

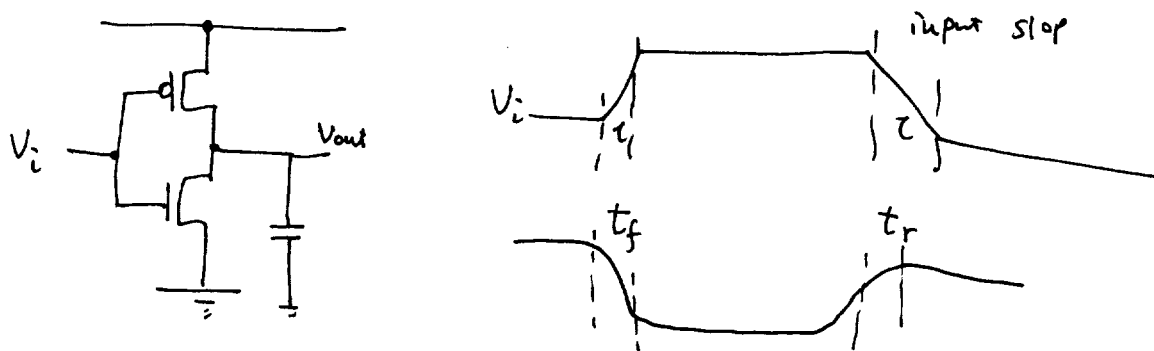
# RT and Algorithmic-level optimization for Low power (10)

- A very new area of research
- Target of the optimization is dynamic (switching) component of power.
- RT-level power optimization
  - \* architecture-driven voltage down scaling
  - \* Reduction of total switching capacitance of the design
- Architecture-driven voltage down scaling

$$P = \frac{1}{2} E_{sw} \cdot C_L \cdot V_{dd}^2 \cdot f$$

∴ Try to reduce  $V_{dd}$ .

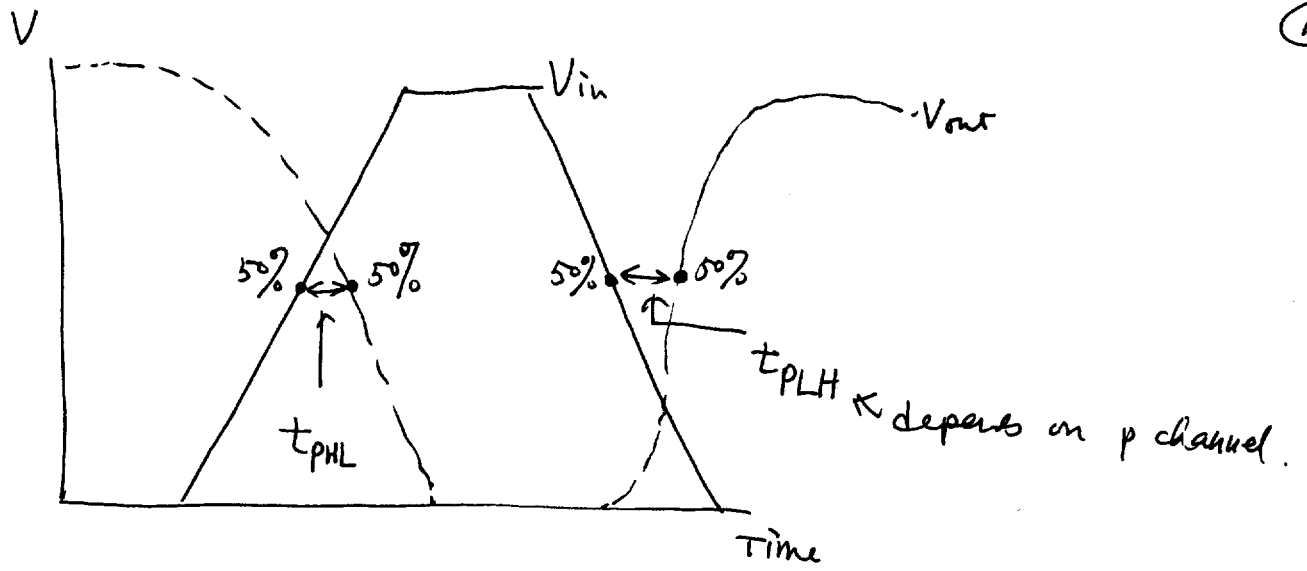
But, this affects circuit speed.



Circuit speed

Three different delays:

$t_f$ ,  $t_r$ ,  $t_{pHL}$



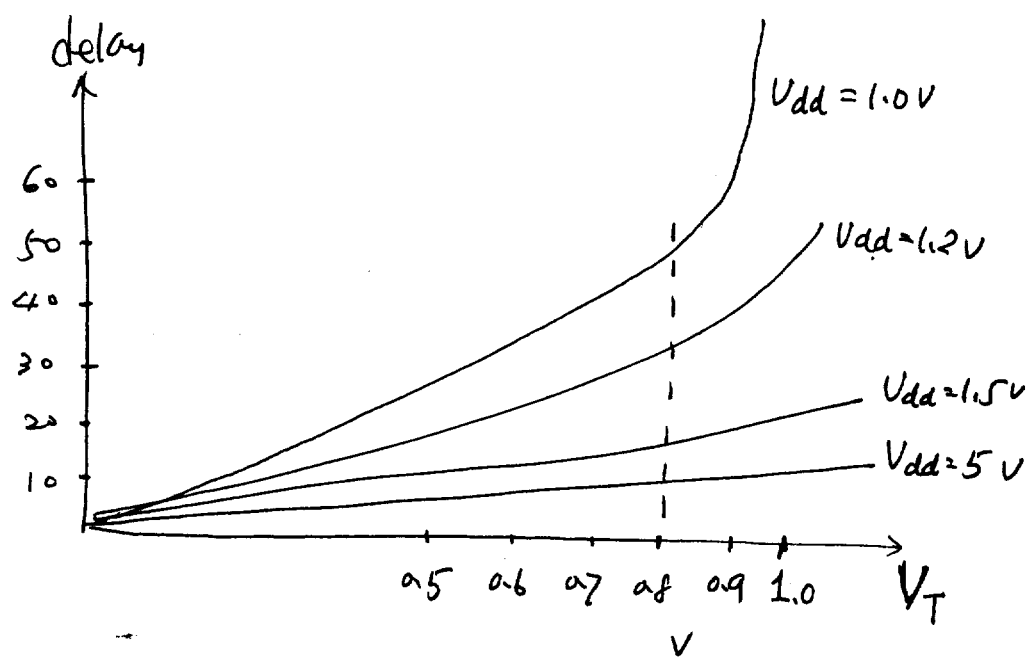
$$t_{PLH} = \frac{C_L}{\beta_p (V_{DD} - V_{TP})} \left[ \frac{2V_{TP}}{V_{DD} - V_{TP}} + \ln \left( \frac{3V_{DD} - 4V_{TP}}{V_{DD}} \right) \right] + \left[ 1 + \frac{2V_{TP}}{V_{DD}} \right] \frac{t}{6}$$

where  $\beta_p = \frac{W_p}{L_p} \frac{\epsilon_{ox} \mu}{t_{ox}}$

$t_{PLH}$  can be approximated by  $\frac{C_L V_{DD}}{K (W/L) (V_{DD} - V_T)^2}$   
 ↑  
 p or n depends on  $t_{PLH}$  or  $t_{PHL}$ .

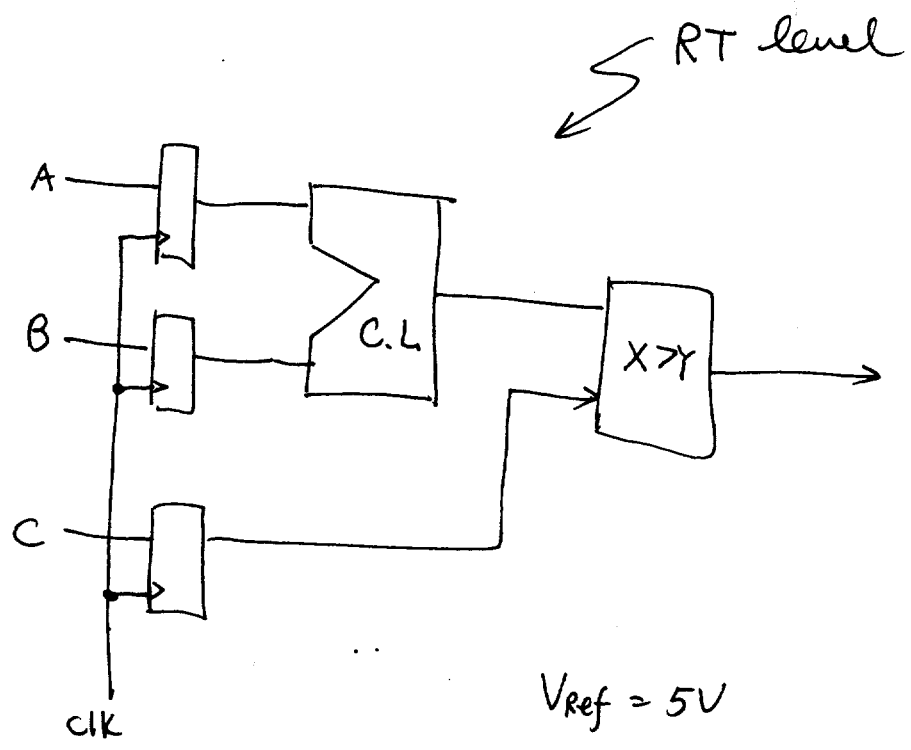
∴  $V_{DD} \downarrow \quad T \uparrow$

when  $V_{DD} \rightarrow V_{TP} \Rightarrow T \uparrow$  dramatically



- The increase in delay due to  $V_{dd}$  reduction must be compensated through architecture modification
- Use parallel/pipeline architecture to maintain the same throughput.

Example:



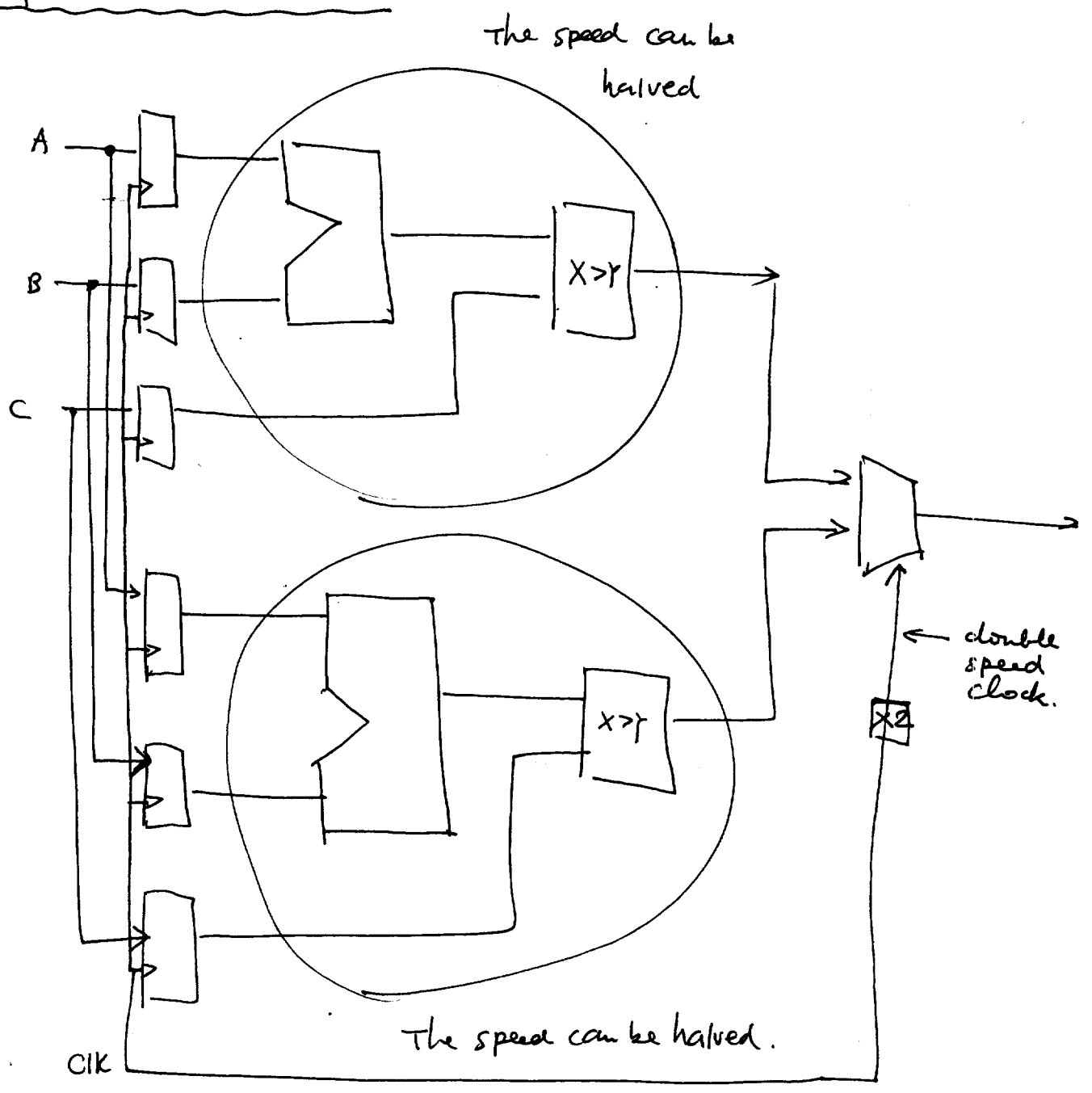
$V_{ref} = 5V$   
 Worst delay = 25 ns  
 $\therefore T_{ref} = 25 ns.$

$$P_{ref} = \frac{1}{2} C_{ref} V_{ref}^2 f_{ref}$$

(Switching activity is not  $\propto$  A, B, C are changed,  $\rightarrow$  ignored). Random #s

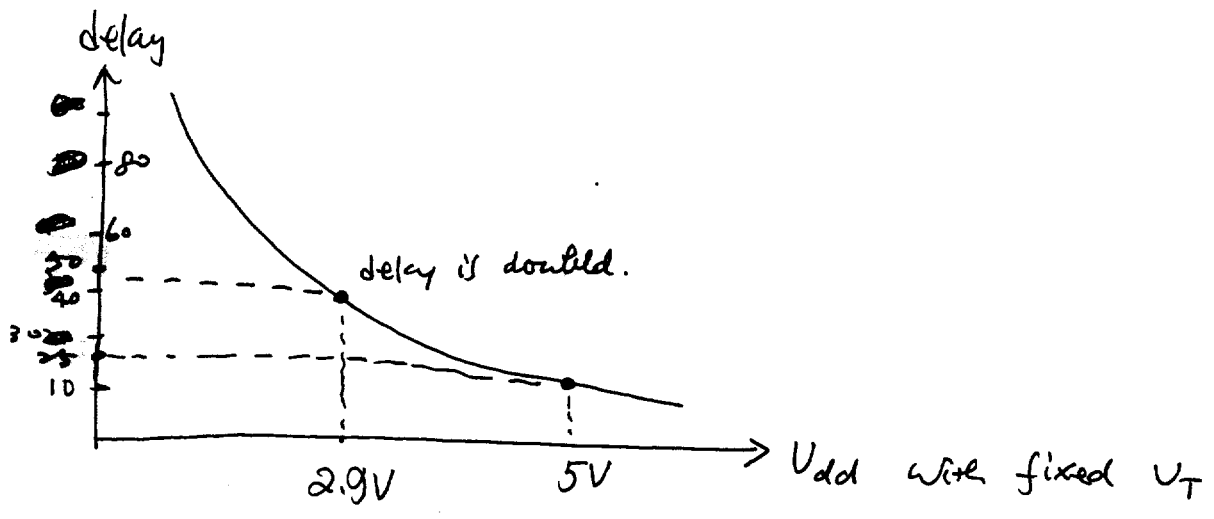
where  $C_{ref}$  is total effective capacitance of the ckt.

Use parallel architecture



$\therefore$  CLK can be reduced to  $T_{ref} = 50 \text{ ns}$  with the same throughput.

How much can the  $V_{DD}$  be reduced?

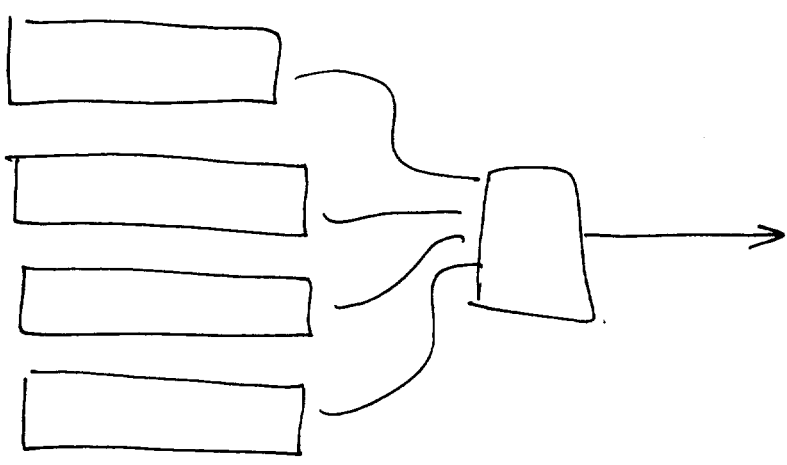


$\rightarrow V_{par} = 2.9V = 0.58 V_{ref}$

$C_{par} \approx 2.15 C_{ref}$

$\rightarrow P_{par} = \frac{1}{2} C_{par} V_{par}^2 f_{par} = \frac{1}{2} (2.15 C_{ref}) (0.58 V_{ref})^2 \cdot \frac{f_{ref}}{2}$   
 $\approx 0.36 P_{ref}$

How much parallelism can be used?

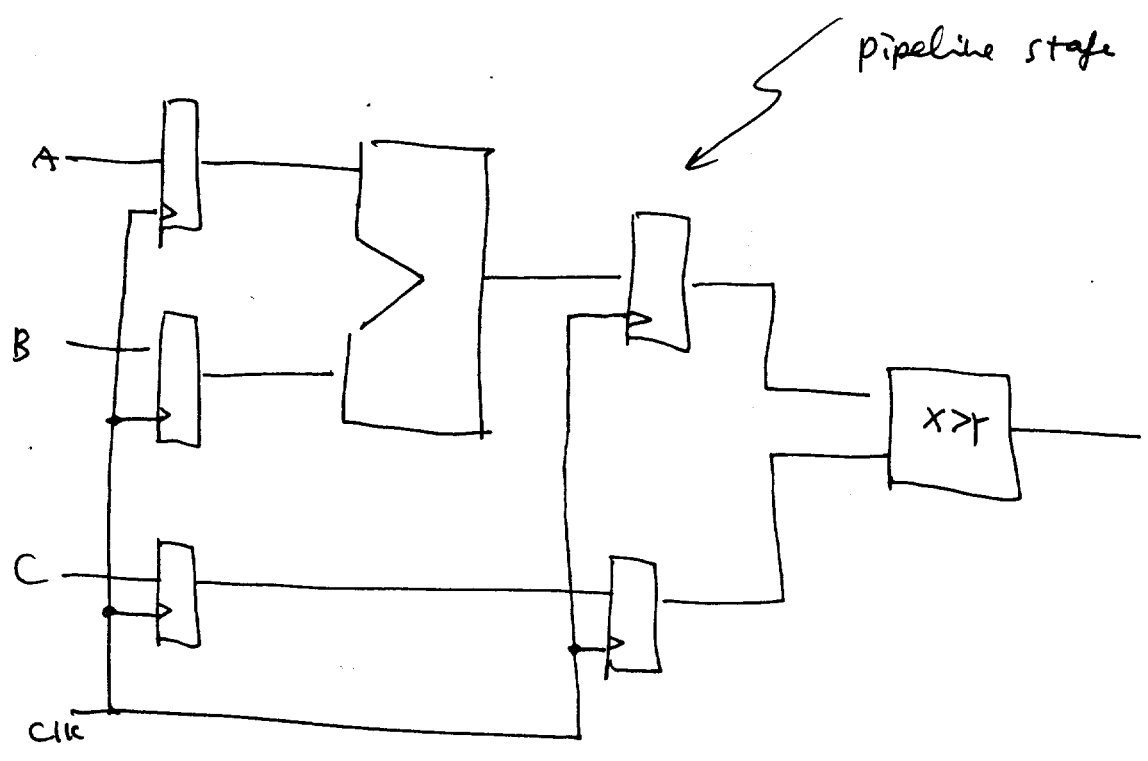


$\nearrow$   
 can be 4-times slower!  
 can  $V_{DD}$  be highly reduced?

Problem: When  $V_{dd}$  approaches  $V_T$ ,

- delay ↑ rapidly.
- much more duplication of the ckt is required
- power ↑ more than power saved!!
- parallelism does not pay off. when  $V_{op} \rightarrow V_T$ .

### Pipeline implementation



• The critical path has been reduced.

•  $Time = \max [T_{adder}, T_{comparator}]$

Worst delay of adder & comparator, respectively.

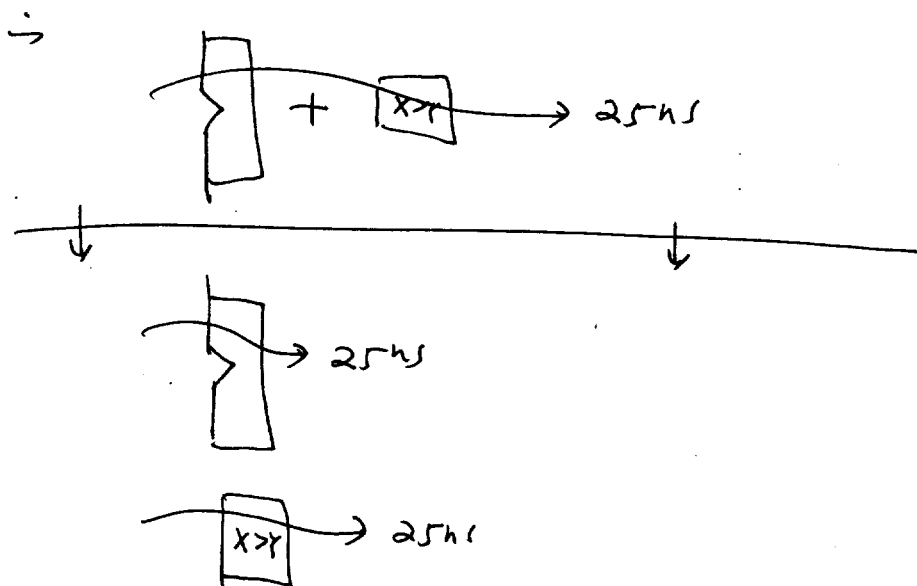
Assume  $T_{adder} = T_{comparator} = 12.5 \text{ ns}$

The clock ~~speed~~ <sup>period</sup> can be  $25 \text{ ns}$ .

But, we just like to keep the same speed

$\rightarrow$  clock speed <sub>period</sub> is still  $25 \text{ ns}$

$\rightarrow$  The ckt speed of each one can be slowed  
by 2 times. ( $25 \text{ ns}$ )



$\rightarrow$  Voltage supply of each ckt can be reduced

Again,  $V_{pipe} = 2.9 \text{ V} = 0.58 V_{ref}$  by checking the curve

$C_{pipe} = 1.15 C_{ref}$  (increased slightly only)

$f_{pipe} = f_{ref}$  (still work with the same speed)

$$\begin{aligned}
 P_{\text{pipe}} &= \frac{1}{2} C_{\text{pipe}} V_{\text{pipe}}^2 f_{\text{pipe}} \\
 &= \frac{1}{2} (1.15 C_{\text{ref}}) (0.58 V_{\text{ref}})^2 f_{\text{ref}} \\
 &\approx 0.39 P_{\text{ref}}
 \end{aligned}$$

∴ power saved is about the same as in parallel architecture

But, area penalty is much smaller.

Further, pipelining causes reduction of ~~set~~ ckt depth

⇒ reduce the power dissipation by glitches.

∴ pipelining is preferred if possible.

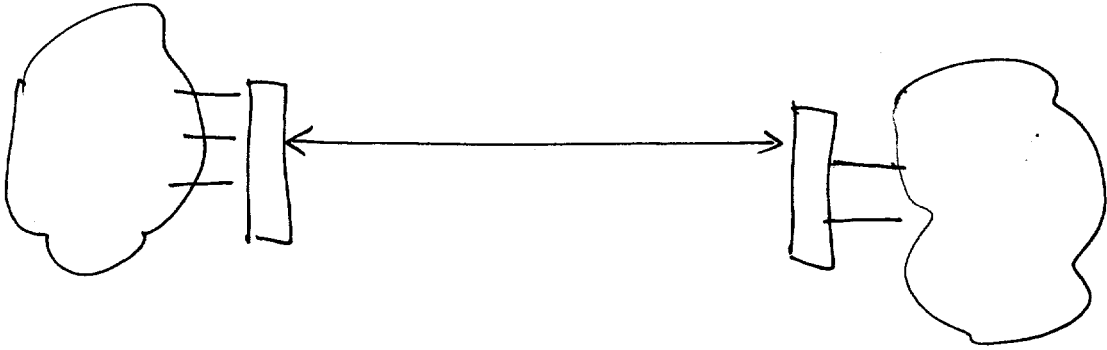
• Merging parallel & pipelining  $\implies$  Reduce even more significant power

$$\underline{\underline{P_{pp}}} \approx \underline{\underline{0.2 P_{ref}}}$$



# Effective Capacitance Reduction

## \* Data representation



- Try to find a good data representation s.t the effective capacitance can be reduced.

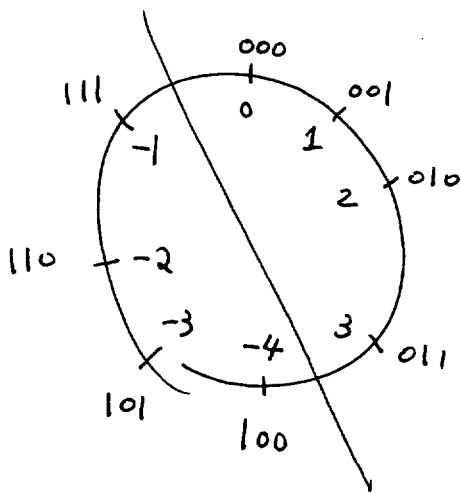
↑

Average switch activity  $\times$  total capacitive load

- Two's Complement representation

$2 - 3 = ?$

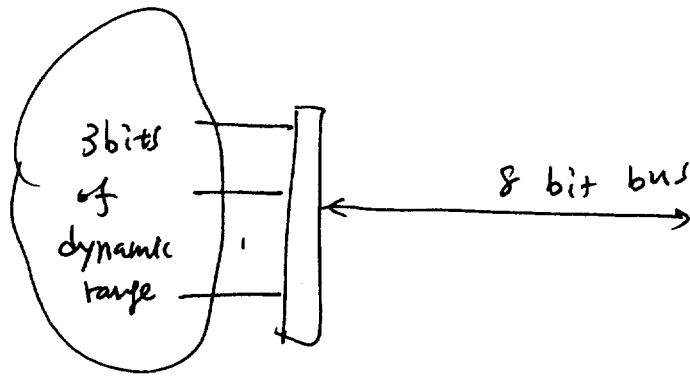
$2 + (-3) = -1$



$$\begin{array}{r} 010 \\ 101 \\ \hline 111 \leftarrow -1 \end{array}$$

- very useful in DSP and up chip
- make arithmetic operations simple.
- sign-bit extension.

Example:



111 = -1 must be extended to

(If no sign extension,  
 ① sign of # will be changed, ② 2's complement  
 sign extension properties cannot be held.)  
 11111111

- This happens when dynamic range  $\ll$  actual bit width.
- If the numbers being transferred or manipulated are switching around zero frequently + do not use the entire bit width  $\Rightarrow$  may incur lots of transitions.

$+1 = 00000001_2$        $-1 = 11111111_2$

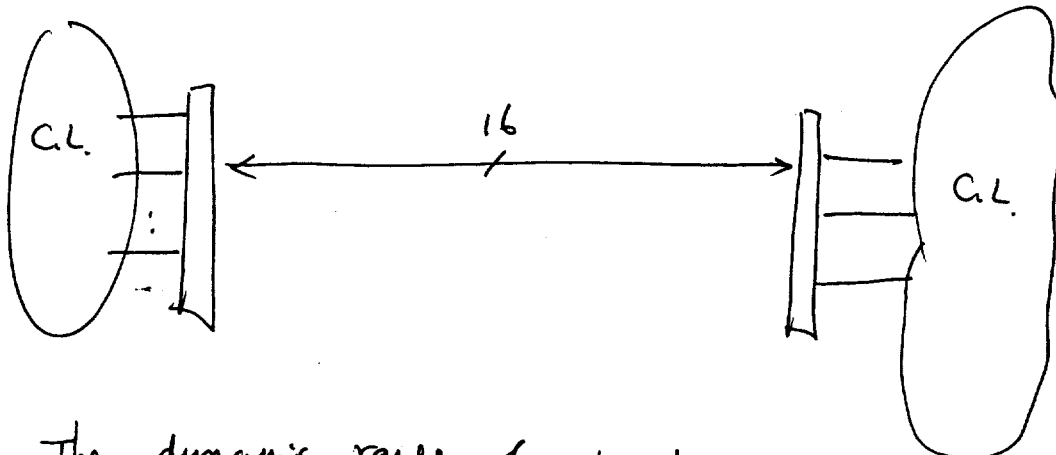
• However, if the #s are represented by sign and magnitude notation  $\Rightarrow$

$+1 = 00000001_2$        $-1 = 10000001_2$

very small # of switchings.

Conclusion: The degree of correlation between data being processed influences the power dissipation.

• Case study.



The dynamic range of numbers is mostly  $[-2^{11}+1, 2^{11}-1]$ .

→ only 12 bits out of the 16 bits are used.

Problems: • Try to determine the transition probability,  $P$ , for each of the 16 bus lines for

- ① 2's complement representation
- ② sign and magnitude representation.

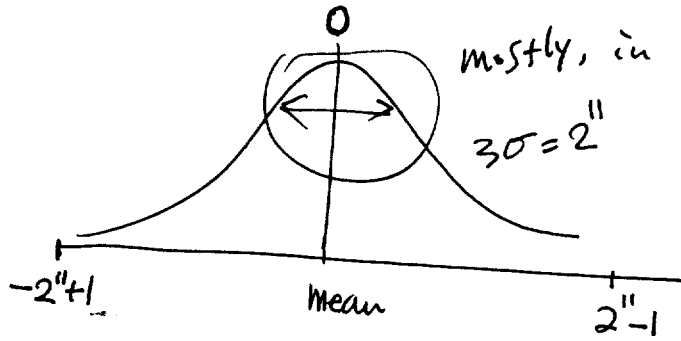
• Data correlations

① no correlation

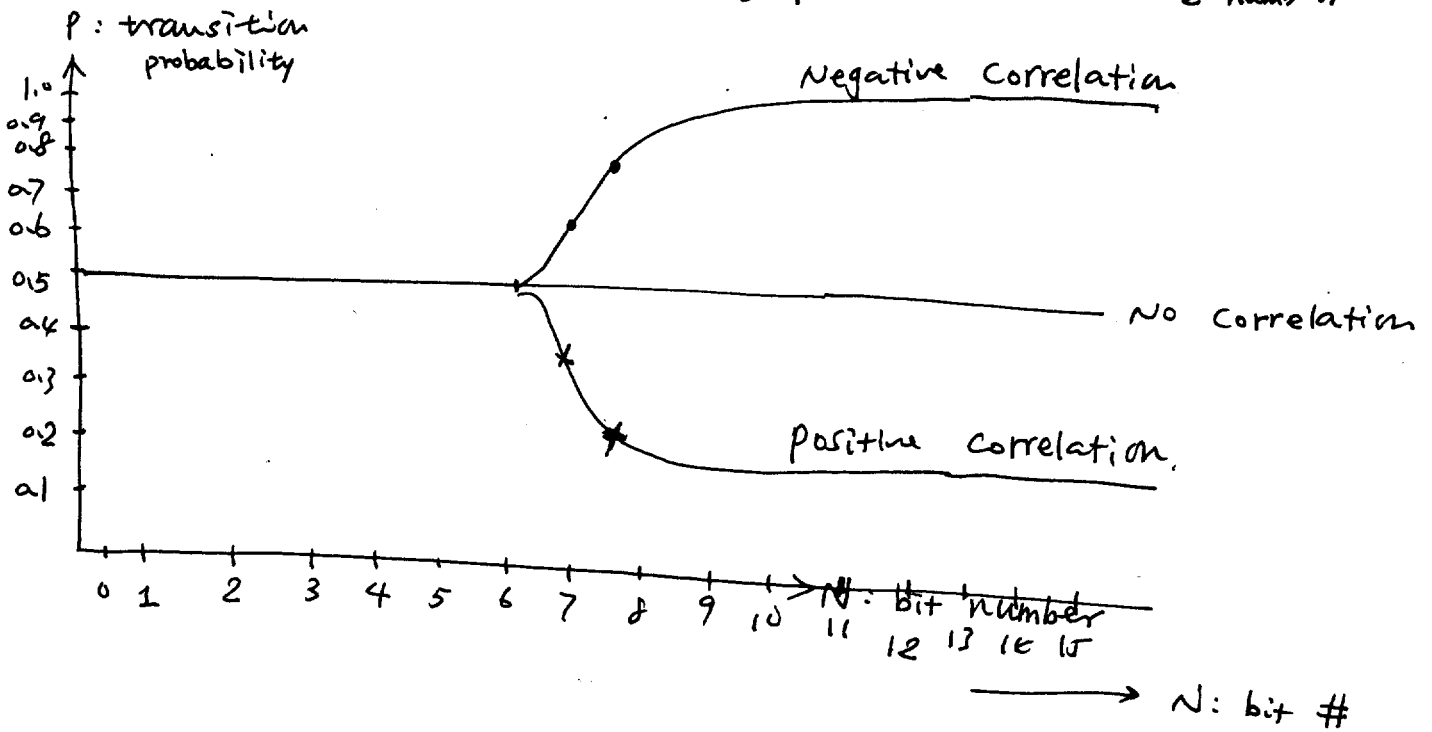
② positive correlation:  $P$  (2 consecutive #'s have same sign) ↑

③ Negative Correlation:  $P(2 \text{ consecutive \#s have same sign}) \downarrow$

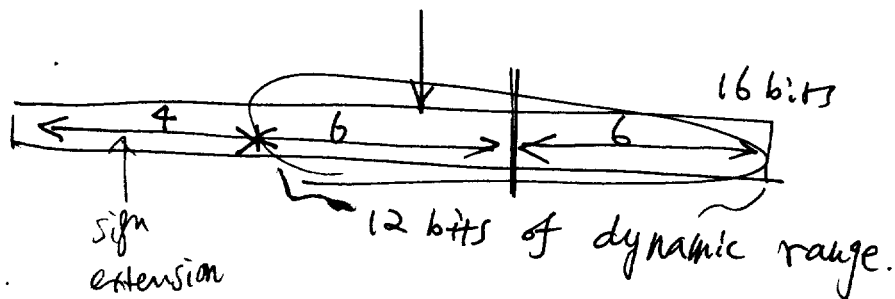
• Data distribution: Gaussian distribution

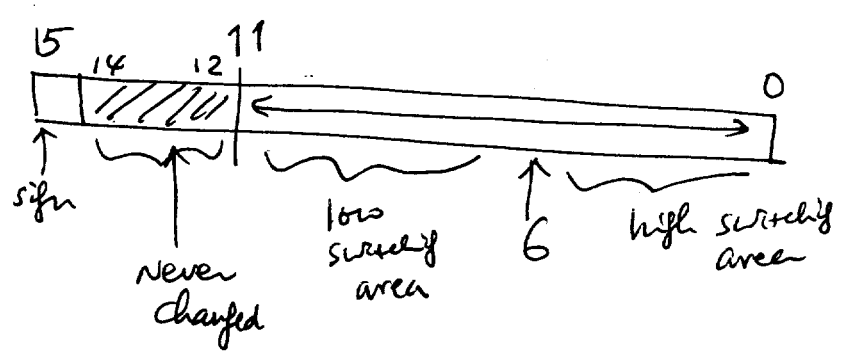
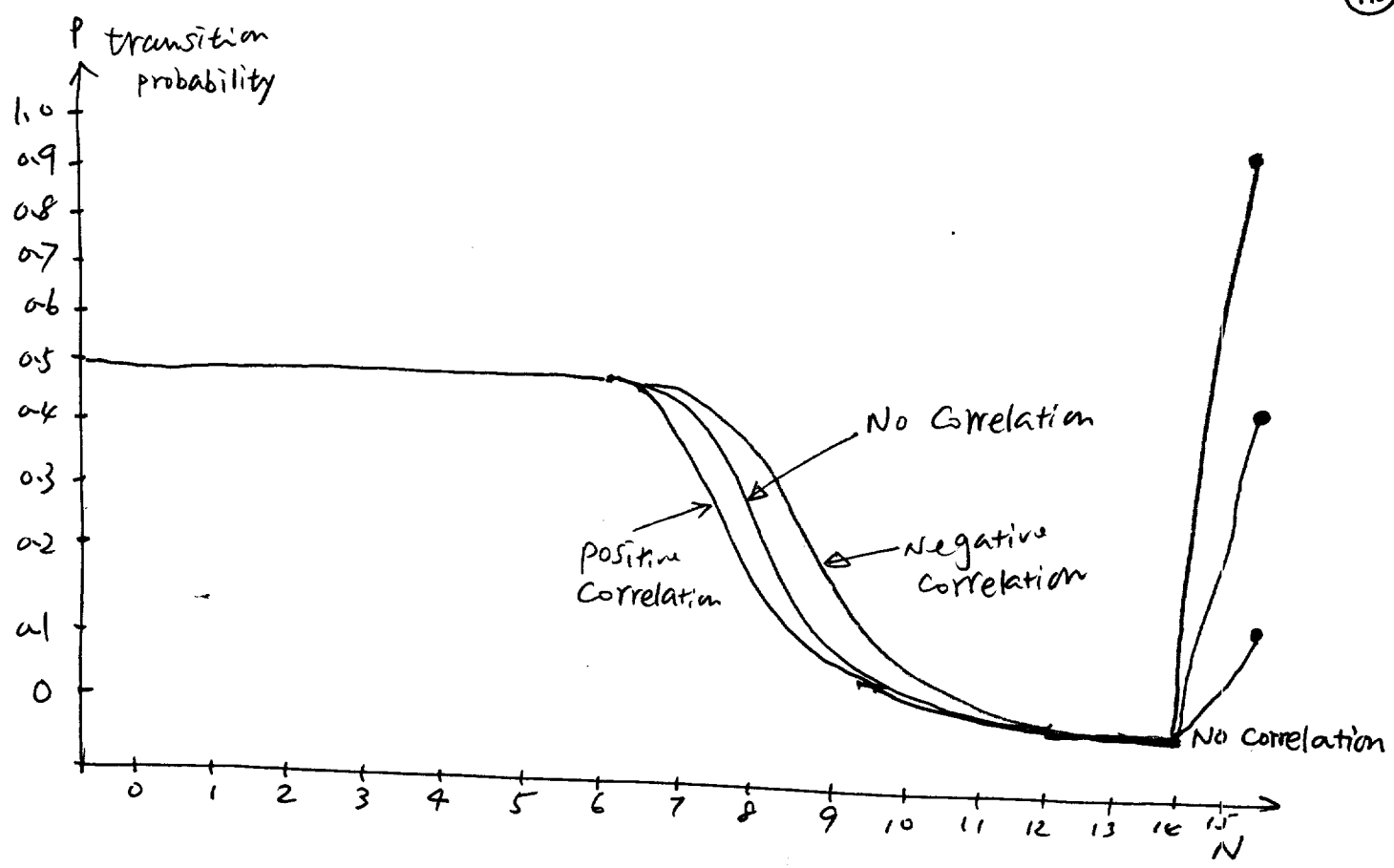


• 93.32% that data will be distributed



2's complement representation.





- The biggest win of sign and magnitude system is
  - ① Dynamic range of the represented numbers ~~is~~ <sup>is</sup> small
  - ② highly negatively correlated between the consecutive numbers.

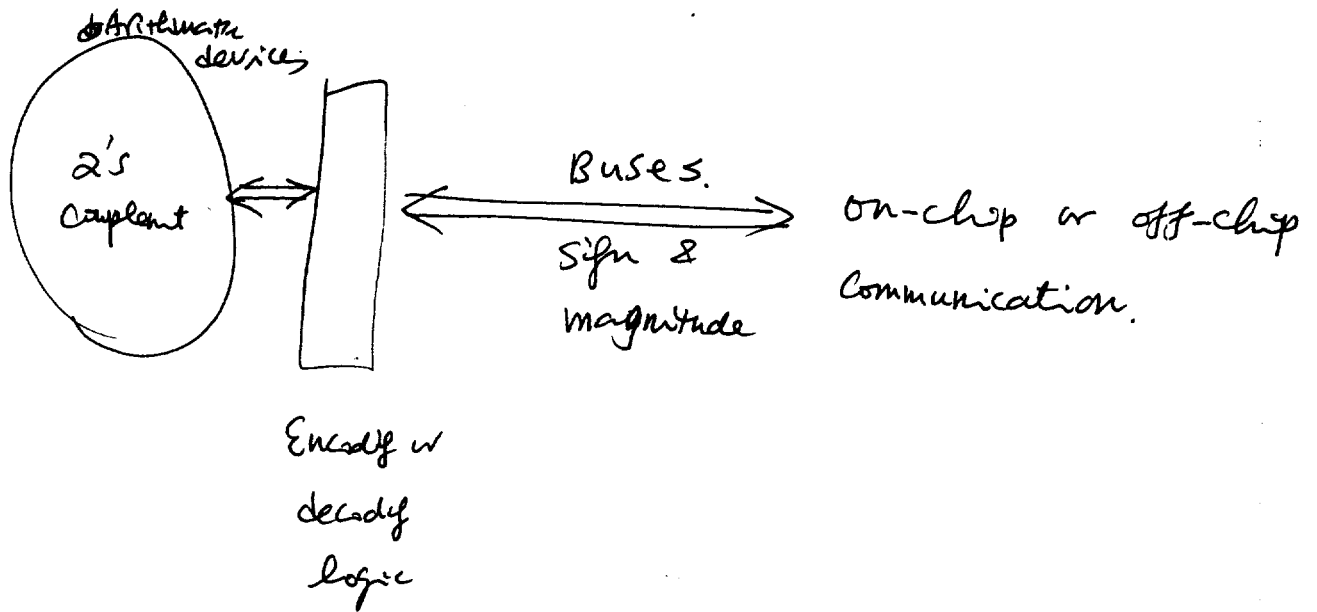
Conclusion: sign & magnitude representation  $\Rightarrow$  good for data transfer devices, e.g., buses

2's Complement representation  $\Rightarrow$  good for arithmetic ccts.

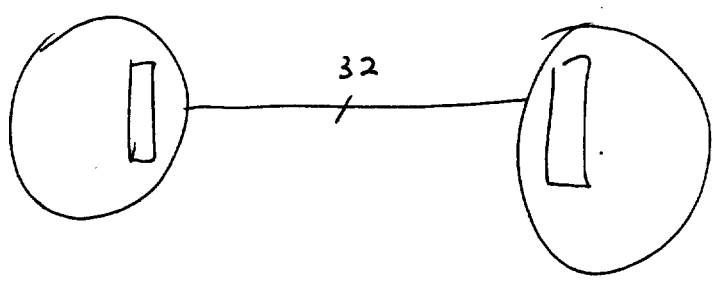
° simple arith  $\Rightarrow$  simple ccts  $\Rightarrow$  smaller ccts

$\Rightarrow$  lower power dissipation.

Best solution:



# Bus Encoding



- Try to use signed magnitude representation for communication
- Can be further improved by bus encoding.

- Motivation:

power dissipation at I/O pins accounts for 5% of overall switching power

∴ pin has very high capacitance.

∴ Should try to reduce the switching activities of I/O pins. by encoding.

- Encoding methods:

Active high encoding:

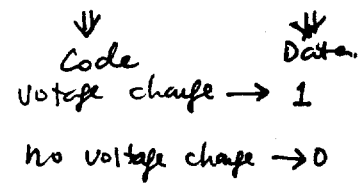
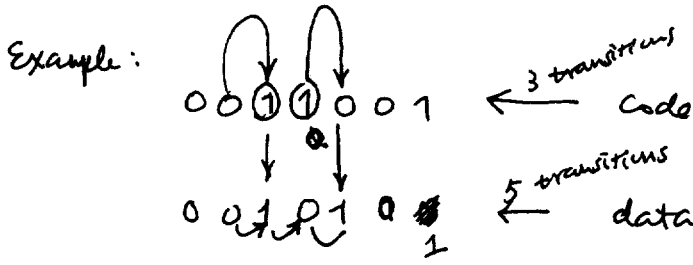
high-level voltage : 1

low-level voltage : 0.

transition-based encoding:

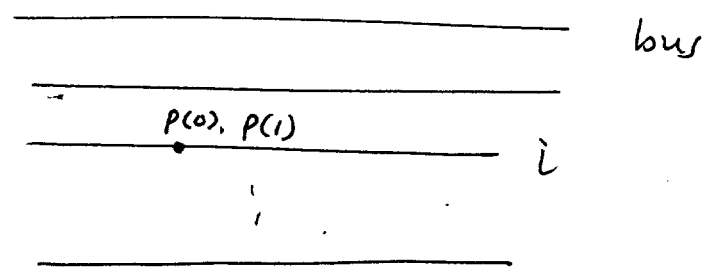
Voltage change: 1

No voltage change: 0.



Transition-based encoding

- Good only when there are many 0's transmitted.



Assume:  $p(0) > p(1)$

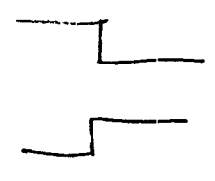
No temporal correlation between consecutive values

Switching activities:

Active-high encoding:

average # of transitions / clock cycle

$$p(0)p(1) + p(1)p(0) = 2p(1)(1-p(1))$$



if  $p(1) = \frac{1}{4}$   
 $p(0) = \frac{3}{4}$

⇒ average switching activity is  $\frac{3}{8}$ .

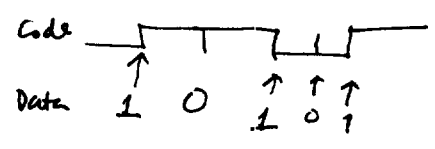


Transition-based encoding:

Transition required only when there is a logic one.

→ average switching activity is  $P(1)$ .

$\frac{1}{4}$  in this case

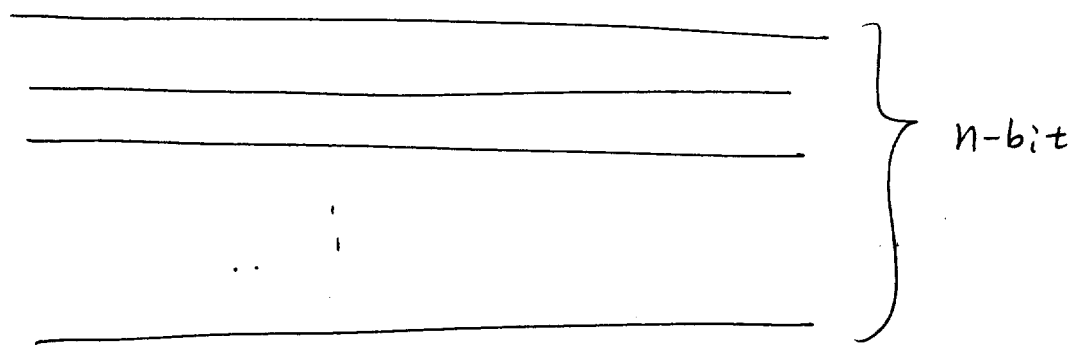


∴ 33% less than the active-high encoding.

• Transition-based encoding is good when inputs are strongly not equiprobable.

\* If  $P(0) > P(1)$  for a line, then the <sup>average</sup> switching activity of the line is smaller (than active high code).

\* Try to use limited-weight code to make sure that  $P(0) > P(1)$  for each bit line.



K-limited-weight code: There are  $k$  1's ~~with~~ in the  $n$  bits.  
 ↑  
 can represent  $\binom{n}{k}$  data

Example:  $k=1$ .

0	0	0	0	1	← code word 1
0	0	0	1	0	← code word 2
0	0	1	0	0	⋮
0	1	0	0	0	⋮
1	0	0	0	0	← code word 5

5 code words

$$\rightarrow P(0) = \frac{n-1}{n} \quad P(1) = \frac{1}{n}$$

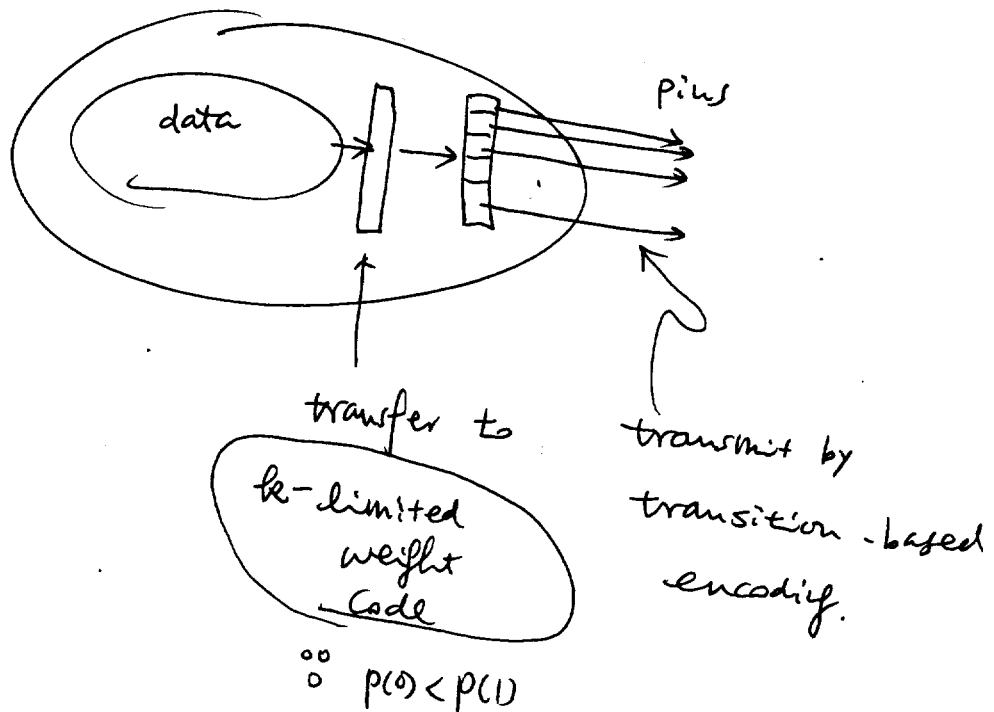
$$k = n/2$$

$$P(0) = P(1) = \frac{1}{2}$$

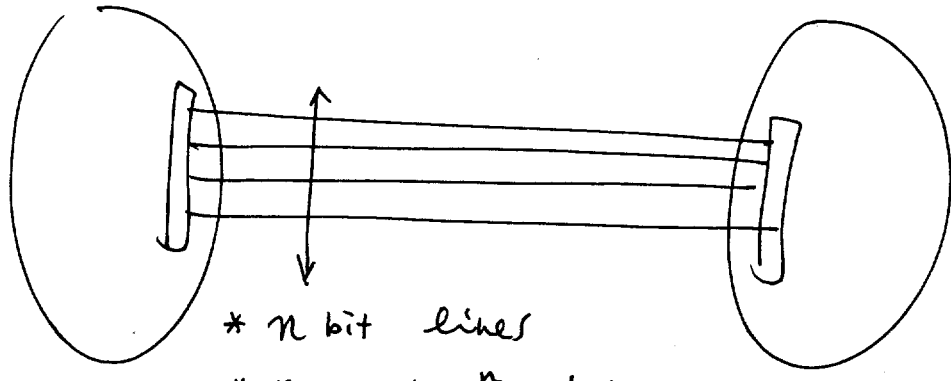
\* for  $1 \leq k < n/2$ , we are sure that  $P(0) > P(1)$

↳ the transition-based encoding benefits.

Conclusion:



Bus inverting encoding

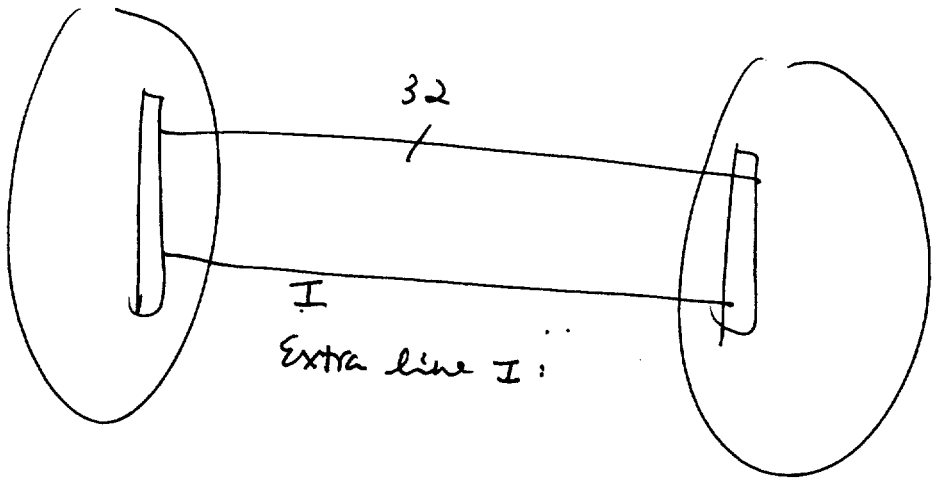


- \*  $n$  bit lines
- \* represent  $2^n$  data
- \*  $p(0) = p(1) = \frac{1}{2}$  for each line

∴ 
$$\frac{\text{Average \# of transitions}}{\text{clock cycle}} = \frac{n}{2}$$

$$\frac{\text{Worst \# of transitions}}{\text{clock cycle}} = n$$

\* Try to reduce both average (worst) # of transitions per cycle.



● Compare each pair of consecutive patterns.

① If Hamming distance between current & next patterns  $\leq n/2$ ,  $I=0$ , send as what it is.

② If Hamming distance between \_\_\_\_\_  $> n/2$ ,  $I=1$ , transmit the data.  
⏟  
 pattern is converted

∴ This guarantee that the max # of transition per clock cycle is  $n/2$ .

\* By statistical analysis, the avege # of transitions per clock cycle is lowered by 25% of the original.

- Hardware overhead for code generation must be invested.
- All above discussions are for data bus.

• How about address bus?

• Addresses generated by a running up are often consecutive.

∴ Try to represent address by <sup>G</sup>Gray Code.

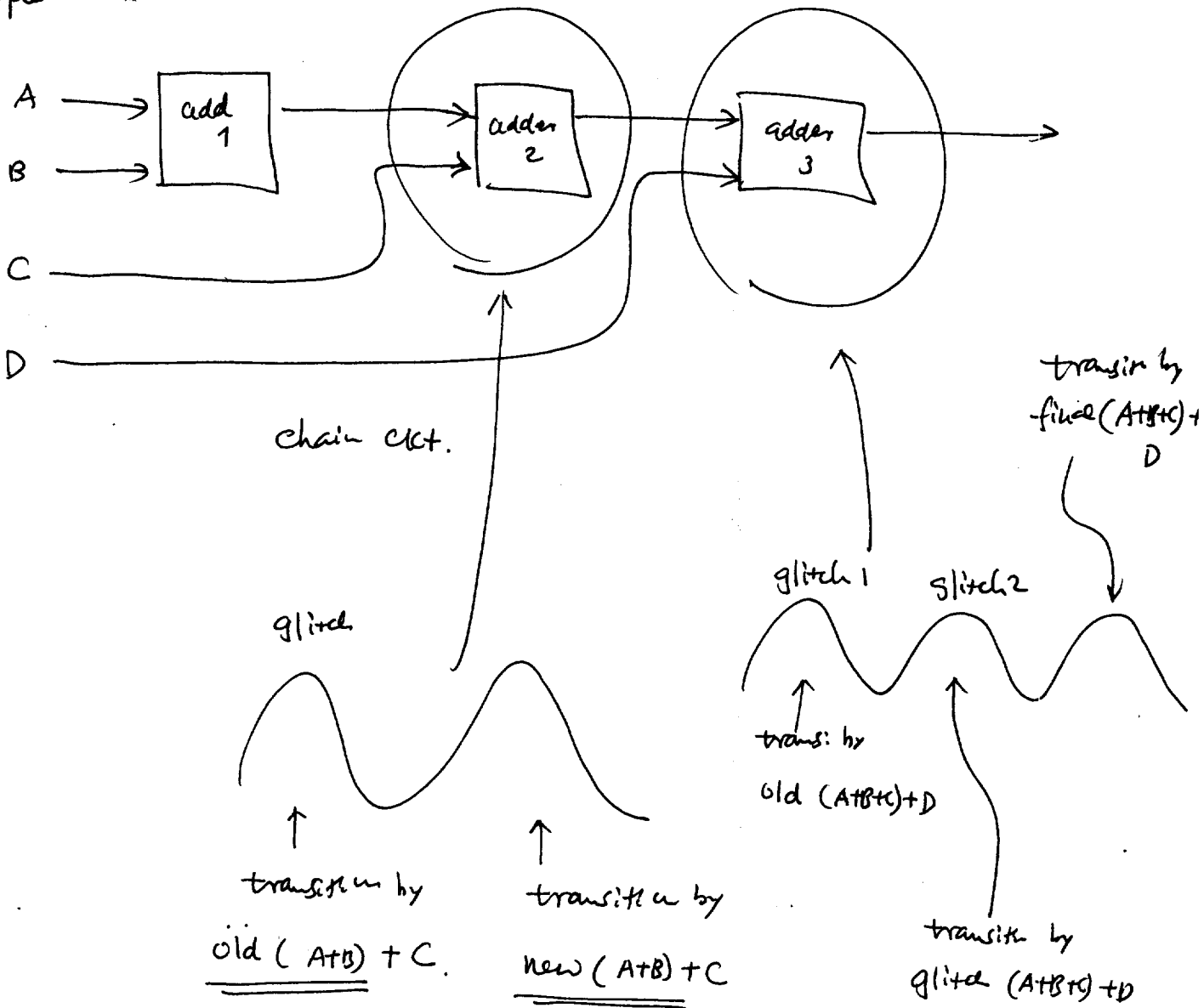
(guarantee single-bit change for current & next address)

• Computer simulation demonstrate about 37% power reduction for several benchmark programs.

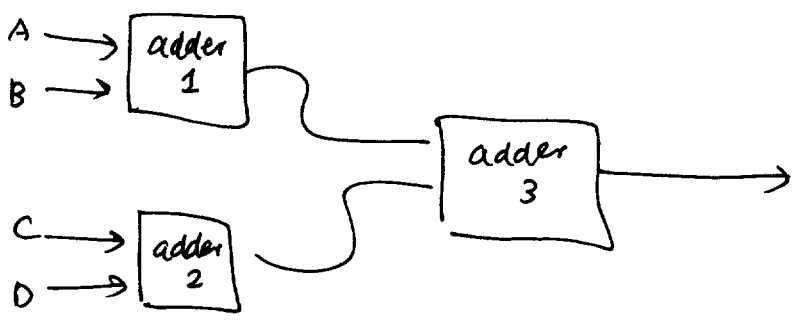
• Signal Synchronization

• Try to reduce the # of glitches in RT-level, by balancing signal paths or reducing logic depth.

Excp: -

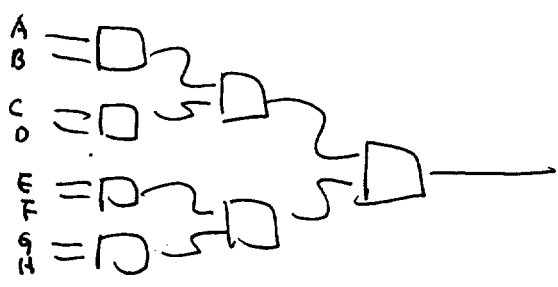
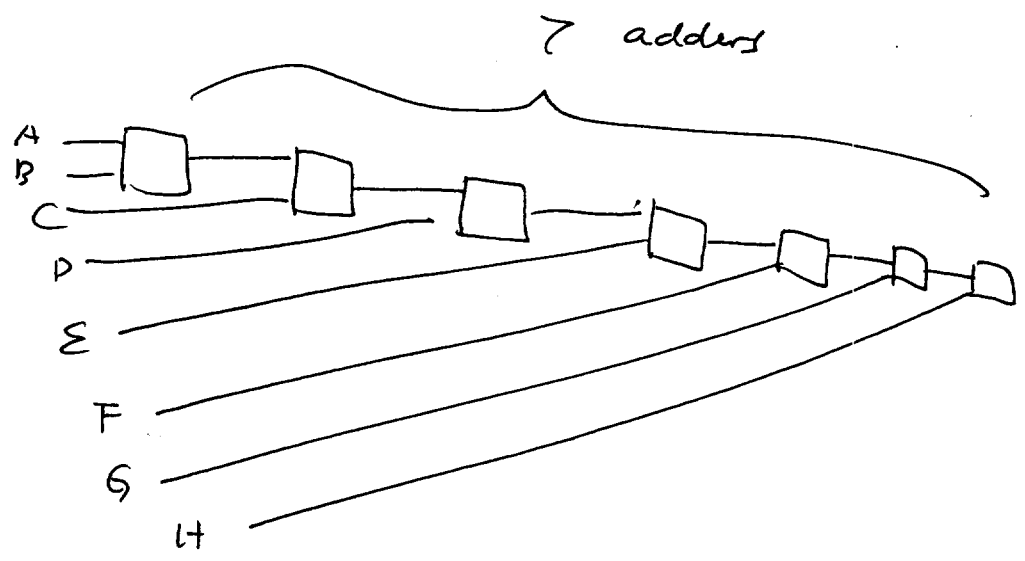
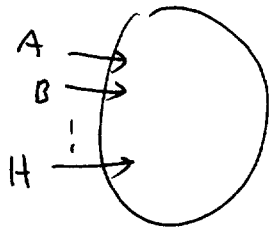


• Low-power design



\* paths are balanced, so the # of glitches is minimized.

• 8 inputs



Decreasing the logic depth increases the # of registers required by the design ? ? → add extra power consumption!!

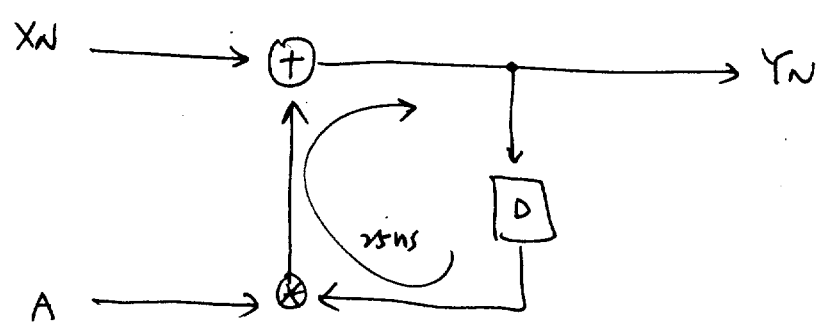
• Resource sharing: ready assignment.

• Algorithmic-Level optimization

• Try to modify the computational structure of the algorithm while preserving I/O behavior

• Objective: optimize the power dissipation with a specific throughput.

• Speed-up transformations



$$Y_N = A \cdot Y_{N-1} + X_N$$

first-order infinite impulse response IIR filter

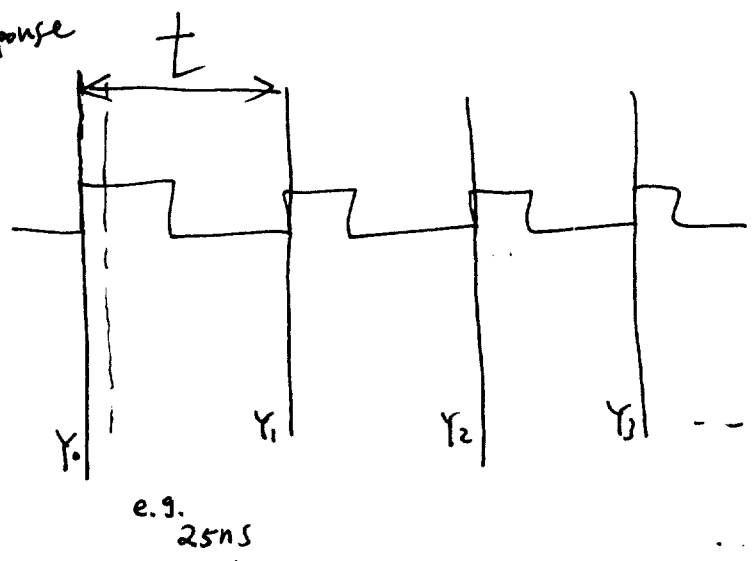
$$C_{eff} = 1$$

$$V_{dd} = 5V.$$

$$Throughput = 1$$

$$Critical\ path = 2$$

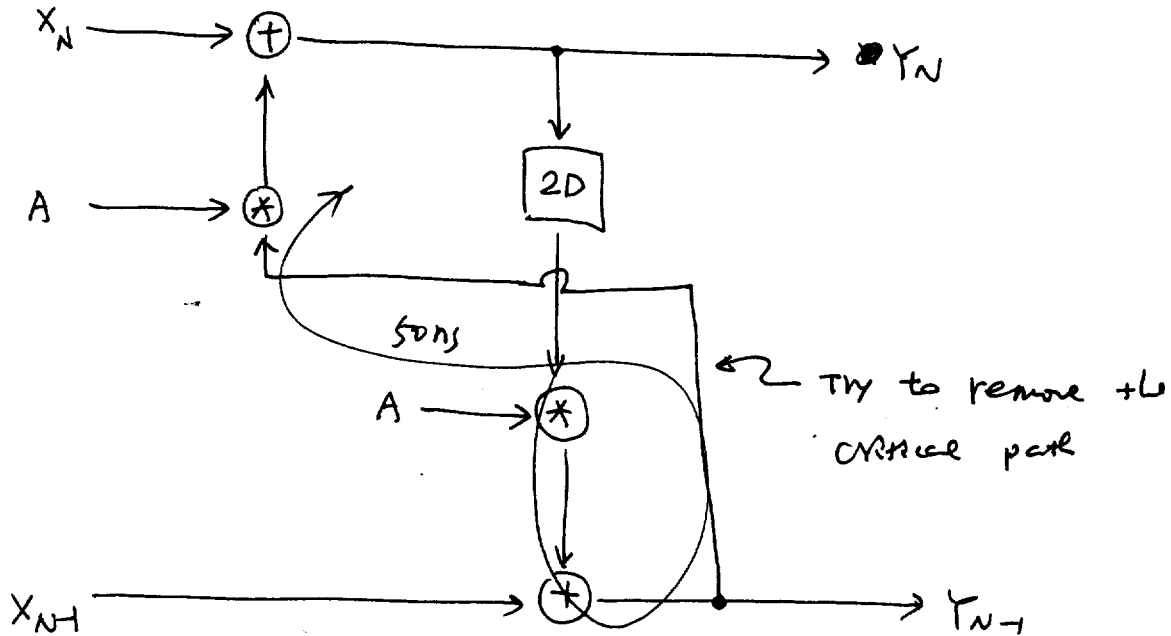
$$C_{eff}^2 \rightarrow power = 25$$



e.g. 25ns

• Can't be further optimized by pipelining, re-timing, - -

• Try to unfold the loop



$$Y_{N-1} = Y_{N-2} \cdot A + X_{N-1}$$

$$Y_N = Y_{N-1} \cdot A + X_N$$

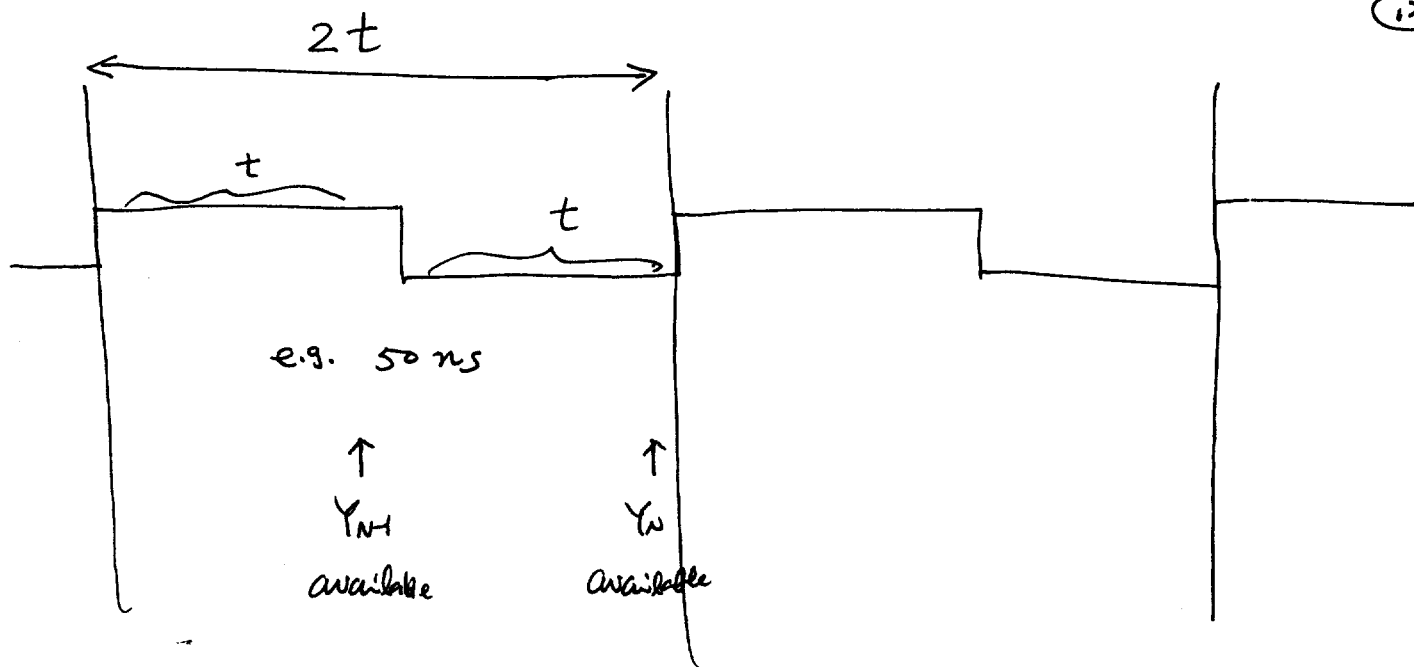
$$= (Y_{N-2} \cdot A + X_{N-1}) A + X_N$$

$$= Y_{N-2} A^2 + X_{N-1} A + X_N$$

for next iteration

∴ function is the same.





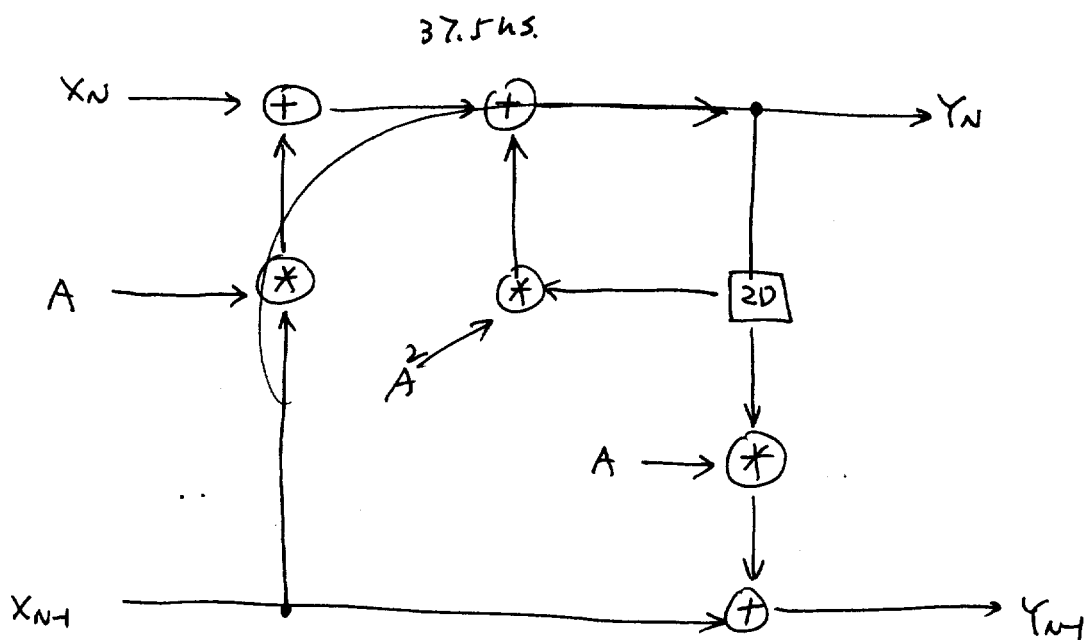
$C_{eff} = 1$  ( $C \uparrow$  but  $E(s.w).f \downarrow$ )

$V_{dd} = 5V$

throughput = 1

critical path = 4, But effective critical path = 2.

power = 25



$$Y_{N-1} = Y_{N-2} \cdot A + X_{N-1}$$

$$Y_N = (X_{N-1} \cdot A + X_N) + A^2 \cdot Y_{N-2}$$

$$= A^2 Y_{N-2} + A X_{N-1} + X_N$$

\* Critical path has been reduced from 4 to 3.

→ Supply voltage can be reduced to keep the same throughput

$C_{eff} = 1.5$  (added more multipliers & adders)  
Also, more devices are switching simultaneously.

$$V_{dd} = 3.7V$$

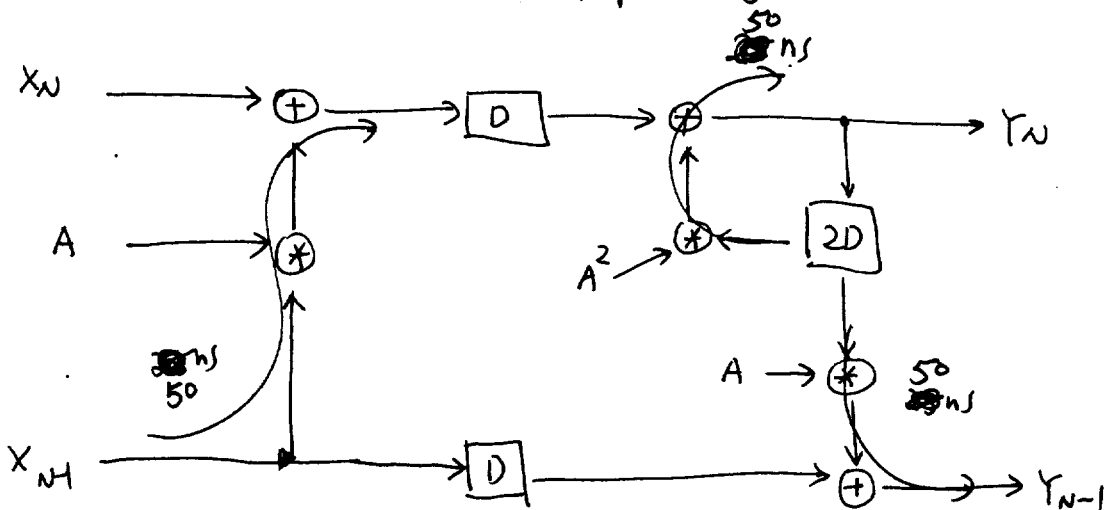
$$\text{Throughput} = 1$$

$$\text{Critical path} = 3$$

$$\text{Power} = 20$$

→ 20% power saving over the original design.

\* Further improvement by pipelining



$$\begin{array}{r}
 42 \\
 34 \\
 \hline
 168 \\
 126 \\
 \hline
 1428
 \end{array}$$

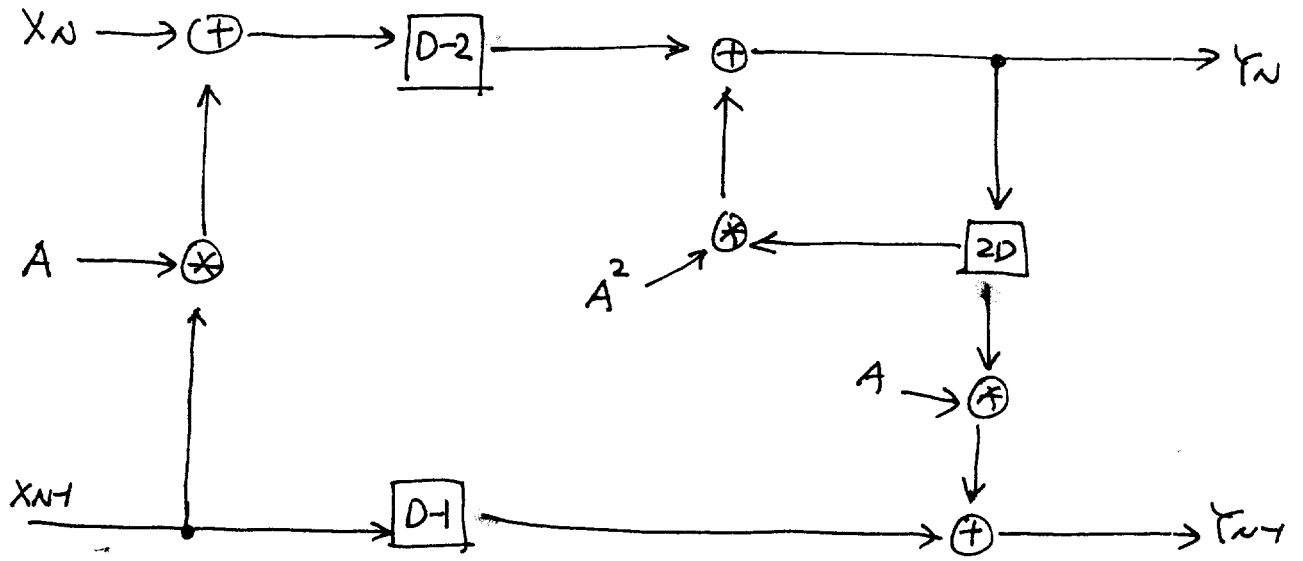
$$Y_0 = X_0 + A \cdot Y_{-1}$$

$$Y_1 = X_1 + A \cdot Y_0 = X_1 + A (X_0 + A \cdot Y_{-1})$$

$$Y_2 = X_2 + A \cdot Y_1 = X_2 + A \cdot X_0 + A^2 \cdot Y_{-1}$$

$$Y_3 = X_3 + A Y_2 = X_3 + A (X_2 + A Y_1)$$

$$= X_3 + A X_2 + A^2 Y_1$$



$x_{n-1}$	$x_n$	D-1	D-2	2D	$y_{n-1}$	$y_n$
$x_0$	$x_1$	x	x	x	x	x

↙

(by reset)  
 $x_0$      $A \cdot x_0 + x_1$      $y_{-1}$      $y_0$      $y_1$

$x_2$      $x_3$

↙

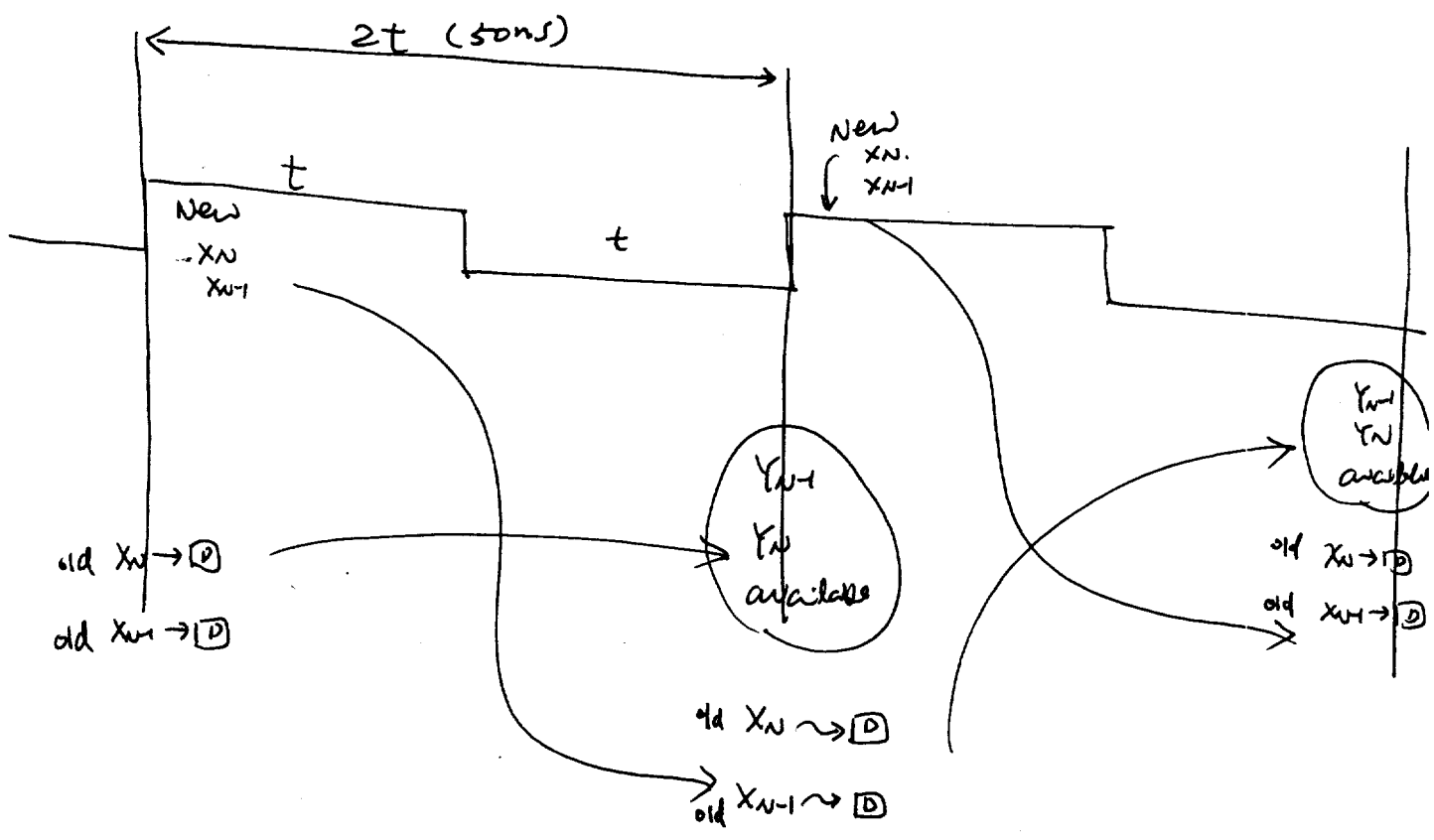
$x_2$      $Ax_2 + x_3$      $y_1$      $y_2 = Ay_1 + x_2$

$$y_3 = x_3 + Ax_2 + A^2 y_1$$

To Get the same throughput, we still use  $2t$

clock cycle (~~50~~ ns)

$\therefore Y_N$  &  $Y_{N-1}$  can be generated simultaneously.



$\therefore$

$C_{eff} = 1.5$

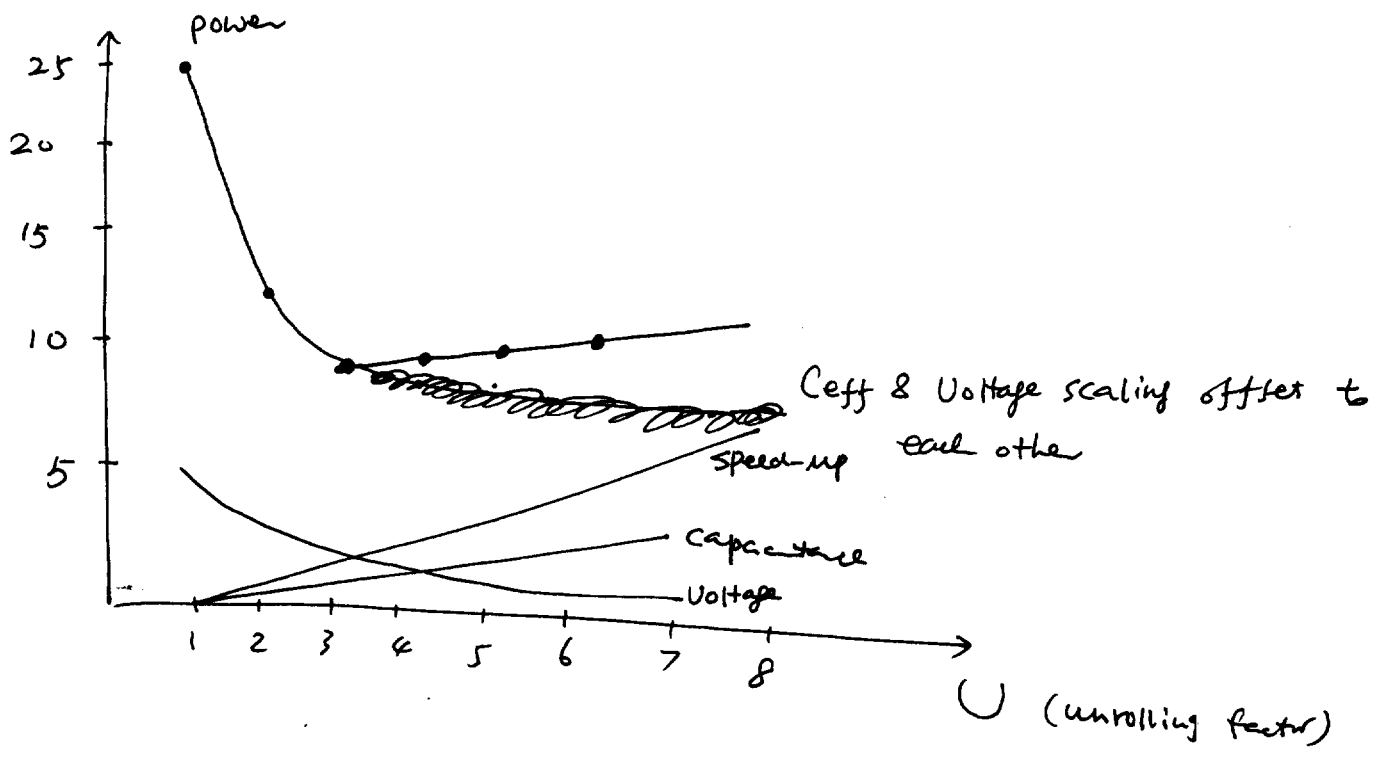
$V_{dd} = 2.9V$  ← delay is changed from 50ns  $\rightarrow$  25ns

Throughput = 1

Critical path = 2

$C_{75}^2$  power = 12.5

50% power reduction



power reduction for unrollif factor 2.3

power increase from factor = 4.