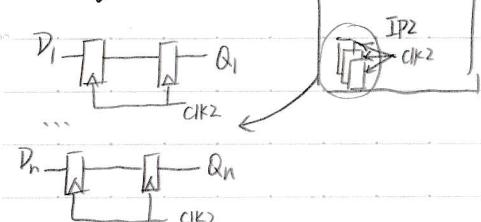
Allow more time for recovery: by forcing a full clock cycle between the two flip-flops



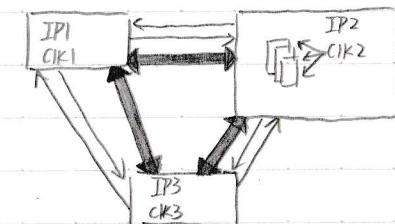
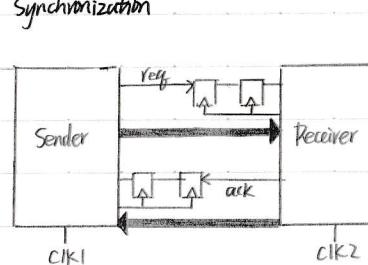
Use more than 2 flip-flops you can.

$$MTBF = \frac{e^{\frac{S}{T}}}{W \times F_c \times F_d} \left\{ \begin{array}{l} S \text{ predetermined time } S \text{ for metastability resolution} \\ f_c, f_d \text{ clock & data frequencies} \\ J, W \text{ technology parameters } (J \approx 2F_0/4 \text{ & } W \approx 4F_0/4 \text{ delays}) \end{array} \right.$$

## # Data Synchronization

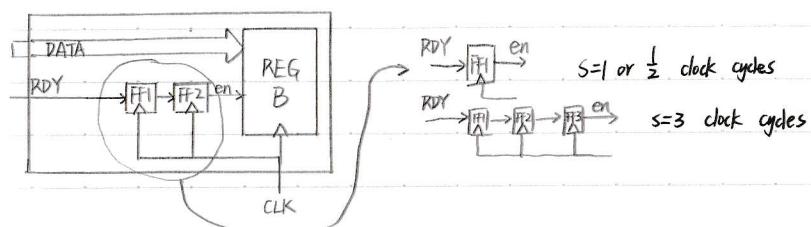


## # Control Synchronization



## # Control Synchronization: How to implement

$S=2$  clock cycles



R

W

$1400 \text{ M/S}$

$180 \text{ M/S}$

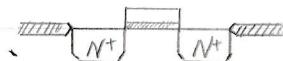
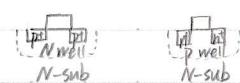
$\times 2$

$380 \text{ M/S}$

$\times 10$

$1900 \text{ M/S}$

## Twin-Tub

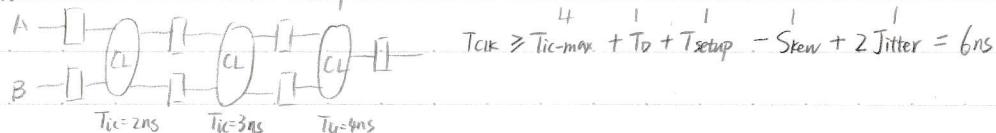


是 NPN 三极管  
以 P 为基，所以上面  
的 NPN 有一个 N-well process  
由于这个 NPN 没有 P-well，  
所以不是 Twin-Tub  
工艺为 N-well process

↑ NMOS with for N-well process

Positive Skew = 1ns

Jitter = 0.5ns,  $T_D = 1\text{ns}$ ,  $T_{Setup} = 1\text{ns}$



0 0	0
1 0	0
0 1	0
1 1	1

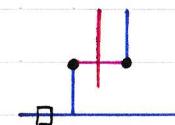
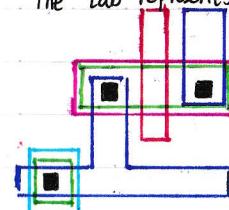
Acknowledge the change  
From 0 → 1 (00,10,01 → 11)  
1 → 0 (11 → 00,10,01)



not shown  
on stick diagram

## # Tabs unfilled black square ✓ □

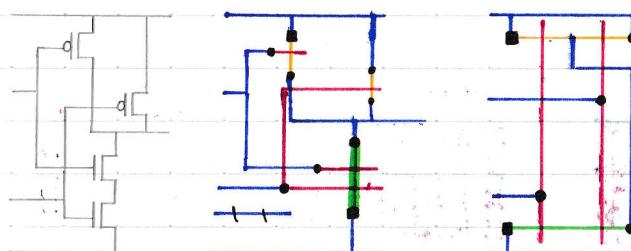
The tab represents a connection to something we can't see. either the N-well or the wafer substrate



## ✓ Filled Black Square ■

We can often save space by using a combined contact and tap.  
Here the tap shares the same active region area as the contact.

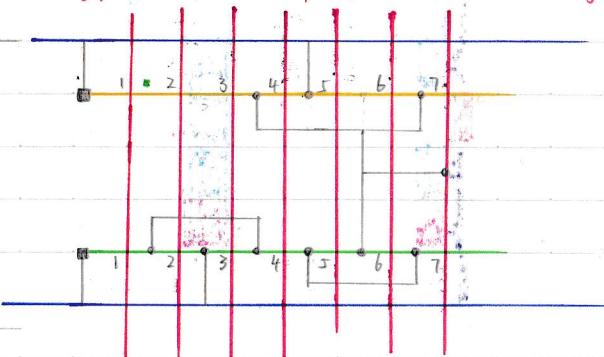
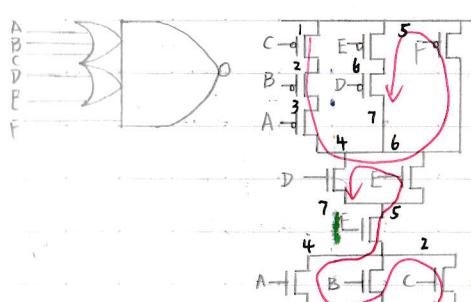
## [Stick Diagram]



- ① Find Euler path
- ② Label poly columns
- ③ Route power nodes
- ④ Route output node
- ⑤ Route remaining nodes
- ⑥ Add taps for PMOS and NMOS

[Euler Path]  $Z = \overline{(A+B+C) \cdot (D+E)} \cdot \bar{F}$ 

(A graph will contain an Euler path if it contains at most two vertices of odd degree)  
(A graph will contain an Euler path if all vertices have even degree)

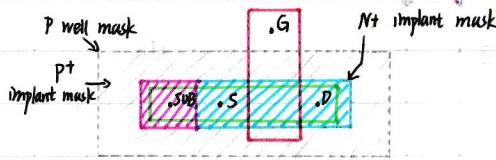
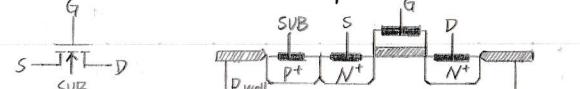


每跨一个桥留一个标记数字

\* A combined contact and tap, ■, may be used only where a power contact exists at the end of a line diffusion. Where this is not the case a simple tap ── should be used

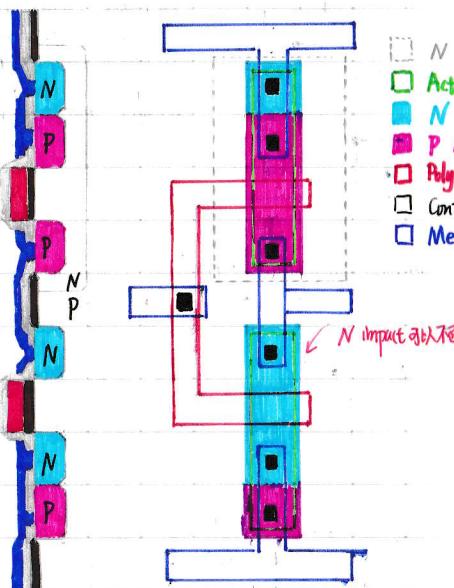
## # CMOS

[NMOS Transistor] - with top substrate connection

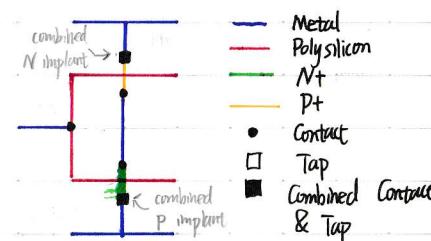


Five masks must be used to define the transistor  
 - P well (for isolation)  
 - Active Area  
 - Polysilicon  
 - N<sup>+</sup> implant  
 - P<sup>+</sup> implant  
 Top substrate connection  
 P<sup>+</sup>/N<sup>+</sup> implants produce good ohmic contacts

## # CMOS Inverter



## [Stick Diagram]



- Metal
- Polysilicon
- N<sup>+</sup>
- P<sup>+</sup>
- Contact
- Tap
- Combined Contact & Tap

## [Resistance]

$$R = \frac{P}{t} \left( \frac{L}{W} \right)$$

$t$  and  $P$  are fixed for a particular mask layer

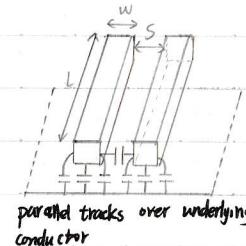
$$R = R_s \left( \frac{L}{W} \right)$$

$R_s$  is the sheet resistance  $0.152/\Omega$  for 170nm thick copper

1	2	3	4	5	6	7	8	9	10
3	18	17	16	15	14	13	12	11	2
4	20	21	22	23	24	25	26	27	5
7	28	29	30	31	32	33	34	35	28
8	36	37	38	39	40	41	42	43	44

(Note)  
 Serpentine shape to save on area

## [Capacitors]



parallel tracks over underlying conductor

$$C = C_a \cdot w \cdot L + 2 C_f \cdot L$$

$C_a, C_f, C$  are constants for a given layer and process  
 In digital designs our only aim is to minimize parasitic capacitance



metal-insulator-metal capacitor

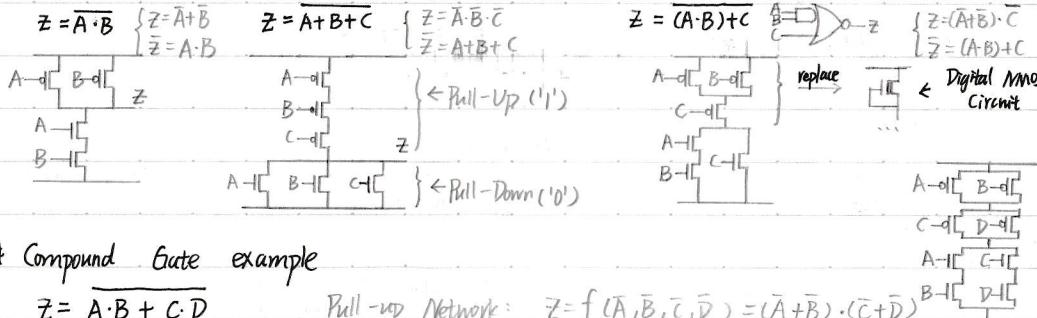
$$C = C_c \cdot 4/s$$

## # CMOS (Complementary Metal-Oxide-Semiconductor)

ALL CMOS devices are enhancement mode

An active PMOS device complements the NMOS device giving  
 - rail to rail output swing  
 - negligible static power consumption

ALWAYS Inverting



## # Compound Gate example

$Z = A \cdot B + C \cdot D$

Pull-up Network:  $Z = f(\overline{A}, \overline{B}, \overline{C}, \overline{D}) = (\overline{A} + \overline{B}) \cdot (\overline{C} + \overline{D})$

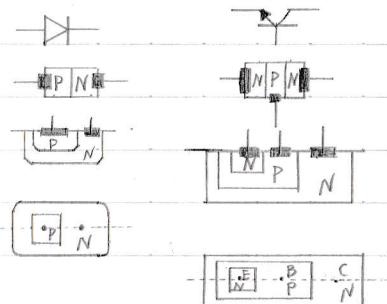
Symbol:

X-Wire

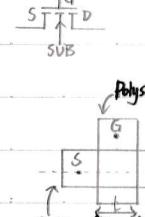
Pull-down Network:  $Z = f(A, B, C, D) = A \cdot B + C \cdot D$

## # Components for IC Design

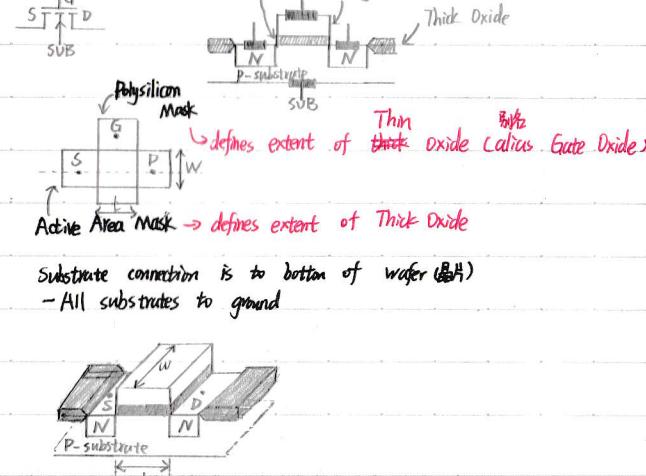
## [Diode]



## [NPN Transistor]



## [NMOS]



## Overview of Technologies

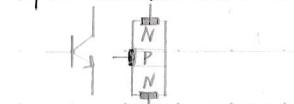
## # Moore's Law

The number of components per chip had doubled every year (1965 - 1975) (1975 - 2017) → Doubling every two years

## # Components for Logic

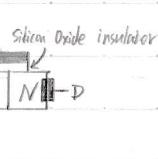
Diode:

Bipolar Transistor:



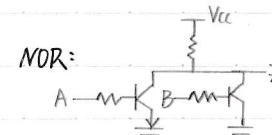
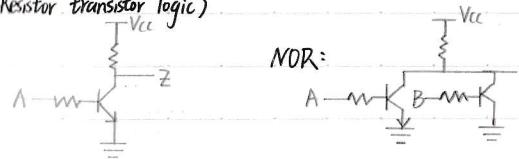
MOSFET:

Enhancement: off at zero gate-source voltage  
 the device is normally ON at zero gate-source voltage  
 Depletion: (For N-type, the threshold voltage might be -3V)



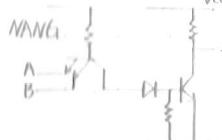
## # RTL (Resistor-transistor logic)

Inverter:

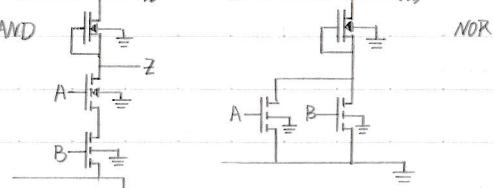


## # TTL (Transistor-transistor logic)

NAND:



## # NMOS



TTL gives faster switching than RTL at the expense of greater complexity.

Depletion transistor acts as non-linear load resistor  
 ↳ Resistance increases as the enhancement device turns on, thus reducing power consumption