

## Optimizing Delay

- Optimizing delay can be broken into two categories
  - Gate Size selection
  - Transistor sizing
- Gate size selection is done in a standard cell design approach in which you have a library that offers multiple drive strength cells and pick the cells sizes that give the highest speed for a design
  - Current synthesis tools do a good job
- Transistor sizing is done in a custom design in which you size individual transistors during the design process to optimize delay
  - quality depends on individual designer
  - some synthesis help available
  - simulation iteration a tempting option but can be time consuming

BR 6/00

1

## Gate Size Selection

- Many algorithms for gate size selection exist
- One iterative approach is known as the *Tilos* algorithm

Assumptions:

1. Can compute the delay along a path of gates
2. Have multiple gate sizes to choose from

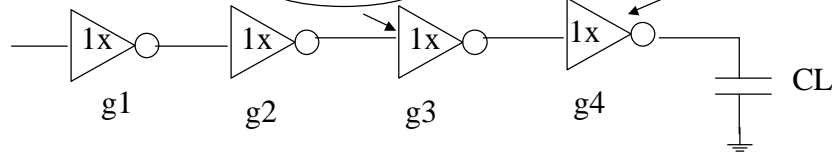
Will yield good results for a path delay

BR 6/00

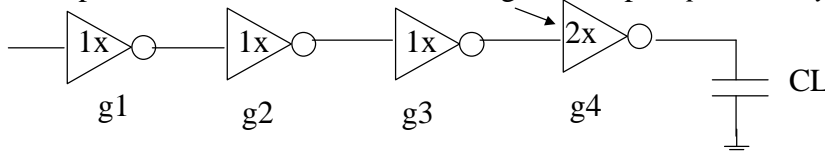
2

## Tilos Algorithm

Step #1: Start with Minimum gate sizes, set *current\_gate* equal to last gate, *driving\_gate* to *current\_gate* - 1.



Step #2a: Increment size of *current\_gate*, compute *path\_delay\_a*

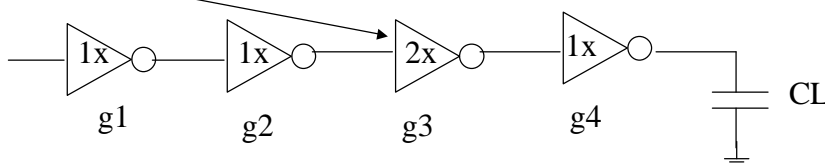


BR 6/00

3

## Tilos Algorithm (cont.)

Step #2b: Restore *current\_gate* size. Increment size of *driving\_gate*, compute *path\_delay\_b*



Step #3: If (*path\_delay\_a* > *path\_delay\_b*) then keep new size of *driving\_gate*, else keep new size of *current\_gate*.

Repeat Steps #2, #3 until no further delay improvement.

Set *current\_gate* to *driving\_gate*, *driving\_gate* to *current\_gate* - 1 and repeat until all gates sized (an exception: the first gate size is considered a FIXED size as in an input buffer).

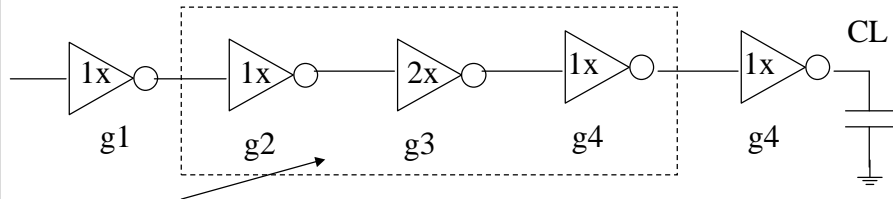
BR 6/00

4

## Some Observations

To save execution time, do have to compute entire path delay.

Computing changes in delay in a 'window' around sized-gate



Compute delay changes here

Also, gate sizes do not have to be exact to get near optimum delay. If optimum gate size happens to be 2.5x, a choice of 2X or 3X will yield good results. This means that rough estimation of gate sizes or transistor sizes can often be satisfactory.

BR 6/00

5

## Rules of Thumb

- Keep fan-in low to keep #transistors in series low (for sub-micron, often  $\leq 3$ ).
- Keep fan-out  $< 5$
- Along a critical path, the minimum delay is achieved if each stage delay is about equal
- Keep rise/fall times about equal

BR 6/00

6

## Estimating Gate Delay, Transistor sizing

- Would be nice to have a “back of the envelope” method of sizing gates/transistors that would be easy to use and would yield reasonable results
- Sutherland/Sproull/Harris book “Logic Effort: Designing Fast CMOS Circuits” introduces a method called “Logical Effort”
- Chapter 1 of the book is posted on the Morgan-Kaufman website ([www.mkp.com](http://www.mkp.com), search for author names)
  - Download this chapter, READ IT!
- We will attempt to apply this method during the semester to the circuits that we will look at.
- Will look at static CMOS application first (these notes taken from that chapter).

BR 6/00

7

## Gate Delay Model

Delay will always be normalized to dimensionless units to isolate effects of fabrication process

$$d_{\text{abs}} = d * \tau$$

Where  $\tau$  is the delay of a minimum sized inverter driving an identical inverter with no gate delay.

Delay of a logic gate is composed of the delay due to *parasitic delay*  $p$  (no load delay) and the delay due to load (*effort delay* or *stage effort*  $f$ )

$$d = f + p$$

BR 6/00

8

## Logical effort, Electrical Effort

The *stage effort*  $f$  (delay due to load) can be expressed as a product of two terms:

$$f = g * h$$

$g$  captures properties of the logic gate and is called the *logical effort*.

$h$  captures properties of the load and is called the *electrical effort*.

On the surface, this does not look different from the model discussed earlier:

$$\text{Gate delay} = \text{no-load delay} + K * C_{\text{load}}$$

Where  $K$  represented the pullup/pulldown strength of the PMOS/NMOS tree.

BR 6/00

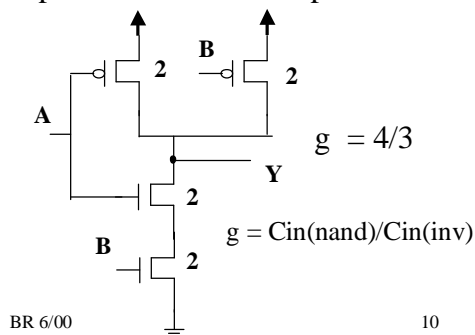
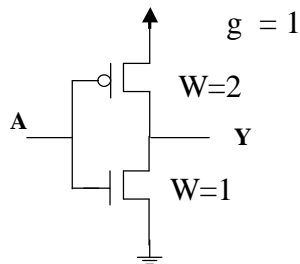
9

## Logical Effort ( $g$ )

In the Sutherland/Sproull model, the logical effort  $g$  factor is normalized to a minimum sized inverter for static CMOS.

So  $g$  for an inverter is equal to 1.

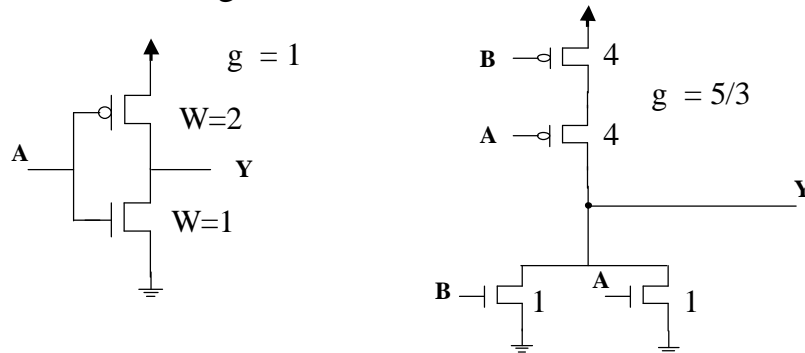
Logical effort  $g$  of other gates represents how much more input capacitance a gate must present to produce the same output current as the inverter.



BR 6/00

10

### Logical Effort inverter vs nor2

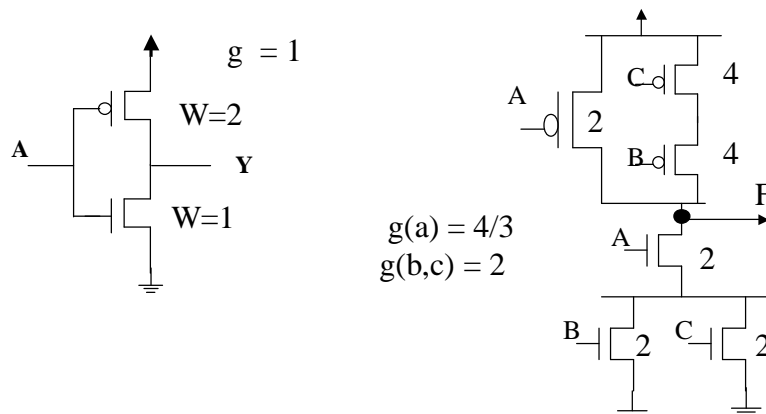


Intuitive result, Nor2  $g$  is higher than Nand2  $g$

BR 6/00

11

### Logical Effort inverter vs Complex gate



Intuitive result, worse case  $g$  of complex gate higher than Nand2 or Nor2.

In general, more inputs, more series transistors, the higher the  $g$  value.

BR 6/00

12

## Logical Effort vs. Electrical Effort

- The value for logical effort  $g$  is independent of transistor size
- The  $g$  value is dependent on number of inputs, and topology
- The electrical effort  $h$  parameter is used to capture the driving capability of the gate via transistor sizing and also the effect of transistor sizes on loading
- Electrical effort  $h$  is defined as  

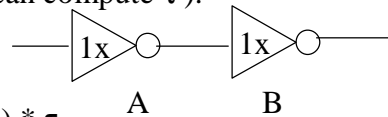
$$C_{out} / C_{in}$$
 where  $C_{out}$  is the load capacitance,  $C_{in}$  is the input capacitance of the gate.
- Note that  $h$  for a gate will reduce as the transistors become wider since  $C_{in}$  increases.

BR 6/00

13

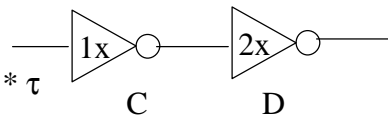
## The Parasitic Delay $p$

- Note that the parasitic delay (no-load)  $p$  is a constant and independent of transistor size; as you increase the transistor sizes the capacitance of the gate/source/drain areas increase also which keeps no-load delay constant
- To measure  $P$  (once  $P$  is known, can compute  $\tau$ ).



$$A\_delay = (g \cdot h + p) \cdot \tau = (1 \cdot 1 + p) \cdot \tau$$

$$\tau = (A\_delay) / (1 + p)$$



$$C\_delay = (g \cdot h + p) \cdot \tau = (1 \cdot 2 + p) \cdot \tau$$

$$C\_delay = (2 + p) (A\_delay) / (1 + p)$$

$$p = (2 \cdot A\_delay - C\_delay) / (C\_delay - A\_delay)$$

BR 6/00

14

## Parasitic Delay of Other Gates

- Normalizing the parasitic delay to that of the inverter can be useful for normalization purposes.
- Some typical values according to Southerland/Sproull:

inverter  $p_{inv} = 1.0$

N-input nand  $n * p_{inv}$

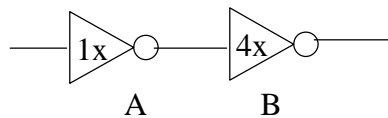
N-input nor  $n * p_{inv}$

Will use these values for example purposes.

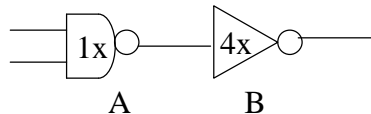
BR 6/00

15

## Delay Estimation



$$\begin{aligned} A\_delay &= g * h + p = 1 * (CinB / CinA) + 1 \\ &= 1 * (4 * CinA / CinA) + 1 = 4 + 1 = 5 \text{ time units} \end{aligned}$$



$$A\_delay = g * h + p = (4/3) * (CinB / CinA) + 2 * 1$$

$$Cin\_B = 4 * 3 = 12. \quad Cin\_A = 4$$

$$A\_delay = (4/3) * (12/4) + 2 = 4 + 2 = 6 \text{ units}$$

Nand2 worse because of higher parasitic delay than inverter.

Note that  $g * h$  term was same for both because NAND2 sized to provide same current drive.

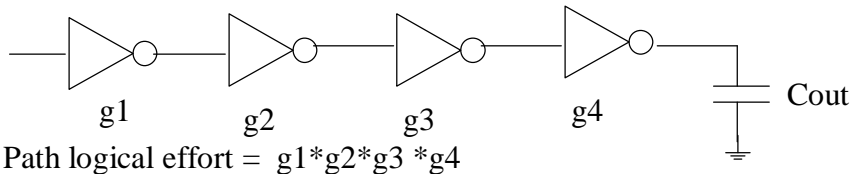
BR 6/00

16



### MultiStage Delay

- Recall rule of thumb that said to balance the delay at each stage along a critical path
- Concepts of logical effort and electrical effort can be generalized to multistage paths



In general, Path logic effort  $G = \prod g(i)$

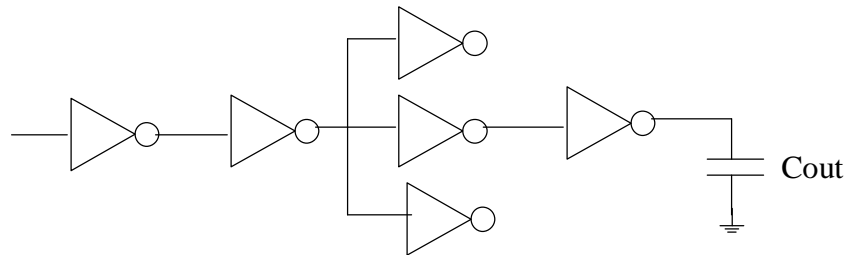
Path electrical effort  $H = C_{out} / C_{in\_first\_gate}$

Must remember that electrical effort only is concerned with effect of logic network on input drivers and output load.

BR 6/00

17

### Off Path Load



Off path load will divert electrical effort from the main path, must account for this. Define a *branching effort*  $b$  as:

$$b = (C_{con\_path} + C_{off\_path}) / C_{con\_path}$$

The branching effort will modify the electrical effort needed at that stage. The branch effort  $B$  of the path is:

$$B = \prod b(i)$$

BR 6/00

18

## Path Effort $F$

Path effort  $F$  is:

$$F = \text{path logic effort} * \text{path branch effort} * \text{path electrical effort} \\ = G * B * H$$

Path branch effort and path electrical effort is related to the electrical effort of each stage:

$$B * H = C_{out}/C_{in} * \prod b(i) = \prod h(i)$$

*Our goal is choose the transistor sizes that effect each stage effort  $h(i)$  in order to minimize the path delay!!!!!!*

BR 6/00

19

## Minimizing Path Delay

The absolute delay will have the parasitic delays of each stage summed together.

However, can *focus on just Path effort  $F$*  for minimization purposes since parasitic delays are constant.

For an N-stage network, *the path delay is least when each stage in the path bears the same stage effort.*

$$f(\min) = g(i) * h(i) = F^{1/N}$$

Minimum achievable path delay

$$D(\min) = N * F^{1/N} + P$$

Note that if  $N=1$ , then  $d = f + p$ , the original single gate equation.

BR 6/00

20

## Choosing Transistor Sizes

Remember that the stage effort  $h(i)$  is related to transistor sizes.

$$f(\min) = g(i) * h(i) = F^{1/N}$$

So

$$h(i) \min = F^{1/N} / g(i)$$

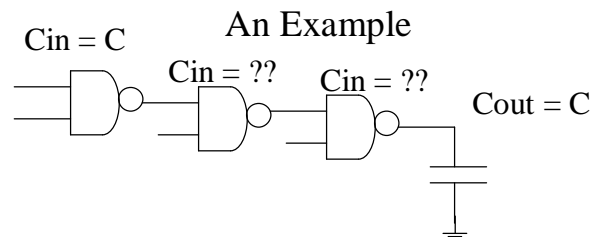
To size transistors, start at end of path, and compute:

$$C_{in}(i) = g_i * C_{out}(i) / f(\min)$$

Once  $C_{in}(i)$  is known, can distribute this among transistors of that stage.

BR 6/00

21



Size the transistors of the nand2 gates for the three stages shown.

$$\text{Path logic effort} = G = g_0 * g_1 * g_2 = 4/3 * 4/3 * 4/3 = 2.37$$

$$\text{Branching effort } B = 1.0 \text{ (no off-path load)}$$

$$\text{Electrical effort } H = C_{out}/C_{in} = C/C = 1.0$$

$$\begin{aligned} \text{Min delay achievable} &= 3 * (G * B * H)^{1/3} + 3 * (2 * p_{inv}) \\ &= 3 * (2.37 * 1 * 1)^{1/3} + 3 * (2 * 1.0) = 10.0 \end{aligned}$$

BR 6/00

22

### An example (cont.)

The effort of each stage will be:

$$f_{\min} = (G \cdot B \cdot H)^{1/3} = (2.37 \cdot 1.0 \cdot 1.0)^{1/3} = 1.33 = 4/3$$

Cin of last gate should equal:

$$\begin{aligned} C_{\text{in last gate (min)}} &= g_i \cdot C_{\text{out (i)}} / f_{\min} \\ &= 4/3 \cdot C / (4/3) = C \end{aligned}$$

Cin of middle gate should equal:

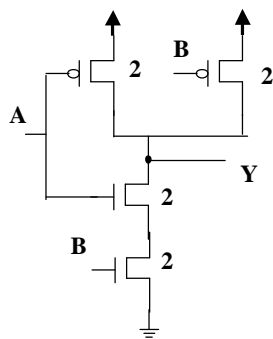
$$\begin{aligned} C_{\text{in middle gate}} &= g_i \cdot C_{\text{in last gate}} / f_{\min} \\ &= 4/3 \cdot C / (4/3) = C \end{aligned}$$

All gates have same input capacitance, distribute it among transistors.

BR 6/00

23

### Transistor Sizes for Example



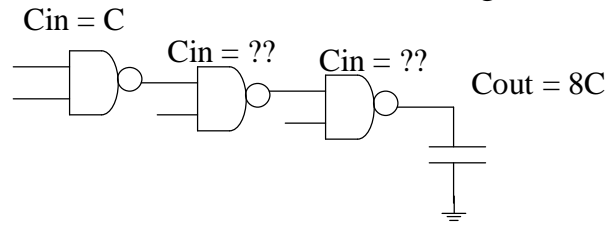
Where gate capacitance of  
 $2 \cdot W \cdot L \text{ Mosfet} = C/2$

Choose W accordingly.

BR 6/00

24

Let Load = 8C, what changes?



Size the transistors of the nand2 gates for the three stages shown.

$$\text{Path logic effort} = G = g_0 * g_1 * g_2 = 4/3 * 4/3 * 4/3 = 2.37$$

Branching effort B = 1.0 (no off-path load)

$$\text{Electrical effort } H = C_{out}/C_{in} = 8C/C = 8.0$$

$$\begin{aligned} \text{Min delay achievable} &= 3 * (G * B * H)^{1/3} + 3 (2 * p_{inv}) \\ &= 3 * (2.37 * 1 * 8)^{1/3} + 3 (2 * 1.0) = 14.0 \end{aligned}$$

BR 6/00

25

### 8C Load Example (cont.)

The effort of each stage will be:

$$f_{min} = (G * B * H)^{1/3} = (2.37 * 1.0 * 8)^{1/3} = 2.67 = 8/3$$

Cin of last gate should equal:

$$\begin{aligned} C_{in \text{ last gate (min)}} &= g_i * C_{out (i)} / f(\min) \\ &= 4/3 * 8C / (8/3) = 4C \end{aligned}$$

Cin of middle gate should equal:

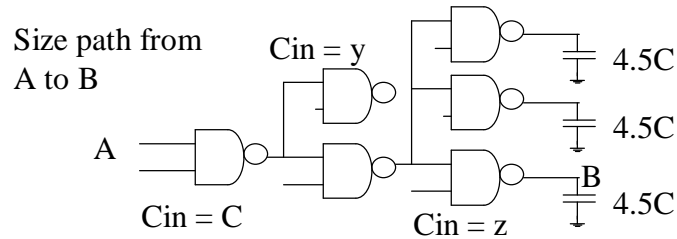
$$\begin{aligned} C_{in \text{ middle gate}} &= g_i * C_{in \text{ last gate}} / f(\min) \\ &= 4/3 * 4C / (8/3) = 2C \end{aligned}$$

Note that each stage gets progressively larger, as is typical with a multi-stage path driving a large load.

BR 6/00

26

### Example 1.6 from Chapter 1



Path logic effort  $G = g_0 * g_1 * g_2 = 4/3 * 4/3 * 4/3 = 2.37$

Branch effort, 1<sup>st</sup> stage =  $(y+y)/y = 2$ .

Branch effort, 2<sup>nd</sup> stage =  $(z+z+z)/z = 3$

Path Branch effort  $B = 2 * 3 = 6$ .

Path electrical effort  $H = C_{out}/C_{in} = 4.5C/C = 4.5$

Path stage effort =  $F = G*B*H = 2.37*6*4.5 = 64$ .

Min delay =  $N(F)^{1/N} + P = 3*(64)^{1/3} + 3(2\text{pinv}) = 18.0 \text{ units}$

BR 6/00

27

### Example 1.6 from Chapter 1 (cont)

Stage effort of each stage should be:

$$f(\min) = (F)^{1/N} = (GBH)^{1/N} = (64)^{1/3} = 4$$

Determine  $C_{in}$  of last stage:

$$C_{in}(z) = g * C_{out} / f(\min) = 4/3 * 4.5C / 4 = 1.5 C$$

Determine  $C_{in}$  of middle stage:

$$C_{in}(y) = g * (3*C_{in}(z)) / f(\min) = 4/3 * (3*1.5C) / 4 = 1.5C$$

Is first stage correct?

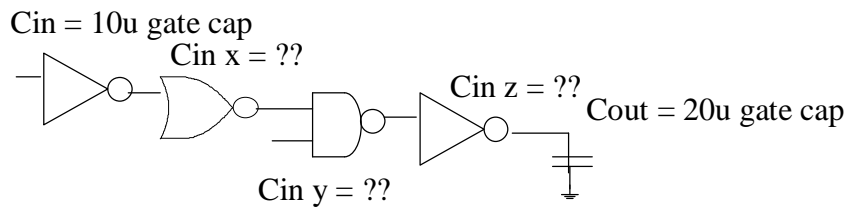
$$C_{in}(A) = g * (2*C_{in}(y)) / f(\min) = 4/3 * (2*1.5C) / 4 = C.$$

Yes, self-consistent.

BR 6/00

28

### Example 1.10 from Chapter 1



Path logic effort  $G = g_0 * g_1 * g_2 * g_3 = 1 * 5/3 * 4/3 * 1 = 20/9$

Path Branch effort  $B = 1$

Path electrical effort  $H = C_{out}/C_{in} = 20/10 = 2$

Path stage effort  $F = G * B * H = (20/9) * 1 * 2 = 40/9$

For Min delay, each stage has effort  $(F)^{1/N} = (40/9)^{1/4} = 1.45$

$z = g * C_{out}/f(\min) = 1 * 20 / 1.45 = 14$

$y = g * C_{in}(z)/f(\min) = 4/3 * 14 / 1.45 = 13$

$x = g * C_{in}(y)/f(\min) = 5/3 * 13 / 1.45 = 15$

BR 6/00

29

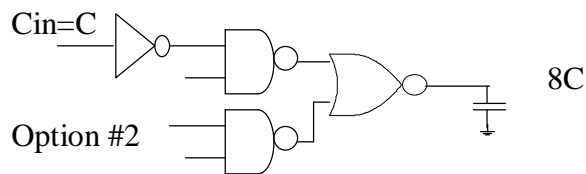
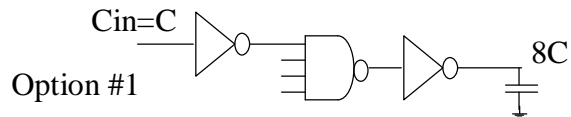
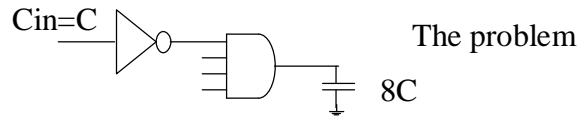
### Misc Comments

- Note that you never size the first gate. This gate size is assumed to be fixed (same as in the Tilos algorithm) – if you were allowed to size this gate you find that the algorithm would want to make it as large as possible.
- This is an estimation algorithm. The author claims that sizing a gate by 1.5x too big or two small still results in path delay within 5% of minimum.

BR 6/00

30

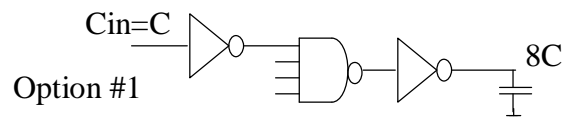
## Evaluating different Structure options



BR 6/00

31

## Option #1



$$\text{Path logic effort } G = g_0 * g_1 * g_2 = 1 * 6/3 * 1 = 2$$

$$\text{Path Branch effort } B = 1$$

$$\text{Path electrical effort } H = C_{out}/C_{in} = 8C/C = 8$$

$$\text{Path stage effort } = F = G * B * H = 2 * 1 * 8 = 16$$

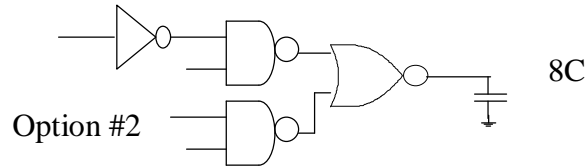
$$\begin{aligned} \text{Min delay: } &= N * (F)^{1/N} + P \\ &= 3 * (16)^{1/3} + (p_{inv} + 4 * p_{inv} + p_{inv}) \\ &= 3 * (2.5) + 6 = 13.5 \end{aligned}$$

BR 6/00

32



### Option #2



Path logic effort  $G = g_0 * g_1 * g_2 = 1 * 4/3 * 5/3 = 20/9$

Path Branch effort  $B = 1$

Path electrical effort  $H = C_{out}/C_{in} = 8C/C = 8$

Path stage effort  $= F = G * B * H = 20/9 * 1 * 8 = 160/9$

Min delay:  $= N * (F)^{1/N} + P$

$= 3 * (160/9)^{1/3} + (p_{inv} + 2 * p_{inv} + 2 * p_{inv})$

$= 3 * 2.6 + 5 = 12.8$

Option #2 appears to be better than Option #1, by a slight margin.

BR 6/00

33