# 18. Design for Low Power

Jacob Abraham

Department of Electrical and Computer Engineering
The University of Texas at Austin

VLSI Design
Fall 2020

October 29, 2020

# Power and Energy

Power is drawn from a voltage source attached to the $V_{DD}$ pin(s) of a chip

**Instantaneous Power:**

$$P(t) = i_{DD}(t)V_{DD}$$

**Energy:**

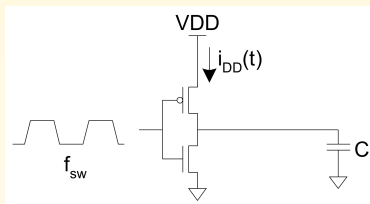$$E = \int_0^T P(t)dt = \int_0^T i_{DD}(t)V_{DD}dt$$

**Average Power:**

$$P_{avg} = \frac{E}{T} = \frac{1}{T}\int_0^T i_{DD}(t)V_{DD}dt$$

**Energy stored in capacitor when it is charged from $0$ to $V_C$,**

$$E_C = \int_0^\infty I(t)V(t)dt = \int_0^\infty C\frac{dV}{dt}V(t)dt = C\int_0^{V_c} V(t)dV = \frac{1}{2}CV_C^2$$

The capacitor releases this energy when it discharges back to 0

# Example – CMOS Inverter Driving a Load Capacitance



- When input switches from 1 to 0, pMOS transistor turns on and charges the load to $V_{DD}$
- Energy stored in the capacitor is $E_c = 1/2 C_L V_{DD}^2$
- Energy delivered from the power supply is

$$E_c = \int_0^\infty C \frac{dV}{dt} V_{DD} dt = CV_{DD} \int_0^{V_{DD}} dV = CV_{DD}^2$$

**Only half of the energy from the power supply is stored in the capacitor**
**The other half is converted to heat (resistance of the pMOS transistor)**

# Sources of Power Dissipation

## Dynamic Power Dissipation

- Charging and discharging of load capacitances
- "Short-circuit" current while both p- and n-MOS networks are partially on
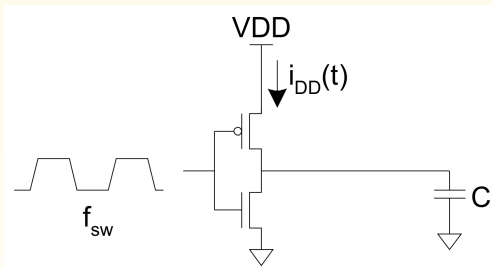
## Static Dissipation

- Subthreshold leakage (through OFF transistors)
- Gate leakage through gate dielectric
- Junction leakage from source/drain diffusion
- Contention current in ratioed circuits

## Dynamic Power

- Dynamic power is required to charge and discharge load capacitances when transistors switch
- One cycle involves a rising and falling output
- On rising output, charge $Q = CV_{DD}$ is required
- On falling output, charge is dumped to GND
- This repeats $T f_{sw}$ times over an interval of T

$$
\begin{aligned}
P_{dynamic} &= \frac{1}{T} \int_0^T i_{DD}(t) V_{DD} dt \\
&= \frac{V_{DD}}{T} \int_0^T i_{DD}(t) dt \\
&= \frac{V_{DD}}{T} \left[ T f_{sw} C V_{DD} \right] \\
&= \boxed{C V_{DD}^2 f_{sw}}
\end{aligned}
$$

# Activity Factor

- Suppose the system clock frequency = f
- Let $f_{sw} = \alpha f$, where $\alpha$ = activity factor
  - If the signal is a clock, $\alpha = 1$
  - If the signal switches once per cycle, $\alpha = 1/2$
  - Dynamic gates: switch either 0 or 2 times per cycle, $\alpha = 1/2$
  - Static gates: depends on design, but typically $\alpha = 0.1$

- Dynamic power: $\boxed{P_{dynamic} = \alpha C V_{DD}^2 f}$

$P_i$: probability that node $i$ is 1 ($1 - P_i$ is probability that it is 0)

Activity factor of node i, $\alpha_i$, is the probability that the node is 0 in one cycle and 1 in the next

If probability is uncorrelated from cycle to cycle, $\alpha_i = \bar{P}_i P_i$

Example: 4-input AND gate



**Tools exist to calculate activity factors, either using probabilities, or by monitoring nodes during simulation**

Where there is reconvergent fanout, calculating probabilities becomes more difficult

Example, glitches in chain of gates and inverters implementing 4-input NAND gate

# Short Circuit ("Crowbar") Current

- When transistors switch, both nMOS and pMOS networks may be momentarily ON at once
- Leads to a blip of "short circuit" current.
- < 10% of dynamic power if rise/fall times are comparable for input and output



Source: EE Times, June 9, 2003

Power reduction depends on the sizes of the driving and driven transistors and the input slew

# Example

- 200 million transistor chip
  - 20M logic transistors, average width: $12\lambda$
  - 180M memory transistors, average width $4\lambda$
  - 1.2 V 100 nm process
  - $C_g = 2$ fF/$\mu$m

## Estimate dynamic power

- Static CMOS logic gates: activity factor $= 0.1$
- Memory arrays: activity factor $= 0.05$ (many banks!)
- Estimate dynamic power consumption per MHz (neglect wire capacitance)

$$C_{logic} = (20 \times 10^6)(12\lambda)(0.05\mu m/\lambda)(2fF/\mu m) = 24nF$$

$$C_{mem} = (180 \times 10^6)(4\lambda)(0.05\mu m/\lambda)(2fF/\mu m) = 72nF$$

$$P_{dynamic} = [0.1C_{logic} + 0.05C_{mem}](1.2)^2 f = 8.6mW/MHz$$

## Static Power

- Static power is consumed even when chip is quiescent.
  - Ratioed circuits burn power in fight between ON transistors
  - Leakage draws power from nominally OFF devices

$$I_{ds} = I_{ds0} e^{\frac{V_{gs} - V_t}{n v_T}} \left[ 1 - e^{\frac{-V_{ds}}{v_T}} \right]$$

$$V_t = V_{t0} - \eta V_{ds} + \gamma \left( \sqrt{\phi_s + V_{sb}} - \sqrt{\phi_s} \right)$$

$\eta$ describes drain-induced barrier lowering (DIBL),

$\gamma$ describes the body effect

For any appreciable $V_{ds}$, the term in brackets approaches unity

## Leakage Example: Estimate Static Power

- Process has two threshold voltages and two oxide thicknesses
- Subthreshold leakage:
  - 20 nA/$\mu$m for low $V_t$
  - 0.02 nA/$\mu$m for high $V_t$
- Gate leakage:
  - 3 nA/$\mu$m for thin oxide
  - 0.002 nA/$\mu$m for thick oxide
- Memories use low-leakage transistors everywhere, and gates use low-leakage transistors on 80% of logic

High leakage: $(20 \times 10^6)(0.2)(12\lambda)(0.05\mu m/\lambda) = 2.4 \times 10^6 \mu m$

Low leakage:

$(20 \times 10^6)(0.8)(12\lambda)(0.05\mu m/\lambda) + (180 \times 10^6)(4\lambda)(0.05\mu m/\lambda) = 45.6 \times 10^6 \mu m$

$I_{static} = (2.4 \times 10^6 \mu m)[(20nA/\mu m)/2 + (3nA/\mu m)] + (45.6 \times 10^6 \mu m)[(0.02nA/\mu m)/2 + (0.002nA/\mu m)] = 32mA$

$P_{static} = I_{static}V_{DD} = 38 \ mW$

If no low-leakage devices used, $P_{static} = 749 \ mW$

Source: Shekhar Borkar, Intel

Source: Shekhar Borkar, Intel

Source: Shekhar Borkar, Intel

# Leakage Becoming A Major Component of Power

- Leakage component to active power becomes significant % of total power
- $\approx$10% in 0.18$\mu$m technology
- Acceptable limit less than $\approx$10%, implies serious challenge in $V_t$ scaling!



Sources: S. Borkar, Intel; Chip Design Magazine

# Low Power Design

- Reduce dynamic power
  - $\alpha$: clock gating, sleep mode
  - $C$: small transistors (especially on clock), short wires
  - $V_{DD}$: lowest suitable voltage
  - $f$: lowest suitable frequency
- Reduce static power
  - Selectively use ratioed circuits
  - Selectively use low $V_t$ devices
  - Leakage reduction: stacked devices, body bias, low temperature

### Use a combination of techniques at different levels

- Algorithm
- Architecture
- Logic/circuit
- Technology/circuit

Data-path operator

Parallel Implementation

$$P_{par} = (2.15C)(0.58V)^2(0.5f) \approx 0.36P$$

Pipelined Implementation

$$P_{pipe} = (1.15C)(0.58V)^2(f) \approx 0.39P$$

Reducing operations, while maintaining throughput

Reducing operations, with lower throughput

Precomputation Architecture

$$f_1 = 1 \implies Z = 1; \quad f_2 = 1 \implies Z = 0$$

# Precomputation-Based Optimization for Low Power, Cont'd



N-bit Comparator

$$f_1 = A(n-1) \cdot \overline{B(n-1)}; \quad f_2 = \overline{A(n-1)} \cdot B(n-1)$$

Adder-comparator circuit

$$f_1 = A(n-1) \cdot B(n-1) \cdot \overline{C(n-1)} \cdot \overline{D(n-1)}$$
$$f_2 = \overline{A(n-1)} \cdot \overline{B(n-1)} \cdot C(n-1) \cdot D(n-1)$$

# Stack Effect – Subthreshold Leakage



Stack effect reduces subthreshold leakage by a factor of $\approx 10$

Stacks with three or more OFF transistors have even lower leakage

Silicon-on-Insulator (SOI) circuits are attractive for low-leakage designs

Affected by voltage across the gate

# Gate and Subthreshold Leakage in NAND3 (nA)

| Input State (ABC) | $I_{sub}$ | $I_{gate}$ | $I_{total}$ | $V_x$ | $V_z$ |
|:---:|:---:|:---:|:---:|:---:|:---:|
| 000 | 0.4 | 0 | 0.4 | stack effect | stack effect |
| 001 | 0.7 | 0 | 0.7 | stack effect | $V_{DD} - V_t$ |
| 010 | 0 | 1.3 | 1.3 | intermediate | intermediate |
| 011 | 3.8 | 0 | 10.1 | $V_{DD} - V_t$ | $V_{DD} - V_t$ |
| 100 | 0.7 | 6.3 | 7.0 | 0 | stack effect |
| 101 | 3.8 | 6.3 | 10.1 | 0 | $V_{DD} - V_t$ |
| 110 | 5.6 | 12.6 | 18.2 | 0 | 0 |
| 111 | 28 | 18.9 | 46.9 | 0 | 0 |

# Controlling Threshold Voltages for Reduced Leakage

## Multiple $V_t$, Longer channels, Oxide thicknesses

- Low-$V_t$ on critical paths, High-$V_t$ on other paths for reduced leakage
- Longer transistors in the caches
- Thicker oxides for I/O transistors

## Body Bias

# Voltage Domains for Low Power



Level Converter



Clustered Voltage Scaling

# Dynamic Voltage Scaling (DVS)

# RAZOR

- Error-tolerant dynamic voltage scaling (DVS) technology which eliminates the need for the voltage margins required for "always correct" circuit operations design
- A different value in the shadow latch shows timing errors
- Pipeline state is recovered after timing-error detection
- Error detection is done at the circuit level
  - The design overhead is large if timing paths are well balanced in the design



Austin et al., 2003

# Direct Monitoring of Critical Path

## Razor Flip-Flop (a) and Architecture using it (b)



- Speculative operation requires an additional pipeline stage
- Design may not be suitable for designs that have many critical paths (increase in area and flip-flop power)

# Indirect Critical Path Monitor

## TEAtime approach



- Use of Critical Path Replicas (CPRs) to control voltage or frequency until one of them fail
- CPRs (1-bit version of potential critical paths) are located near potential critical paths to monitor them
- 1-bit detector may result in "oscillations"

Critical Path Monitor (CPM)

| P[2:0] | Delay of CPRs | Frequency Control |
|--------|---------------|-------------------|
| {0,0,0} | Fast | ↑ |
| {0,0,1} | Appropriate | — |
| {0,1,1} | Slow (Safety Margin) | ↓ |
| {1,1,1} | | |

C-element and logic function



Configuration of 8 CPRs

## Simulation Results

- MIPS core implemented in 45nm process
- Optimized to meet target frequency of 1.5GHz
  - Many critical paths
- Power results from HSPICE, PrimeTime and PrimeTimePX



Delay changes in critical paths and CPMs



Maximum improvement in Energy-Delay product

- Given a microprocessor design and an instruction
  - Identify the instruction-driven slice
  - Shut off the rest of the circuitry
- This might include
  - Gating out parts of different blocks
  - Gating out floating point units during integer ALU execution
  - Turning off certain FSMs in different control blocks since exact constraints on their inputs are available due to instruction-driven slicing



OR1200-RTL Reduction in Power Dissipation

# Low Power by Design: StrongArm 110

Start with Alpha 21064: 200 MHz @ 3.45V, Power = 26 W

Vdd reduction:       Power reduction = 5.3X  $\implies$  4.9W

Reduce functions:   Power reduction = 3X  $\implies$  1.6W

Scale process:       Power reduction = 2X  $\implies$  0.8W

Clock load:          Power reduction = 1.3X  $\implies$  0.6W

Clock rate:          Power reduction = 1.25X  $\implies$  0.5W

Source: D. Dobberpuhl

## LongRun Technology Demonstration

| MHz | Voltage | % Full Power |
|-----|---------|--------------|
| 700 | 1.65 | 100% |
| 400 | 1.4 | 41% |
| 333 | 1.2 | 25% |

**Power = C x V$^2$ x F = 400MHz/700MHz * 1.4V$^2$/1.65V$^2$ = 41%**

- **Crusoe processor starts off at 700MHz**
- **DVD movie requires between 333 and 400MHz**
- **Power is reduced to 25 or 41% of full power**
- **The result is extended DVD playtime**

Source: Doug Laird

**LongRun Technology in Operation**

- Crusoe processor starts off at 700MHz
- Code Morphing software detects user activity
- The software dynamically adjusts MHz and voltage to the most efficient power level

Crusoe Processor AC/DC Modes

| MHz | Voltage |
|-----|---------|
| 700 | 1.65 |
| 667 | 1.65 |
| 633 | 1.60 |
| 600 | 1.60 |
| 566 | 1.55 |
| 533 | 1.55 |
| 500 | 1.50 |
| 466 | 1.50 |
| 433 | 1.45 |
| 400 | 1.40 |
| 366 | 1.35 |
| 333 | 1.30 |
| 300 | 1.25 |
| 266 | 1.20 |
| 233 | 1.15 |
| 200 | 1.10 |

Source: Doug Laird

**Dynamic Software Execution (2nd Pass)**

x86 Memory

Translation Cache

Execute →

| | VLIW code for block #1 | |
|---|---|---|
| | VLIW code for block #2 | |
| | VLIW code for block #3 | |
| | | |

Code Morphing software actions
- Translation already exists
- Execute VLIW code immediately

Source: Doug Laird

**Processor Thermal Comparison**

Pentium III
Playing DVD

Crusoe Processor
Playing DVD

105.5º C
221.9º F

48.2º C
118.8º F

Source: Doug Laird