`inst.eecs.berkeley.edu/~eecs151`

# EECS151 : Introduction to Digital Design and ICs
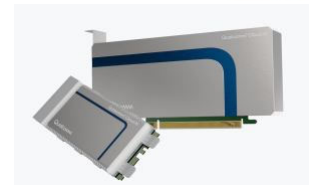
## Lecture 13 – CMOS Logic

### Bora Nikolić

**Qualcomm Takes on Nvidia for MLPerf Inference Title**

October 1, 2021, EETimes, Sally Ward-Foxton - The latest round of MLPerf AI inference benchmark scores are in. Nvidia has dominated both MLPerf training and inference results since the beginning, but in this round Qualcomm appears to be close on Nvidia's tail when it comes to data center/edge server inference.

Qualcomm Cloud AI100 PCIe and M.2 cards (Source: Qualcomm)

EETimes

1

---

## Review

- CMOS process is used for producing chips
  - Planar bulk process used up to 28nm node
  - finFET, FDSOI used below the 22nm node

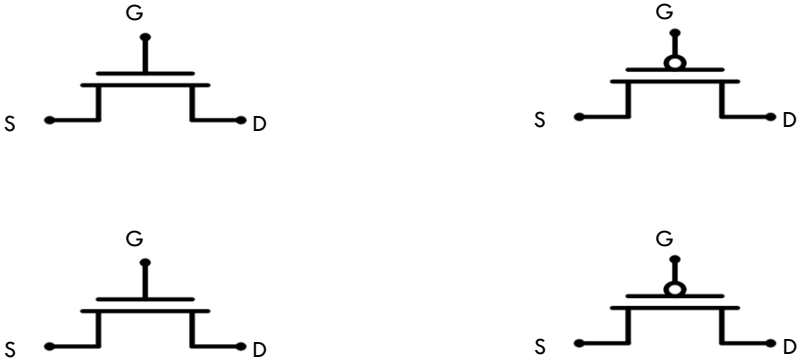- Switch-level abstraction for MOS transistors

2

MOS Switch

3



MOS Switch

4

## CMOS Inverter



Nikolić Fall 2021   5   Berkeley

5

---

## CMOS Inverter

- Simple DC behavior
  - Schematic

  

  $W_p/L$

  $W_n/L$

- Switch model

  

  $V_{in} = V_{DD}$        $V_{in} = 0$

  $$V_{OL} = 0$$
  $$V_{OH} = V_{DD}$$

Nikolić Fall 2021   6   Berkeley

6

3

## Voltage Transfer Characteristic (VTC)



$V_A = V_{GS,n} = V_{DD} - V_{SG,p}$

$I_{DS,n} = I_{SD,p}$

$V_{Out} = V_{DS,n} = V_{DD} - V_{SD,p}$

EECS151 L12 CMOS2  Nikolić Fall 2021  7

7

## Voltage Transfer Characteristic (VTC)



• Can we change switching point ($V_A$ for which $V_{out} = V_{DD}/2$)?

EECS151 L12 CMOS2  Nikolić Fall 2021  8

8

4

# Digital Circuits

- One logic representation

$$Out = \overline{A}$$

Truth table

| A | Out |
|---|-----|
| 0 | 1   |
| 1 | 0   |

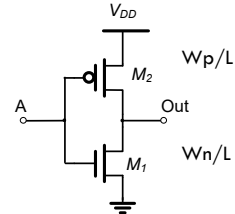- Multiple libraries

  - Layouts
    - Number of metal 'tracks'
    - More tracks, faster, but larger
    - Less tracks – more compact, but slower

  - Transistor thresholds ($V_{Th}$) (for each track height):
    - Regular (RVT)
    - Low (LVT)
      - Faster, higher power
      - Slower, lower power
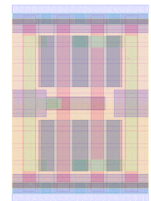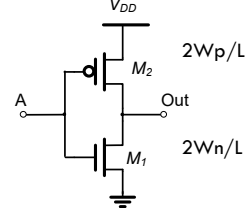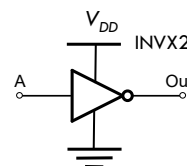    - High (HVT)
  - Transistor lengths

- Multiple gate sizes within a library

  - Symbol

    $V_{DD}$ INVX1

    A     Out

    $V_{DD}$ INVX2

    A     Out

    INVX3,

    INVX4,…

  - Schematic

    $V_{DD}$

    A    $M_2$ Wp/L

    Out

    $M_1$ Wn/L

    $V_{DD}$

    A    $M_2$ 2Wp/L

    Out

    $M_1$ 2Wn/L

  - Layout

9

# Administrivia

- Homework 5 due this week

- Lab 6 (last) this week

- Projects start next week

10

# CMOS Logic

11

---

# Building logic from switches

**Series**    X ———A——— B——— Y     **AND**

$Y = X$ if A AND B

**Parallel**    X ——[A / B]—— Y     **OR**

$Y = X$ if A OR B

(output undefined if condition not true)

12

# Logic using inverting switches

**Series**



A                    B

X —⟋— • —⟋— Y

**NOR**

$Y = X$ if $\overline{A}$ AND $\overline{B}$

$= \overline{A + B}$

**Parallel**



A

X • —⟋— • Y

B

**NAND**

$Y = X$ if $\overline{A}$ OR $\overline{B}$

$= \overline{AB}$

(output undefined if condition not true)

13

---

# Static Complementary CMOS

$V_{DD}$

$In_1$
$In_2$  **Pull-up Network**     Inverting switches
$In_N$

→ $F(In_1, In_2, \ldots In_N)$

$In_1$
$In_2$  **Pull-down Network**   Non-Inverting switches
$In_N$

PUN and PDN are dual logic networks
PUN and PDN functions are complementary
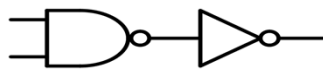
Dual Graphs

14

## Complementary CMOS Logic Style

❑ PUN is the **dual** to PDN
(can be shown using DeMorgan's Theorems)

$$\overline{A + B} = \overline{A}\,\overline{B}$$

$$\overline{AB} = \overline{A} + \overline{B}$$
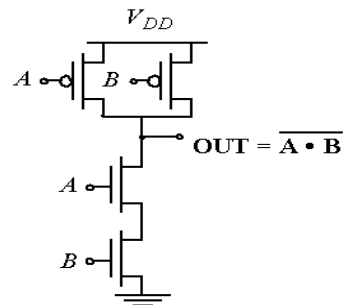
❑ Static CMOS gates are always inverting



**AND = NAND + INV**

15

---

## Example Gate: NAND



| A | B | Out |
|---|---|-----|
| 0 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

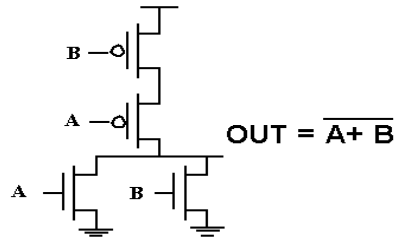Truth Table of a 2 input NAND gate

$V_{DD}$

$\text{OUT} = \overline{A \cdot B}$

❑ PDN: G = AB ⟹ Conduction to GND
❑ PUN: F = $\overline{A}$ + $\overline{B}$ = $\overline{AB}$ ⟹ Conduction to $V_{DD}$

❑ $\overline{G(In_1, In_2, In_3, \ldots)} \equiv F(\overline{In_1}, \overline{In_2}, \overline{In_3}, \ldots)$

16

8

# Example Gate: NOR

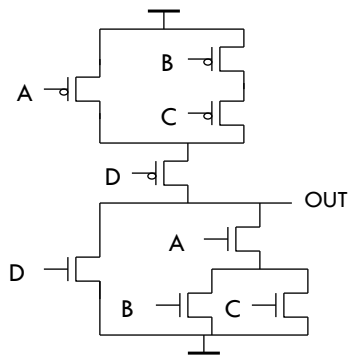| A | B | Out |
|---|---|-----|
| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 0 |

**Truth Table of a 2 input NOR gate**

$$OUT = \overline{A + B}$$

17

# Complex CMOS Gate

$$OUT = \overline{D + A \cdot (B + C)}$$

- Note: In scaled processes max #inputs is 3-4
  - Max stack height is 2 or 3

18

9

# Stick Diagrams

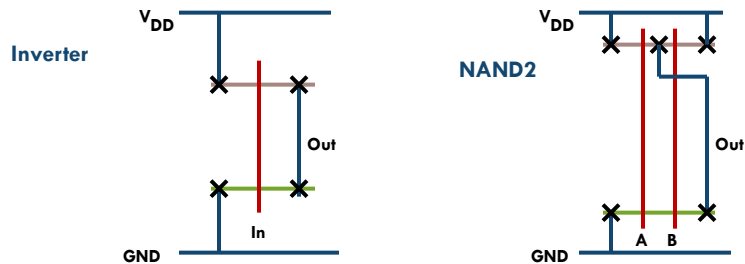**Contains no dimensions**
**Represents relative positions of transistors**

Inverter

NAND2

Nikolić Fall 2021

19

19

# Stick Diagrams

$X = \overline{C \cdot (A + B)}$

Logic Graph

PUN

PDN

Nikolić Fall 2021

20

20

10

## Two Versions of $\overline{C \cdot (A + B)}$

Nikolić Fall 2021

21

---

## Consistent Euler Path



$X = \overline{C \cdot (A + B)}$
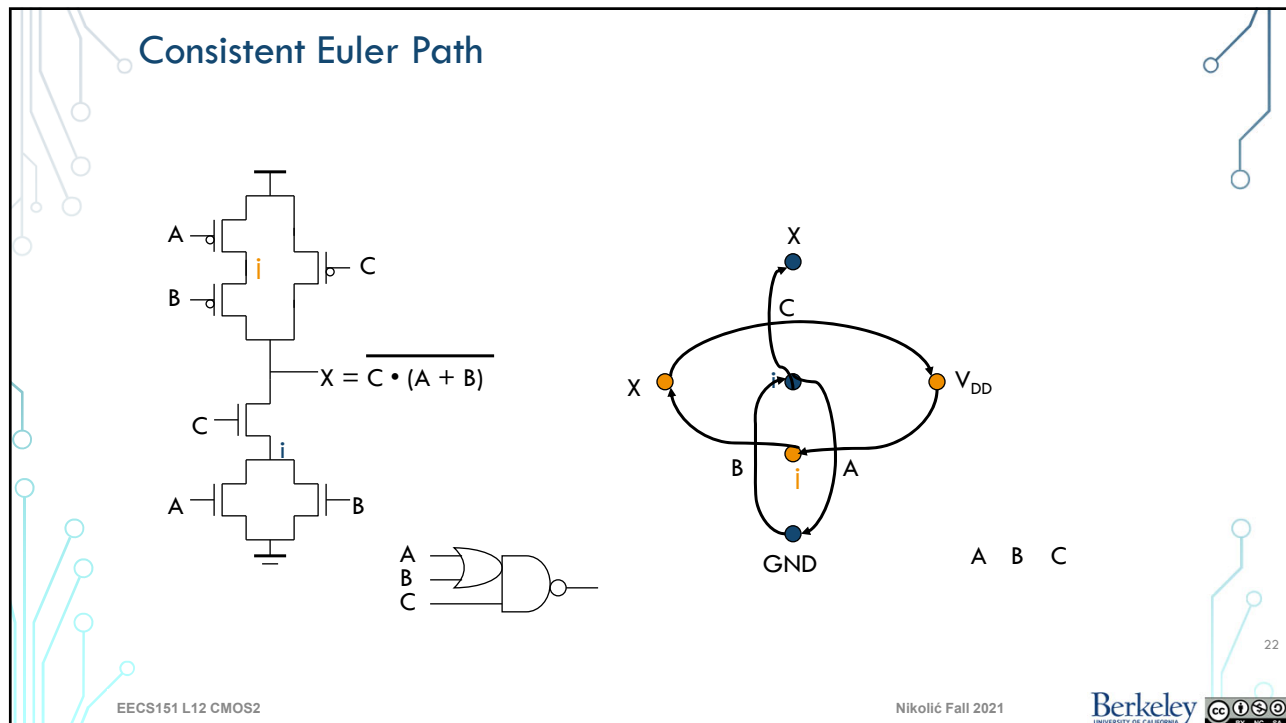
Nikolić Fall 2021

22

# OAI22 Logic Graph



$$X = \overline{(A+B) \cdot (C+D)}$$

PUN

PDN

Nikolić Fall 2021

# Example: $x = \overline{ab+cd}$



(a) Logic graphs for $\overline{(ab+cd)}$

(b) Euler Paths $\{a\ b\ c\ d\}$

(c) stick diagram for ordering $\{a\ b\ c\ d\}$

Nikolić Fall 2021

## Switch Limitations

$V_{DD}$

S

Good 1

$0 \to V_{DD}$

D

$C_L$

$V_{DD}$

D

Bad 1

$V_{DD}$

$V_{GS}$

S

$0 \to V_{DD} - V_{Tn}$

$C_L$

$V_{DD} \to 0$

D

$C_L$

$V_{DD}$

S

Good 0
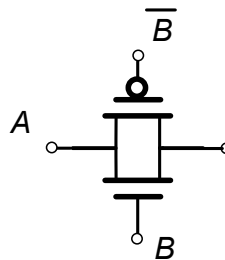
$V_{GS}$

S

$V_{DD} \to |V_{Tp}|$

$C_L$

D

Bad 0

25

---

## Transmission Gate

- Transmission gates are the way to build "switches" in CMOS.
- In general, both transistor types are needed:
  - ❏ nFET to pass zeros.
  - ❏ pFET to pass ones.
- The transmission gate is 'non-isolating'.

$\overline{B}$

$A$

$B$

26

## Transmission-Gate Multiplexer

• Implementation



| Sel | Y |
|-----|---|
| 0 | A |
| 1 | B |

```
module comb(input a, b, sel,
    output reg y);
  always @(*) begin
    case (sel)
      1b'0: y <= a;
      1b'1: y <= b;
    eendcase
  end
endmodule
```

27

## CMOS Multiplexer

| Sel | Y |
|-----|---|
| 0 | A |
| 1 | B |

28

14

CMOS Sizing
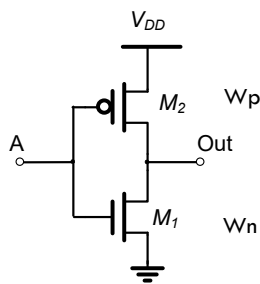
29

# Transistor Sizing

- Optimal Wp/Wn



- In the past, Wp > Wn (see Rabaey, 2$^{nd}$ ed)

- In modern processes (finFET), Wp = Wn

30

15

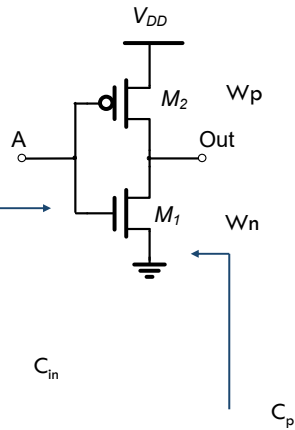## Gate Sizing

- Doubling the gate size (by doubling Ws):



- Doubles $C_{in}$

- Halves equivalent gate resistance

- Doubles $C_p$

---



# CMOS Delay

## Inverter Delay

- How to time this?



$R_{eq}, C_p$

$C_{in}$

- Each gate has an $R_{eq}$ and drives $C_{in}$ of the next gate

$R_{eq}, C_p$

In — Out

$C_{load} = C_{in}$ (next gate)

$V_{DD}$

$R_{eq,p}$

Out

$C_p + C_L$

$R_{eq,n}$

33

---

## Inverter Delay

- High-to-low

$V_{DD}$

In

$R_{eq,p}$

Out

$C_p + C_L$

$R_{eq,n}$

Out

$C_p + C_L$

$R_{eq,n}$

In

In — Out

$V_{in}$

$V_{DD}$

0

$V_{out}$

$V_{DD}$

$V_{DD}/2$

0   $t_p$

$$V_{out} = VDD \, e^{\frac{-t}{\tau}}$$

$$t_{p,HL} = (\ln 2)\tau = 0.7 \, R_{eq,n}(C_p + C_L)$$

$$\tau = R_{eq,n}(C_p + C_L)$$

34

## Inverter Delay

In → Out

$$V_{out}=VDD\left(1-e^{\frac{-t}{\tau}}\right)$$

$$t_{p,LH}= (\ln 2)\tau = 0.7\, R_{eq,p}(C_p+C_L)$$

EECS151 L12 CMOS2 — Nikolić Fall 2021 — 35 — Berkeley

35

## Capacitances

- $C_{in}$ is largely set by the gate cap
  - ~WL
  - $2\times W = 2\times C_{in}$
  - It is non-linear, but we will ignore that

- $C_p$ is largely set by the drain cap
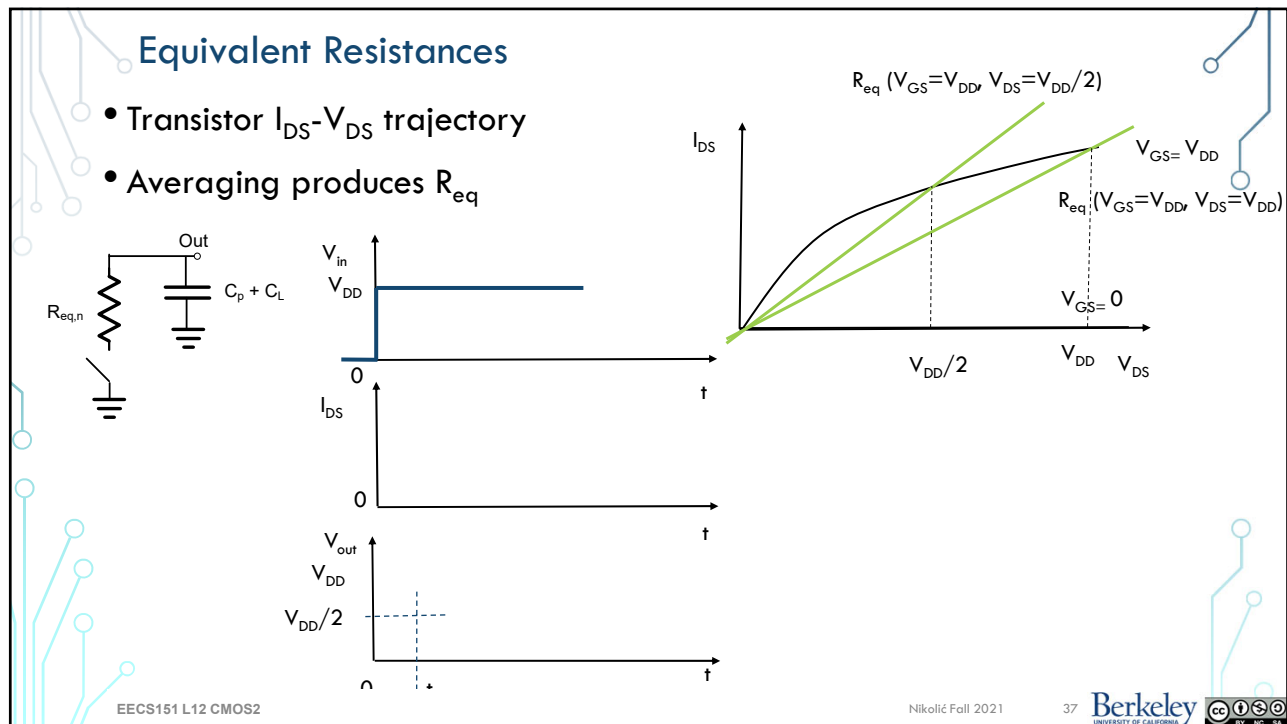  - ~W (drain area/perimeter)
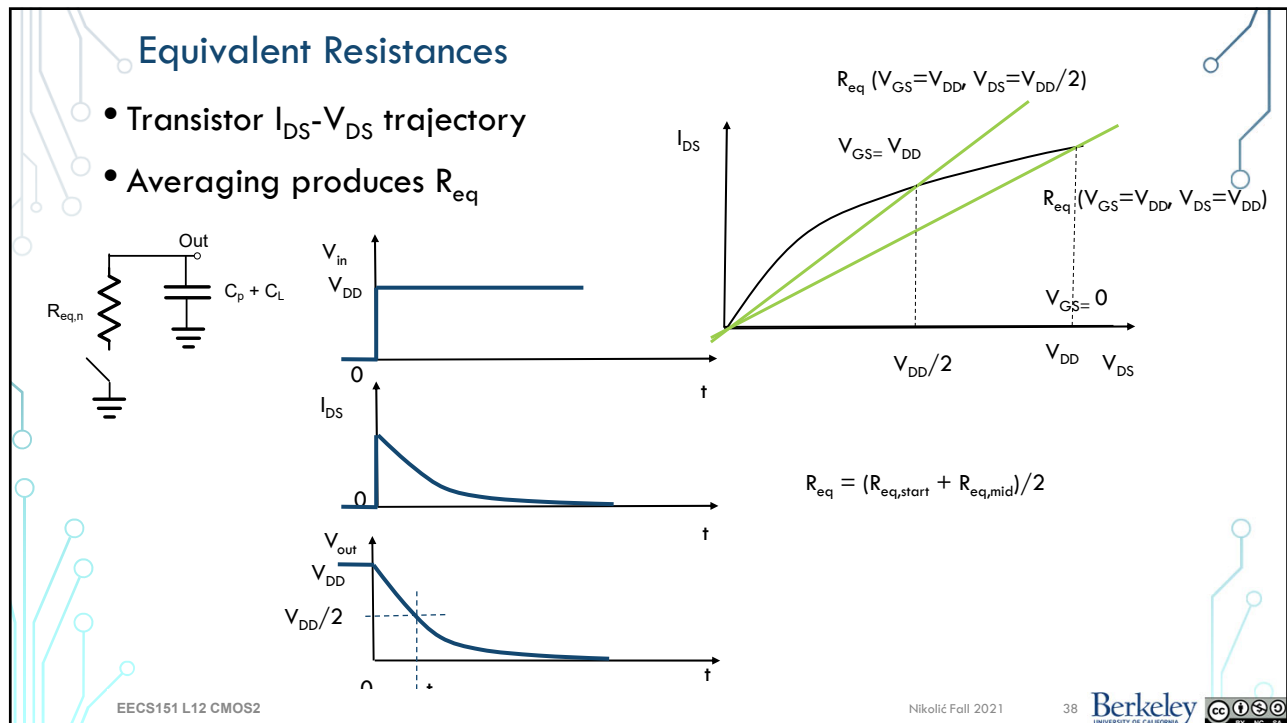  - $2\times W = 2\times C_p$

$$C_p = \gamma C_{in}$$

Gate (G), Source (S), Drain (D), L, W

EECS151 L12 CMOS2 — Nikolić Fall 2021 — 36 — Berkeley

36

18

**Equivalent Resistances**

- Transistor $I_{DS}$-$V_{DS}$ trajectory
- Averaging produces $R_{eq}$

37



**Equivalent Resistances**

- Transistor $I_{DS}$-$V_{DS}$ trajectory
- Averaging produces $R_{eq}$

$$R_{eq} = (R_{eq,start} + R_{eq,mid})/2$$

38

# Impact of Rise/Fall times

- Impacts the $I_{DS}$-$V_{DS}$ trajectory

39

# Impact of Rise/Fall times

- Impacts the $I_{DS}$-$V_{DS}$ trajectory

40

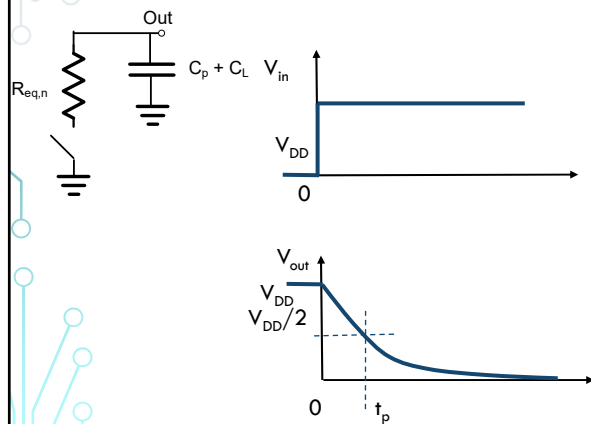## Impact of Supply Voltage

- Lowering VDD, slows down the circuit

41

## Quiz: Inverter Delay

- If we double the load capacitance, assuming the default Vout shown in blue, which of the following waveforms shows the new Vout?



**A**

**B**

**C**

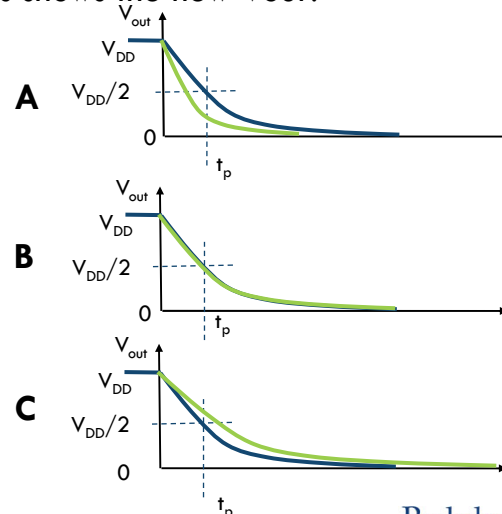42

21

# Summary

- CMOS allows for convenient switch level abstraction

- CMOS pull-up and pull-down networks are complementary
  - Graph models for CMOS gates

- Transistor sizing affects gate performance

- Delay is a linear function of R and C

43