

---

## Lecture 11

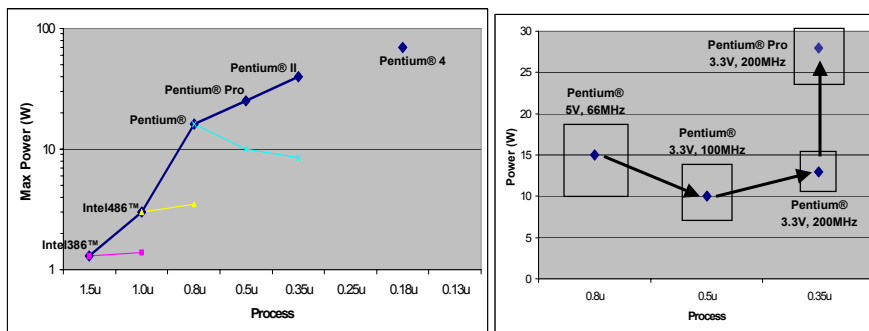
### Low Power and Power Efficient Circuits

Intel Corporation  
jstinson@stanford.edu

(many thanks to Intel University for much of the material)

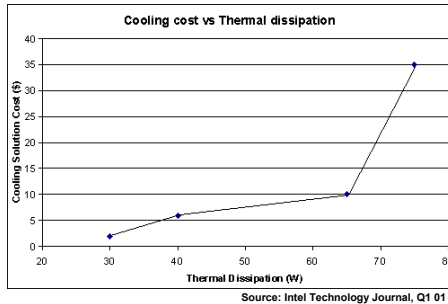
---

## Power Trends



- Power historically has grown at almost 2X/generation
  - Frequency growing at ~2X/generation
  - Xtor count growing at ~2X/generation
- Voltage scaling and process temporary relief within a uproc generation
  - Has not prevented growth thru architectural generations

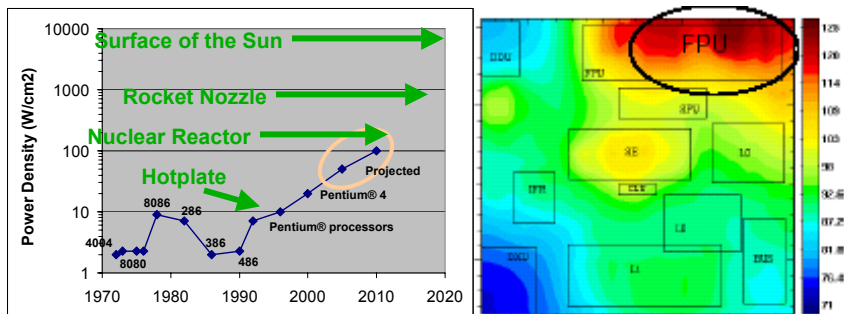
# Power Cost



- Cost of delivery power and thermal solutions is growing at an exponential rate
  - 25000 sq. ft. data center, ~8000 servers = 2,000,000 Watts!!
  - Cost driven by rack height, cooling air flow, power delivery, maintenance

“What matters most to the computer designers at Google is not speed, but power—low power, because data centers can consume as much electricity as a city.”  
 -Eric Schmidt, CEO Google (NYT, 09/29/02)

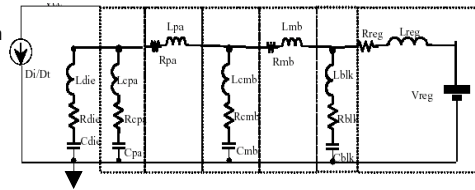
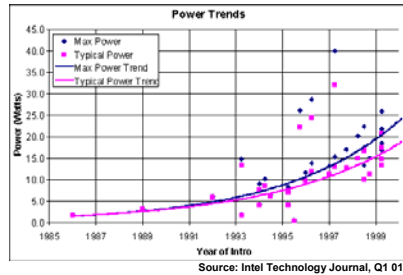
# Power/Thermal Density



- Power is going up...yet die size is going down
  - Power isn't evenly distributed across the die
    - Makes local hotspots even more difficult
- Designs must be reliable at worst case temperature and voltage
  - Speed, noise, leakage, oxide wearout, EM, SH, etc.

## di/dt

- Changes in current (A) create large di/dt events
  - Min power state to max power state can occur in less than 4-5 cycles
    - Min power state is typically 30-40% lower than max power state
- Architectural efforts to exploit parallelism increase di/dt
  - Putting more hardware on die to do work in parallel
    - Turn it off to save power ☺
    - Creates large discrepancy between hi and lo utilization power ☺
- Low power efforts can increase di/dt
  - Turning off hardware creates current (A) switching events



## Dynamic Power Analysis

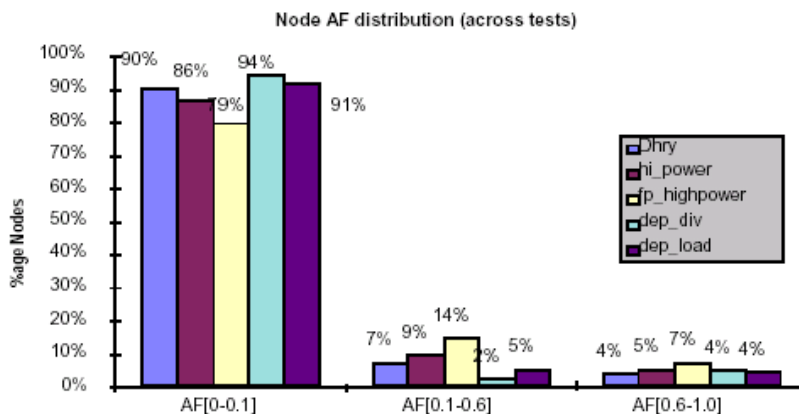
$$P = \underbrace{AF * C}_{\text{Greatest flexibility/control}} * \underbrace{V^2 * F}_{\text{Usually fixed by the project}}$$

- Activity Factor (AF) – switching or toggle rate
  - Measurement of how often a node switches
  - Logic and circuit design can control
- Capacitance
- Voltage and Frequency
  - Typically set by the product spec
- Majority of power reduction techniques focus on AF and C

## Where does the dynamic power go?

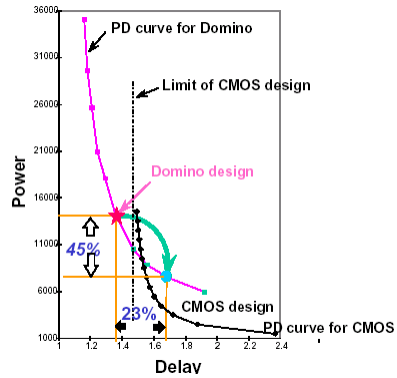
- Majority of power consumed in the clock/clocked elements
  - Clock distribution, sequentials, domino, enables, clocked logic
  - 5-10% of the node capacitance—close to 50% of the power!
    - AF makes the difference
- Large I/O and bus drivers
  - Large capacitances
  - Not too many...tends to be localized issues more than total power
- Datapath
  - Tend to upsize drivers to meet frequency requirements
  - High utilization
- Memory
  - Tend to be low power (dynamic)
  - Can be largest capacitance on die, but lowest power

## Node Activity Factors

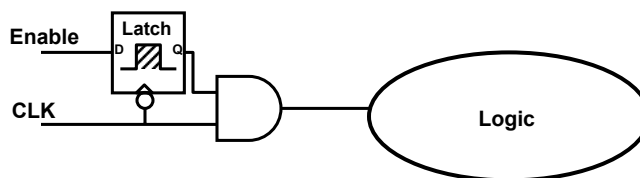


## Device Sizing and Design Style

- Most designs have optimum point for power/performance
  - Historically, designers always went for the speed....
  - Low power designs must choose the most optimum point
- Domino vs. CMOS
  - Domino node cap lower than CMOS (typically 65-75%)
  - Clock loading and AF usually make domino significantly higher power
    - Design domino logic so that precharge state is “most common”



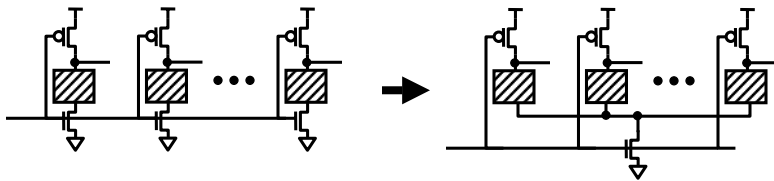
## Clock Gating



- Look for opportunities in the design to “turn off” the clock
  - Unused hardware, guaranteed to evaluate in a certain direction, don’t care
- Advantages
  - Reduces activity factor for clocks
  - Prevents downstream data from switching
- Disadvantages
  - Can be difficult to identify opportunities
    - Goal of most architectures is to keep everything “busy”
  - Often creates speedpaths
    - Easier to identify gating opportunity “closest to the event”
  - di/dt problems

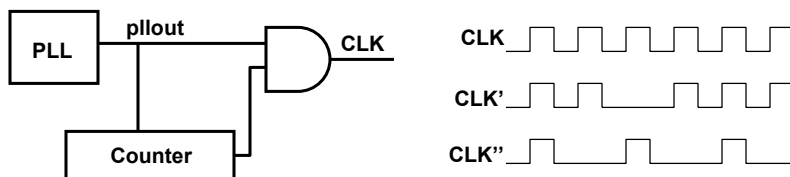
## Clock Collapsing

- Combine local clocked xtors into grouped single xtor
  - Reduces total driven xtors width
  - Reduces “generation” overhead (pulses, enables)
    - Locally generated pulses/enables require distributed logic—power hungry
  - Can cause a performance impact
    - Group domino logic that is unlikely to evaluate together
  - Can increase clock skew (esp. on pulses)
    - Tradeoff power reduction vs. race margin



## Clock Throttling

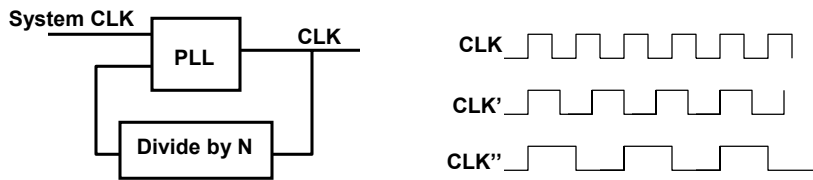
- Dynamically change the frequency of the design
  - Directly impacts the dynamic power of the design
  - Need effective means of changing the frequency on-the-fly
  - Need effective means of triggering the frequency change
- “Missing Teeth”
  - Globally disable every  $N^{\text{th}}$  high phase of the clock
  - Need architecture that understands on-the-fly frequency variation
  - Design must still meet constraints of full frequency clock



## Clock Throttling (cont'd)

---

- Re-sync the PLL
  - Change the fundamental frequency coming from PLL
  - Design doesn't have to meet full frequency constraints during throttle
    - Opportunity to dynamically lower voltage
  - Must wait for PLL to lock to new frequency
    - Limits fast transition in and out of throttled mode



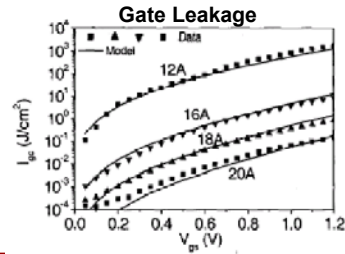
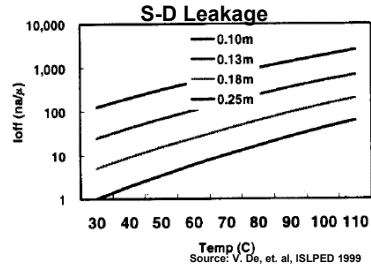
## Low Swing Circuits

---

- Separate power supply for portion of the design
  - Using either an explicit power supply or on-die regulation
  - Needs separate supply lines
  - Need to separate out well taps
  - Usually share Vss, split Vcc
    - Cuts the overhead in half
- Lower supply = lower performance
  - Identify portions of the design that are not speed critical
    - Caches are one of the best targets
    - Localized logic within larger blocks (hard)
  - Identify structures that are less impacted by lower supply voltage
    - Differential, low swing circuit families (long wires)

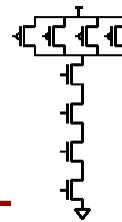
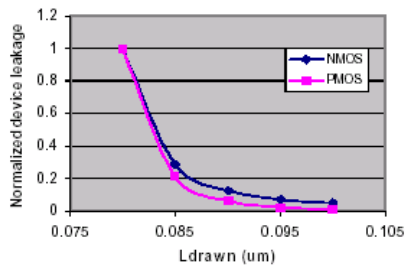
# Leakage

- Leakage increasing at a ~3-4X per process generation
  - loff inversely proportional to  $\exp(V_{th})$ 
    - $V_{th}$  is reducing for performance
  - Gate leakage inversely proportional to  $\exp(T_{ox})$ 
    - Gate leakage quickly approaching same levels as S-D leakage
- Modern microprocessors have up to 20-30% leakage budget
  - Within the large caches, leakage can account for up to 50% total power
  - Burnin conditions (stress testing) can increase leakage to over 90% of total power
    - Strong temperature, voltage dependence



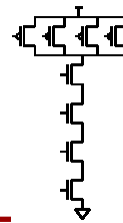
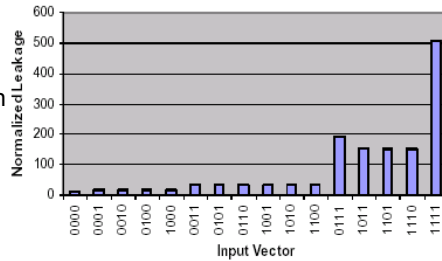
# Long Le

- Increasing channel length dramatically reduces leakage
  - Insignificant increase in gate leakage
  - Insignificant increase in gate capacitance (performance)
- Corresponding performance loss
  - $I_{dsat}$  reduced with longer channel
  - Must identify circuits with no performance impact
    - Memories good candidates



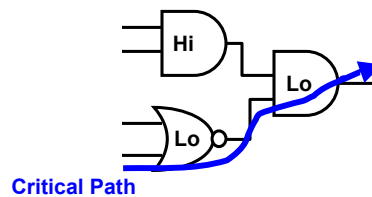
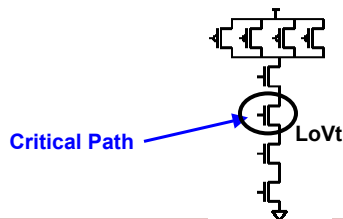
## Transistor Stacking

- Stacked devices reduce leakage
  - Intermediate voltages between stacked devices help combat leakage
  - Only true for "off" xtors in series!!
- Rely on statistical nature of inputs to reduce overall leakage
  - May not be true depending on real distribution of input vectors



## Multiple Vt

- Strong trend towards multi-Vt processes
  - Provides two types of FETs to designer:
    - HiVt: geared towards non-critical transistors to reduce leakage
    - LoVt: geared towards most critical transistors
  - Usage/success highly dependent on ability to detect critical paths
    - Large complex designs not always easy to identify the right locations
    - \*Can\* run into situation where wrong application of LoVt devices actually slows down the device!!



## Body Biasing

---

- Reverse body-bias xtors to increase  $V_t$ 
  - Leakage reduced significantly (S-D;  $I_{gate}$  reduced slightly)
  - Performance goes down accordingly
    - Can also forward bias xtors to increase performance (and leakage)
- Need to route separate supply for body-biasing
  - Can enable during normal operation
    - Either to decrease leakage OR increase performance
  - Can also enable only during special operating modes
    - Burnin – performance doesn't matter, leakage is everything