inst.eecs.berkeley.edu/~eecs251b

# EECS251B : Advanced Digital Circuits and Systems

## Lecture 20 – Low Power Design

### Borivoje Nikolić, Vladimir Stojanović, Sophia Shao

**IEEE MICRO, Nov/Dec 2021**
Microprocessor at 50: Looking Back and Looking Forward Special issue on 50 years of a microprocessor

Advertisement in the Electronics News Weekly in November 1971 announcing the Intel 4004.

EECS251B L20 LOW-POWER DESIGN

Berkeley

1

---

## Recap

- Power is a primary design constraint
  - In both cloud and edge systems

- Excess performance traded off for power savings

2

# Architectural Optimizations

3

---

## Optimal Processors

- Processors used to be optimized for performance
  - Optimal logic depth was found to be 8-11 FO4 delays in superscalar processors
  - 1.8-3 FO4 in sequentials, rest in combinatorial
    - Kunkel, Smith, ISCA'86
    - Hriskesh, Jouppi, Farkas, Burger, Keckler, Shivakumar, ISCA'02
    - Harstein, Puzak, ISCA'02
    - Sprangle, Carmean, ISCA'02
- But those designs have very high power dissipation
  - Need to optimize for both performance and power/energy

4

## From System View: What is the Optimum?

- How do sensitivities relate to more traditional metrics:
  - Power per operation (MIPS/W, GOPS/W, TOPS/W)
  - Energy per operation (Joules per op)
  - Energy-delay product
- Can be reformatted as a goal of optimizing power x delay$^n$
  - $n = 0$ – minimize power per operation
  - $n = 1$ – minimize energy per operation
  - $n = 2$ – minimize energy-delay product
  - $n = 3$ – minimize energy-$(delay)^2$ product

5

5

## Optimization Problem

- Set up optimization problem:
  - Maximize performance under energy constraints
  - Minimize energy under performance constraints
- Or minimize a composite function of $E^n D^m$
  - What are the right n and m?
- $n = 1$, $m = 1$ is EDP – improves at lower $V_{DD}$
- $n = 1$, $m = 2$ is invariant to $V_{DD}$
  - $E \sim C V_{DD}^2$
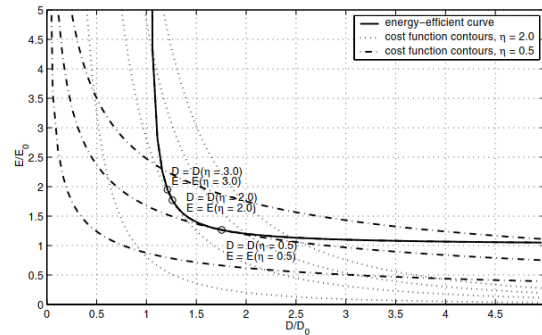  - $D \sim 1/V_{DD}$

6

6

3

# Hardware Intesnity

- Introduced by Zyuban and Strenski in 2002.

- Measures where is the design on the Energy-Delay curve

- Parameter in cost function optimization

$$F_c = (E/E_0)(D/D_0)^\eta \qquad 0 \le \eta < +\infty,$$

$$\eta = -\left.\frac{D \partial E}{E \partial D}\right|_v$$

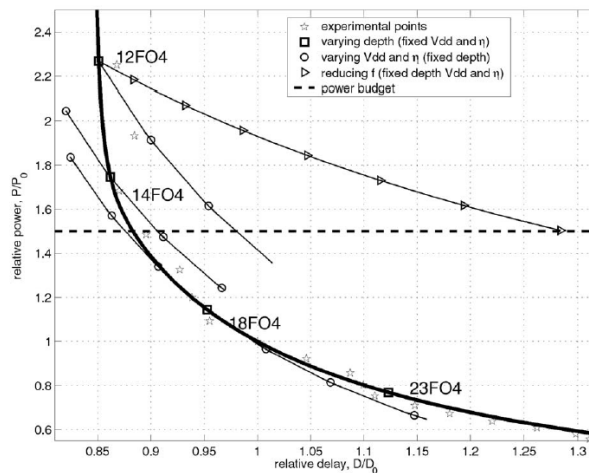**Slope of the optimal E-D curve at the chosen design point**

7

---

# Optimum Across Hierarchy Layers



Zyuban et al, TComp'04

**Optimal logic depth in pipelined processors is ~18FO4**
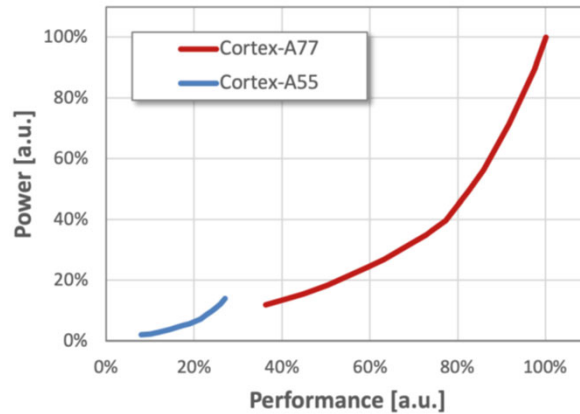Relatively flat in the 16-22FO4 range

8

## Architectural Tradeoffs

- H, Mair, ISSCC'20
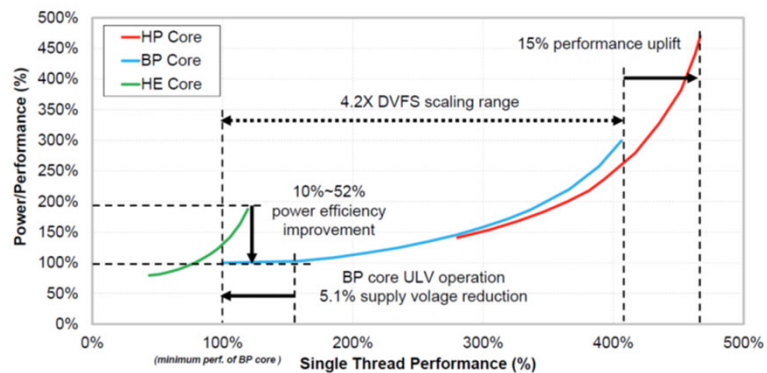
9

## Architectural Tradeoffs: Tri-Gear

- HP: High performance (ARM Cortex A78, optimized for speed, 3.0GHz)
- BP: Balanced performance (ARM Cortex A78, optimized for power, 2.6GHz))
- HE: High efficiency (ARM A55, 2.0GHz)

10

## Announcements

- Quiz 2 today

- Homework 3 due next week

11

11



## Circuit-Level Tradeoffs

12

12

## Alpha-Power Based Delay Model



In ———▷o——▷o— - - - - ▷o— Out

1      2        N    $C_L$

$$t_{pi} = \frac{K_d V_{DD}}{(V_{DD} - V_{Th})^\alpha}\left(1 + \frac{C_{L,i}}{C_{in,i}}\right)$$

$$D = \sum t_{pi} = \sum \frac{K_d V_{DD}}{(V_{DD} - V_{Th})^\alpha}\left(1 + \frac{W_{L,i}}{W_{in,i}}\right)$$
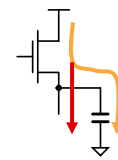
EECS251B L20 LOW-POWER DESIGN      13

13

## Energy Models

♦ **Switching**

$$E_{Sw} = \alpha_{0\to1}\left(C_{L,i} + C_{int,i}\right)V_{DD}^2$$

♦ **Leakage**

$$E_{Lk} = W_{In} I_0 e^{\frac{-(V_{Th} - \gamma V_{DD})}{n V_t}} V_{DD} D$$

EECS251B L20 LOW-POWER DESIGN      14

14

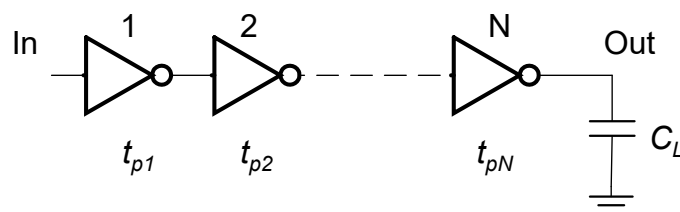## Sizing, Supply, Threshold Optimization

- Transistor sizing can yield large power savings with small delay penalties
  - Gate sizing
  - Beta-ratio adjustments $\qquad \beta = Wp/Wn$
  - (Stack resizing)
- Supply voltage affects both active and leakage energy
- Threshold voltage affects primarily the leakage

15

---

## Apply to Sizing of an Inverter Chain



Unconstrained energy: find min D = $\Sigma t_{pi}$

$$C_{gin,j} = \sqrt{C_{gin,j-1}C_{gin,j+1}} \qquad W_j = \sqrt{W_{j-1}W_{j+1}}$$

Constrained energy: find min D, under E < $E_{max}$
Where E = $\Sigma e_i$

16

8

## Constrained Optimization

- Find min($D$) subject to $E = E_{max}$
  - *Constrained function minimization*

- E.g. Lagrange multipliers                    Or dual:

$$\Lambda(x) = D(x) + \lambda(E(x) - E_{max})$$    $$K(x) = E(x) + \lambda(D - D_{max})$$

$$\frac{\partial \Lambda}{\partial x} = 0$$

- Can solve analytically for $x = W_i, V_{DD}, V_{Th}$

EECS251B L20 LOW-POWER DESIGN                                                    17

17

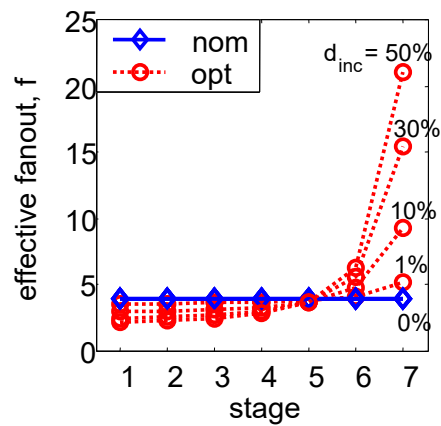## Inverter Chain: Sizing Optimization

EECS251B L20 LOW-POWER DESIGN                                                    18

18

## Inverter Chain: Sizing Optimization



$$W_j = \sqrt{\frac{W_{j-1}W_{j+1}}{1 + \lambda W_{j-1}}}$$

[Ma, Franzon, *IEEE JSSC*, 9/94]

$$\lambda = -\frac{2KV_{DD}^2}{\tau_{nom}S_W}$$

$e_i$ – energy per stage
$f_i$ – fanout per stage

$$S_W \propto \frac{e_j}{f_j - f_{j-1}}$$

Stojanovic, ICCAD'02

- **Variable taper achieves minimum energy**
- **Reduce number of stages at large d$_{inc}$**

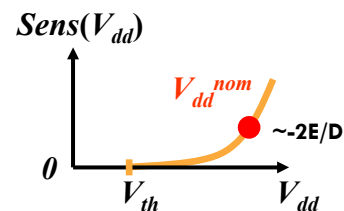EECS251B L20 LOW-POWER DESIGN

19

19

## Sensitivity to Sizing and Supply

- Gate sizing ($W_i$)

$$-\frac{\partial E_{sw}/\partial W_j}{\partial D/\partial W_j} = \frac{e_j}{\tau_{nom}(f_j - f_{j-1})}$$

∞ for equal $f_{eff}$
($D_{min}$)

- Supply voltage ($V_{dd}$)

$$-\frac{\partial E_{sw}/\partial V_{DD}}{\partial D/\partial V_{DD}} = \frac{E_{sw}}{D} 2 \frac{1 - x_v}{\alpha - 1 + x_v}$$



$$x_v = (V_{Th} + \Delta V_{Th})/V_{dd}$$
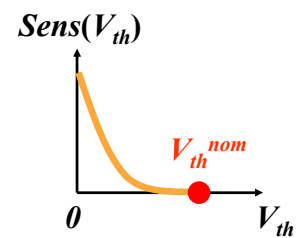
EECS251B L20 LOW-POWER DESIGN

20

20

10

## Sensitivity to $V_{th}$

- Threshold voltage ($V_{th}$)

$$-\frac{\partial E/\partial \Delta V_{Th}}{\partial D/\partial \Delta V_{th}} = P_{Lk}\left(\frac{V_{DD} - V_{Th} - \Delta V_{Th}}{\alpha n V_t} - 1\right)$$

**Low initial leakage**

$\Rightarrow$ **speedup comes for "free"**

21

---



## Scaling Supplies

22

## Reducing $V_{dd}$

$$P \times t_d = E_t = C_L * V_{dd}^2$$

$$\frac{E_{(Vdd=2)}}{E_{(Vdd=5)}} = \frac{(C_L) * (2)^2}{(C_L) * (5)^2}$$

$$E_{(Vdd=2)} \approx 0.16 \; E_{(Vdd=5)}$$

- **Strong function of voltage ($V^2$ dependence).**
- **Relatively independent of logic function and style.**
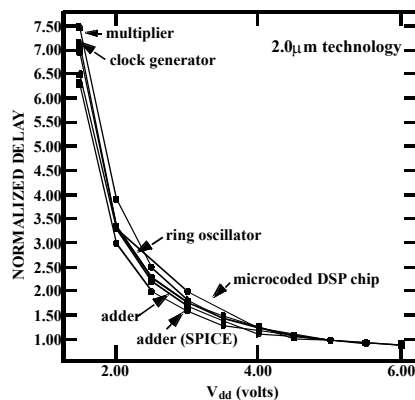- **Power Delay Product Improves with lowering $V_{DD}$.**

Chandrakasan, JSSC'92

EECS251B L20 LOW-POWER DESIGN

23

23

## Lower $V_{DD}$ Increases Delay

$$T_d = \frac{C_L * V_{dd}}{I}$$

$$I \sim (V_{dd} - V_t)^2$$

$$\frac{T_{d(Vdd=2)}}{T_{d(Vdd=5)}} = \frac{(2) * (5 - 0.7)^2}{(5) * (2 - 0.7)^2}$$

$$\approx 4$$

- **Relatively independent of logic function and style.**

EECS251B L20 LOW-POWER DESIGN
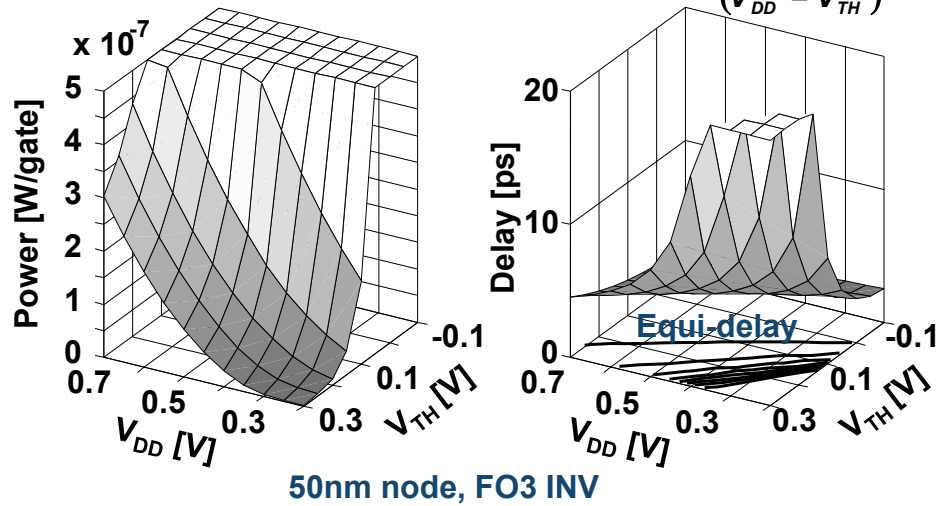
24

24

## Trade-off Between Power and Delay

$$Power = a \cdot f \cdot C \cdot V_{DD}^2 + I_0 \cdot 10^{-\frac{V_{TH}}{s}} \cdot V_{DD}$$

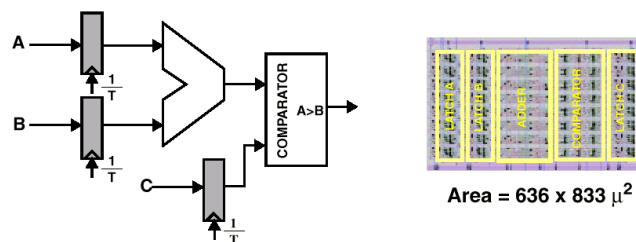$$Delay \; \propto \; \frac{C \cdot V_{DD}}{(V_{DD} - V_{TH})^{1.3}}$$

**50nm node, FO3 INV**

25

---

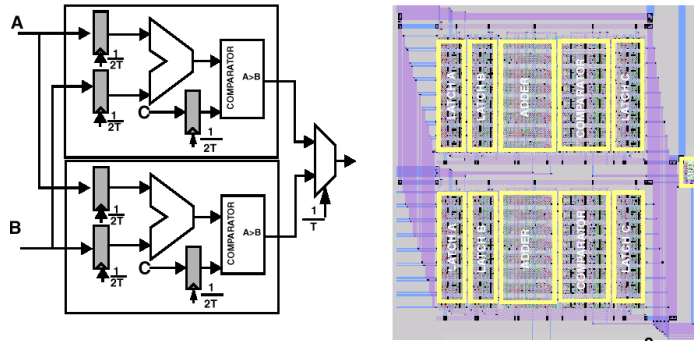## Architecture Trade-off for Fixed-rate Processing Reference Datapath

**Area = 636 x 833 $\mu^2$**

- Critical path delay $\Rightarrow T_{adder} + T_{comparator}$ (= 25ns)
  $\Rightarrow f_{ref}$ = 40Mhz

- Total capacitance being switched = $C_{ref}$

- $V_{dd} = V_{ref}$ = 5V

- Power for reference datapath = $P_{ref} = C_{ref} V_{ref}^2 f_{ref}$

  **from [Chandrakasan92] (*IEEE JSSC*)**

26

# Parallel Datapath



**Area = 1476 x 1219 μ²**
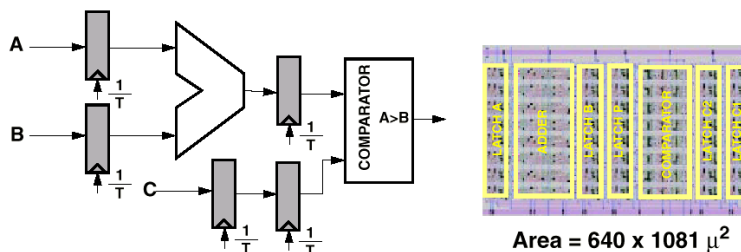
- The clock rate can be reduced by half with the same throughput $\Rightarrow f_{par} = f_{ref} / 2$
- $V_{par} = V_{ref} / 1.7$, $C_{par} = 2.15 C_{ref}$
- $P_{par} = (2.15 C_{ref}) (V_{ref}/1.7)^2 (f_{ref}/2) \approx 0.36 P_{ref}$

EECS251B L20 LOW-POWER DESIGN 27

27

# Pipelined Datapath



**Area = 640 x 1081 μ²**

- Critical path delay is less $\Rightarrow \max [T_{adder}, T_{comparator}]$

- Keeping clock rate constant: $f_{pipe} = f_{ref}$
  Voltage can be dropped $\Rightarrow V_{pipe} = V_{ref} / 1.7$

- Capacitance slightly higher: $C_{pipe} = 1.15 C_{ref}$

- $P_{pipe} = (1.15 C_{ref}) (V_{ref}/1.7)^2 f_{ref} \approx 0.39 P_{ref}$

EECS251B L20 LOW-POWER DESIGN 28

28

## A Simple Datapath: Summary

| Architecture type | Voltage | Area | Power |
|---|---|---|---|
| Simple datapath (no pipelining or parallelism) | 5V | 1 | 1 |
| Pipelined datapath | 2.9V | 1.3 | 0.39 |
| Parallel datapath | 2.9V | 3.4 | 0.36 |
| Pipeline-Parallel | 2.0V | 3.7 | 0.2 |

EECS251B L20 LOW-POWER DESIGN

29

29

# Multiple Supplies

EECS251B L20 LOW-POWER DESIGN

30

30

# Multiple Supply Voltages

- Block-level supply assignment ("power domains" or "voltage islands")
  - Higher throughput/lower latency functions are implemented in higher $V_{DD}$
  - Slower functions are implemented with lower $V_{DD}$
  - Often called "Voltage islands"
  - Separate supply grids, level conversion performed at block boundaries
- Multiple supplies inside a block
  - Non-critical paths moved to lower supply voltage
  - Level conversion within the block
  - Physical design challenging
  - (Not used in practice)

31

# Power Domains

32

## Practical Examples

- Intel Skylake (ISSCC'16)
  - Four power planes indicated by colors

33

## Practical Examples

- Intel 28-core Skylake-SP (ISSCC'18)



- 9 primary VCC domains are partitioned into 35 VCC planes
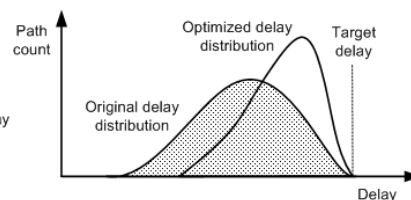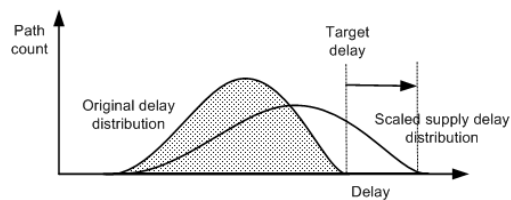
34

## Leakage Issue

- Driving from $V_{DDL}$ to $V_{DDH}$

› Level converter

35

## Multiple Supplies Within A Block

- Downsizing, lowering the supply on the critical path will lower the operating frequency
- Downsize (lowering supply) non-critical paths
  - Narrows down the path delay distribution
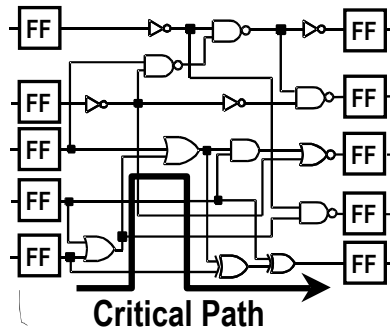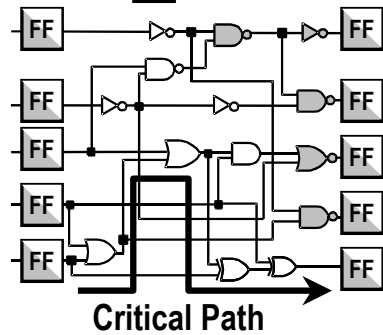  - Increases impact of variations

36

## Multiple Supplies in a Block

**Conventional Design**

**CVS Structure**

FF Level-Shifting F/F



Critical Path

Critical Path

Lower $V_{DD}$ portion is shaded

**M.Takahashi, ISSCC'98.**

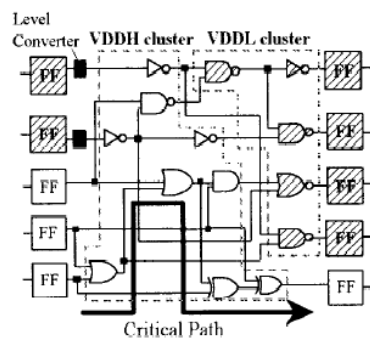"Clustered voltage scaling"

EECS251B L20 LOW-POWER DESIGN

37

37

## Multiple Supplies in a Block

CVS

Layout:



Usami'98

EECS251B L20 LOW-POWER DESIGN

38

38

19

## Level-Converting Flip-Flop

39

## Summary

- Power-performance tradeoffs
  - Sizing
  - Supplies
  - Thresholds
- Lowering supplies
- Multiple supply voltages
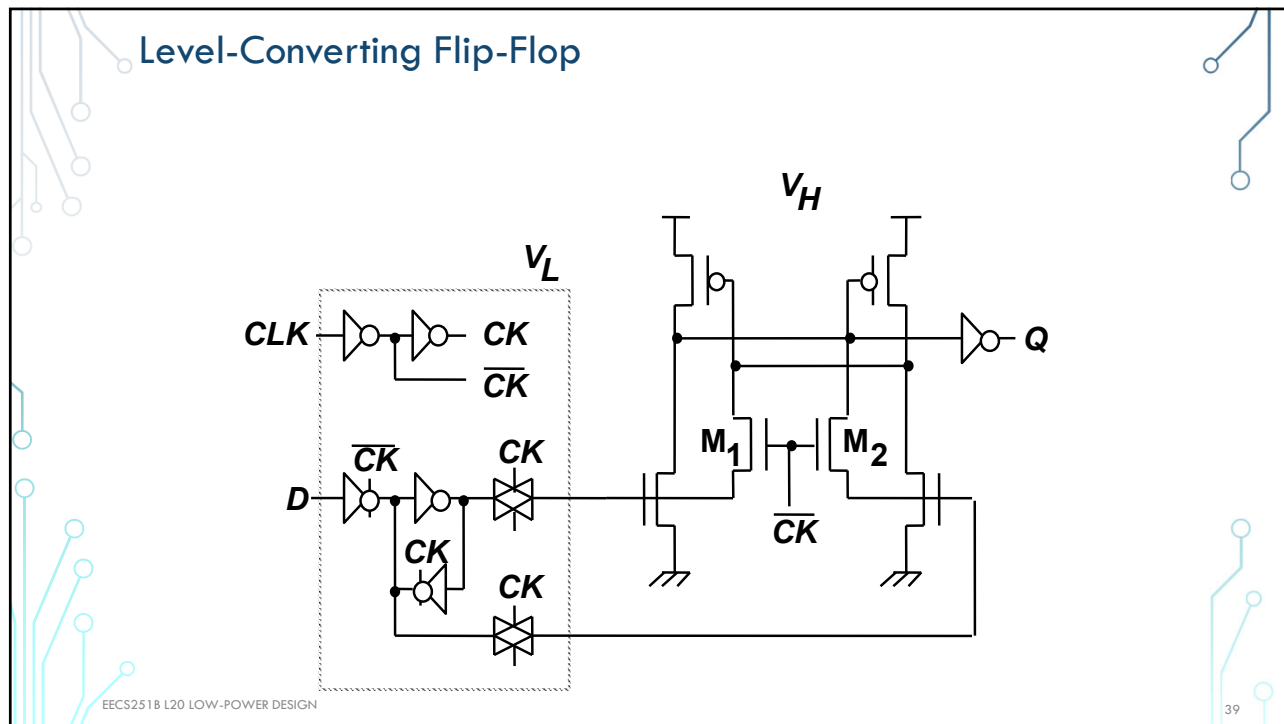
40

20

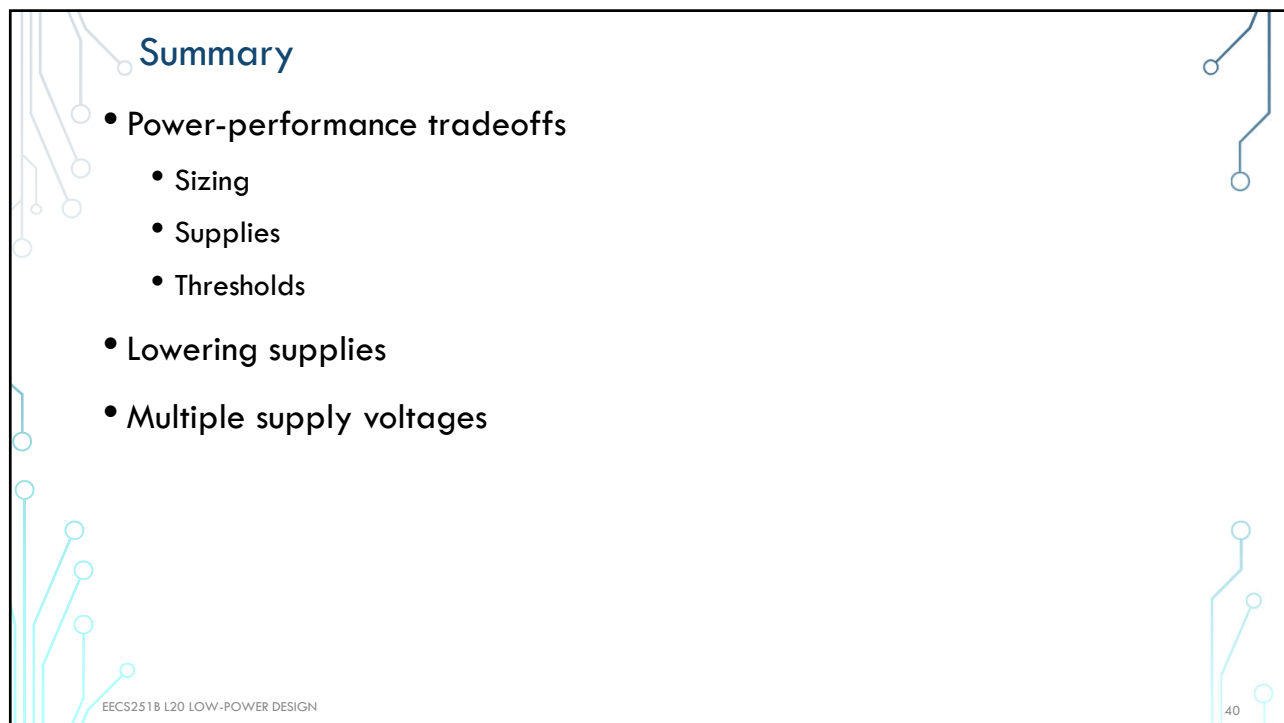## Next Lecture

- Low-power design
  - Dynamic voltage-frequency scaling
  - Clock gating

41

41