



University of Pittsburgh

CS 1541 Introduction

Technology Advances

Wonsun Ahn

Department of Computer Science

School of Computing and Information



Technology Advances



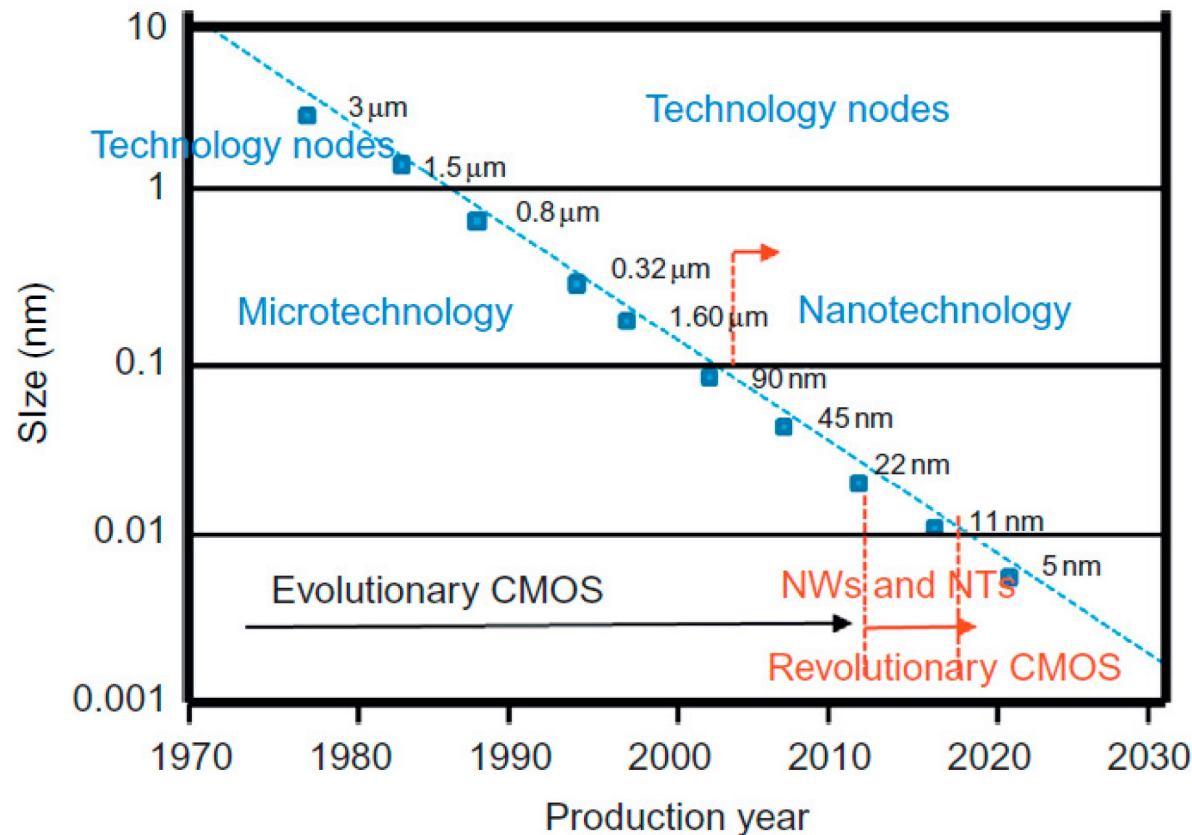


Advances in Technology

- Technology has been advancing at lightning speed
- Architecture and IT as a whole were beneficiaries
- Technology advance is summarized by *Moore's Law*
 - You probably heard of it at some point. Something about ...
 - "X doubles every 18-24 months at constant cost"
- Is X:
 - CPU performance?
 - CPU clock frequency?
 - Transistors per CPU chip?



Miniaturization of Transistors



Data source: Radamson, H.H.; He, X.; Zhang, Q.; Liu, J.; Cui, H.; Xiang, J.; Kong, Z.; Xiong, W.; Li, J.; Gao, J.; Yang, H.; Gu, S.; Zhao, X.; Du, Y.; Yu, J.; Wang, G. Miniaturization of CMOS. *Micromachines* **2019**, *10*, 293.

- Moore's Law has been driven by transistor miniaturization
 - CPU chip area hasn't changed much



Future of Moore's Law

- The semiconductor industry has produced roadmaps
 - Semiconductor Industry Association (SIA): 1977~1997
 - International Technology Roadmap for Semiconductors (ITRS): 1998~2016
 - International Roadmap for Devices and Systems (IRDS): 2017~Present

■ IRDS Lithography Projection (2020)

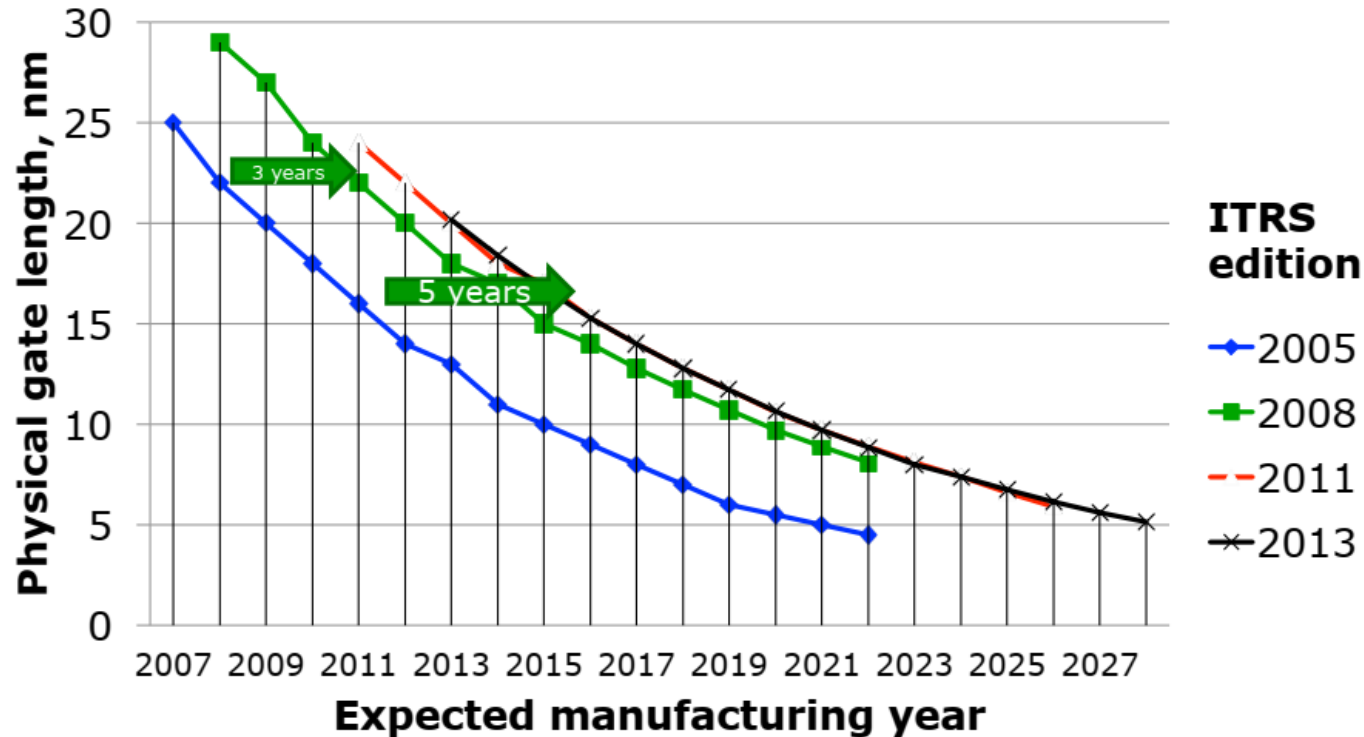
Year of Production	2018	2020	2022	2025	2028	2031	2034
Technology Node (nm)	7	5	3	2.1	1.5	1.0	0.7

- Moore's Law will continue into foreseeable future
- IRDS does not project significant increase in CPU chip size
- Increases in transistors will come from *transistor density*



IRDS isn't Perfect

- ITRS (predecessor of IRDS) has made corrections before



- After all, you are trying to predict the future
- But architects rely on the roadmap to design future processors



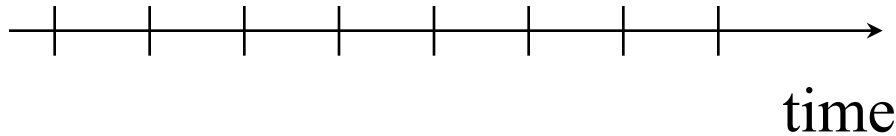
Moore's Law and Performance

- Did Moore's Law result in higher performance CPUs?
- When you decide on a CPU for your desktop, what number(s) do you look at to see how fast it is?
- Best way: try running your favorite apps on it.
- Everyone has different favorite apps so instead publish
 - Results from running a suite of benchmark apps (e.g. SPEC CPU)
 - Several important components of performance



Components of Execution Time

- Processor activity happens on clock “ticks” or cycles



- On each tick, bits flow through logic gates and are latched

- Execution time = $\frac{\text{seconds}}{\text{program}}$

$$\begin{aligned}\frac{\text{seconds}}{\text{program}} &= \frac{\text{cycles}}{\text{program}} \times \frac{\text{seconds}}{\text{cycle}} \\ &= \frac{\text{instructions}}{\text{program}} \times \frac{\text{cycles}}{\text{instruction}} \times \frac{\text{seconds}}{\text{cycle}}\end{aligned}$$



How to improve Execution Time

$$\frac{\text{instructions}}{\text{program}} \times \frac{\text{cycles}}{\text{instruction}} \times \frac{\text{seconds}}{\text{cycle}}$$

■ Reduce $\frac{\text{seconds}}{\text{cycle}}$:

- Clock frequency = $\frac{\text{cycles}}{\text{second}}$ = reverse of $\frac{\text{seconds}}{\text{cycle}}$
- Higher clock frequency (GHz) leads to shorter exec time

■ Reduce $\frac{\text{cycles}}{\text{instruction}}$:

- Also known as CPI (Cycles Per Instruction)
- IPC (Instructions Per Cycle) = $\frac{\text{instructions}}{\text{cycles}}$ = reverse of $\frac{\text{cycles}}{\text{instructions}}$
- Higher IPC leads to shorter execution time

■ Reduce $\frac{\text{instructions}}{\text{program}}$:

- Less instructions leads to shorter execution time
- ISAs that do a lot of work with one instruction shortens time



Moore's Law impacts two layers

- Did Moore's Law result in higher performance CPUs?
- Law impacts both architecture and physical layers

Instruction Set Architecture

Processor Organization

Transistor Implementation

Computer
Architecture

Physical Layer

- Processor Organization: many more transistors to use in design
- Transistor Implementation: smaller, more efficient transistors



Moore's Law Impact on Architecture

- So where did architects use all those transistors?
- Well, we will learn this throughout the semester 😊
 - Pipelining
 - Parallel execution
 - Prediction of values
 - Speculative execution
 - Memory caching
 - In short, they were used to improve frequency or IPC
- Let's go on to impact on the physical layer for now



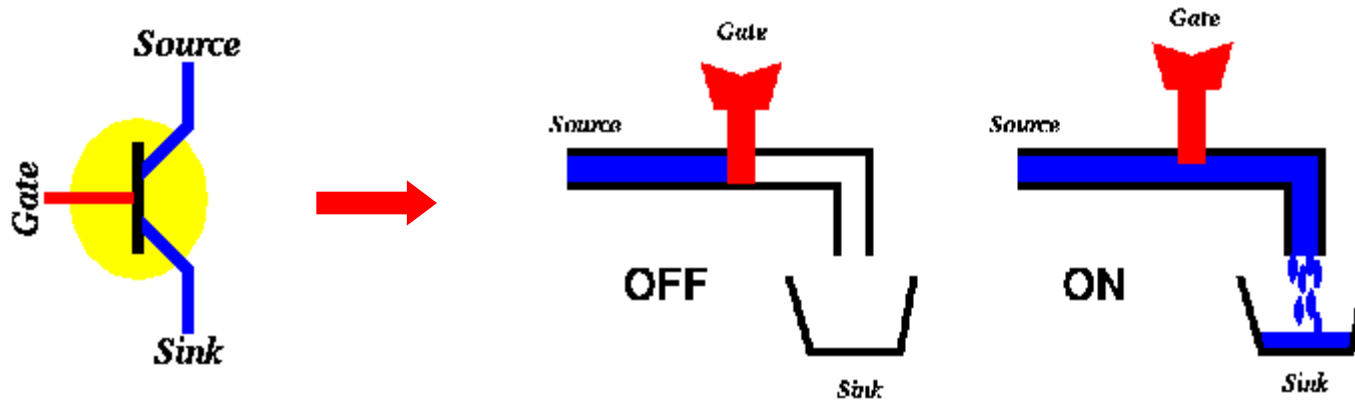
Moore's Law Impact on Physical Layer

- CPU frequency is also impacted by transistor speed
 - As well as how many transistors are in between clock ticks (which is determined by processor organization)
- So did Moore's Law result in faster transistors?
 - In other words, are smaller transistors faster?



Speed of Transistors

■ Transistor 101: Transistors are like faucets!



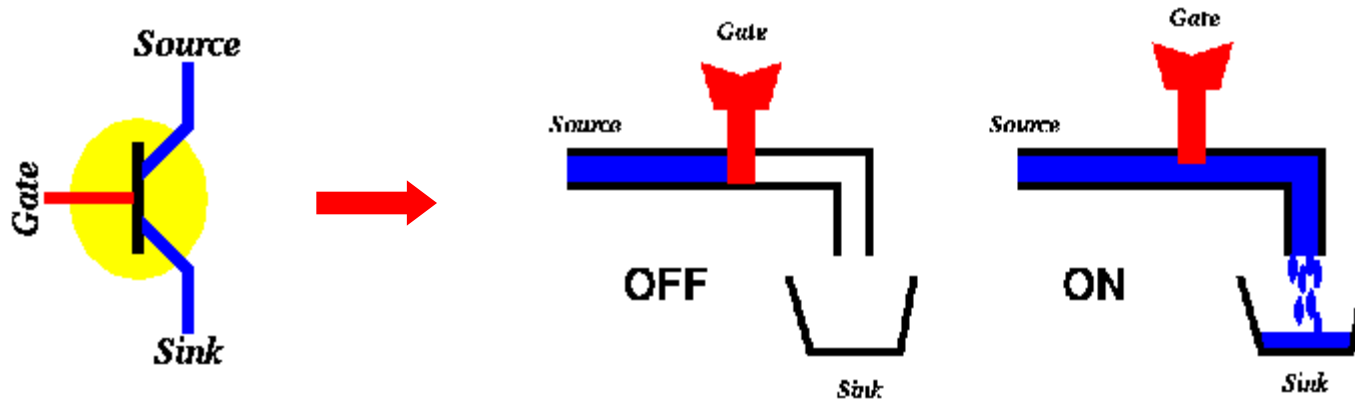
■ To make a transistor go fast, do one of the following:

- Reduce distance from source to sink (*channel length*) ↓
- Reduce bucket size (*capacitance*) ↓
- Increase water pressure (*supply voltage*) ↑



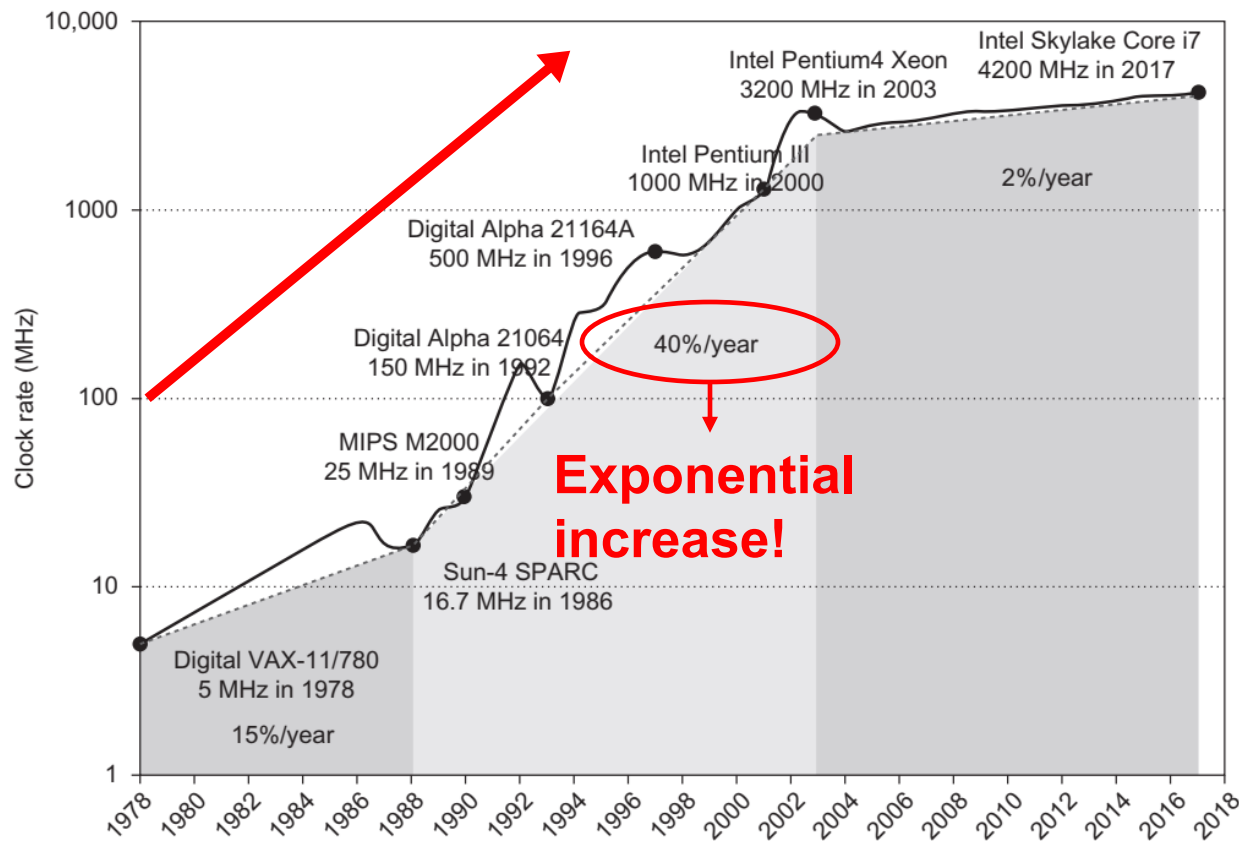
Smaller Transistors are Faster!

- Transistor 101: Transistors are like faucets!



- When a transistor gets smaller:
 - *Channel length L is reduced* ↓
 - *Capacitance C is reduced* ↓
- So, given the same *supply voltage*, smaller is faster!
- So, did Moore's Law enjoy faster and faster frequencies?

Yes, for a while ...



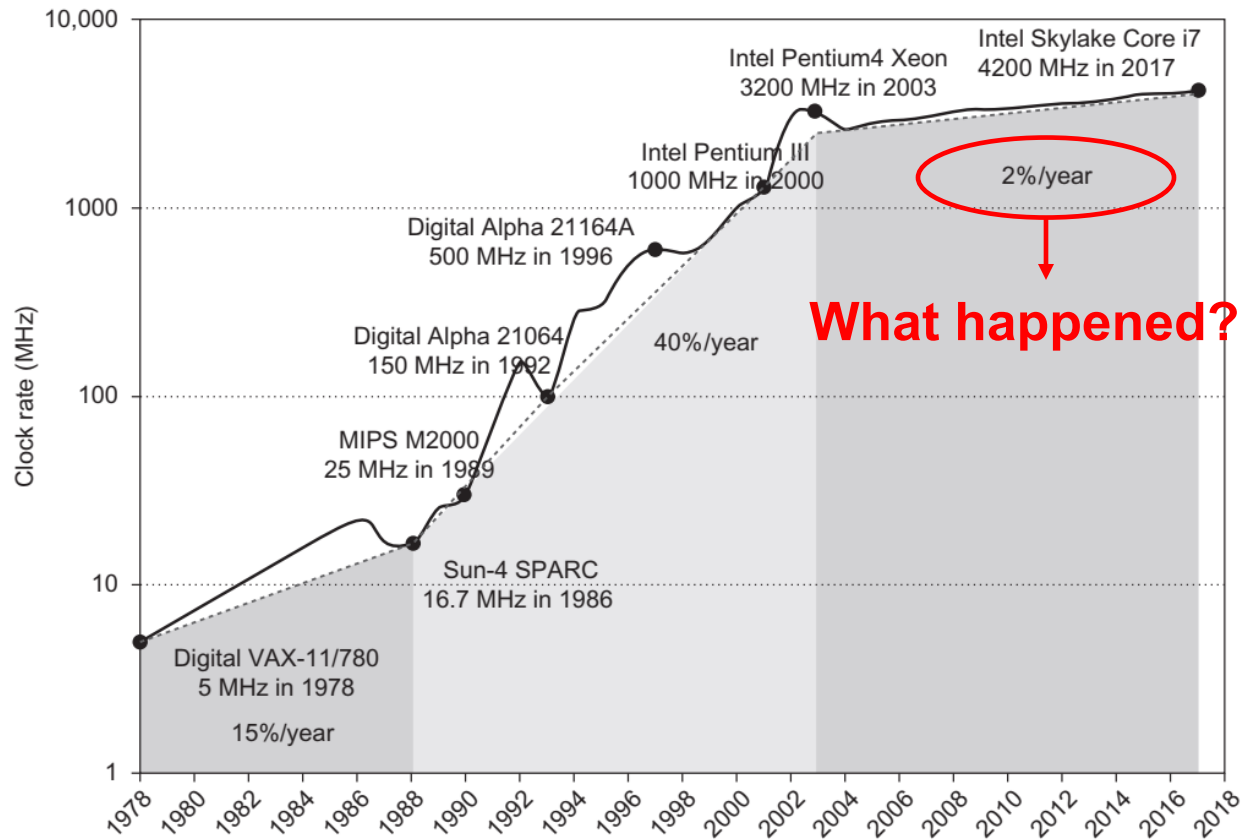
Source: Computer Architecture, A Quantitative Approach (6th ed.) by John Hennessy and David Patterson, 2017

■ Improvements in large part due to transistors

- Processor design also contributed but we'll discuss later



But not so much lately

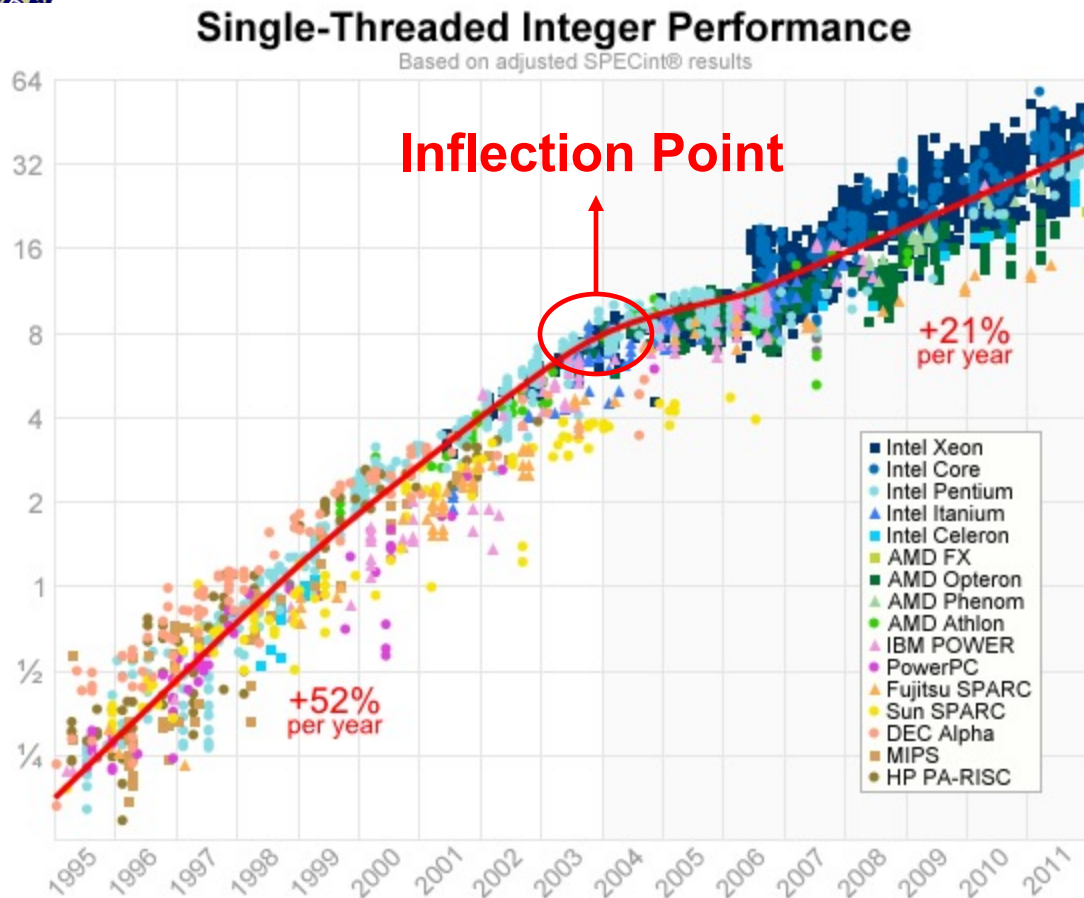


Source: Computer Architecture, A Quantitative Approach (6th ed.) by John Hennessy and David Patterson, 2017

- Suddenly around 2003, frequency scaling stops



Dent in CPU Performance



Source: <https://preshing.com/20120208/a-look-back-at-single-threaded-cpu-performance/>

- This caused a big dent in CPU performance at 2003
- Improvements henceforth only came from architecture
 - From improvements to IPC (instructions per cycle)



So What Happened? TDP.

■ *TDP (Thermal Design Power):*

- Maximum power (heat) that CPU is designed to generate
- Capped by the amount of heat cooling system can handle
- Cooling system hasn't improved much over generations

■ CPU Power = $A * N * CFV^2$ **must be < TDP**

- A = Activity factor (% of transistors with activity)
- N = Number of transistors
- C = Capacitance
- F = Frequency
- V = Supply Voltage



■ What happens to each factor with Moore's Law?



TDP and Moore's Law

- Change in CPU Power $\propto A * N * CFV^2$,
with reduction in d (one dimension of transistor):
 - A = Activity factor
 - V = Supply Voltage (water pressure)
 - N = Chip area / (x-dimension * y-dimension) $\propto 1/d^2$ ↑ ↑
 - C = Size of bucket $\propto d$ ↓
 - L = Channel length $\propto d$ ↓
 - $F \propto V/(C * L) \propto 1/d^2$ ↑ ↑
 - Decrease in C cannot offset increases in both N and F
 - That means F (frequency) needs to be decreased to meet TDP
- Q) So how did CPU frequency keep increasing up to 2003?



Dennard Scaling Maintains TDP

- By reducing **Supply Voltage** proportional to d , change in CPU Power $\propto A * N * CFV^2$ is:
 - A = Activity factor
 - V = **Supply Voltage** ($\propto d$) $\downarrow \rightarrow V^2 \downarrow \downarrow$
 - N = **Number of transistors** $\propto 1/d^2 \uparrow \uparrow$
 - C = **Capacitance** $\propto d \downarrow$
 - L = **Channel length** $\propto d \downarrow$
 - $F \propto V/(C * L) \propto 1/d \uparrow$
- Factors balance each other out to keep **power constant**
 - $V^2 \downarrow \downarrow$ cancels out $N \uparrow \uparrow$, $C \downarrow$ cancels out $F \uparrow$
- **Dennard Scaling:** Above **recipe** for scaling up frequency, while reducing supply voltage to keep power constant

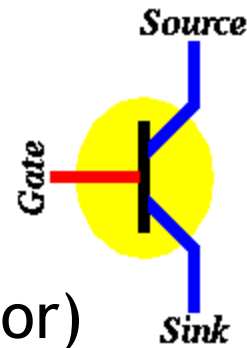


Dennard Scaling and V_{th}

- So, it's that easy? Just reduce V until you meet TDP?
- No, it's not that simple ☹.
- Reducing V_{dd} (supply voltage) affects CPU operation
 - As V_{dd} is reduced, CPU becomes slower and slower
 - Eventually, CPU stops working altogether
- CPU (specifically transistors) needs redesigning
 - V_{th} (threshold voltage) needs to be reduced along with V_{dd}
 - To understand this, we need a 101 on MOSFETs

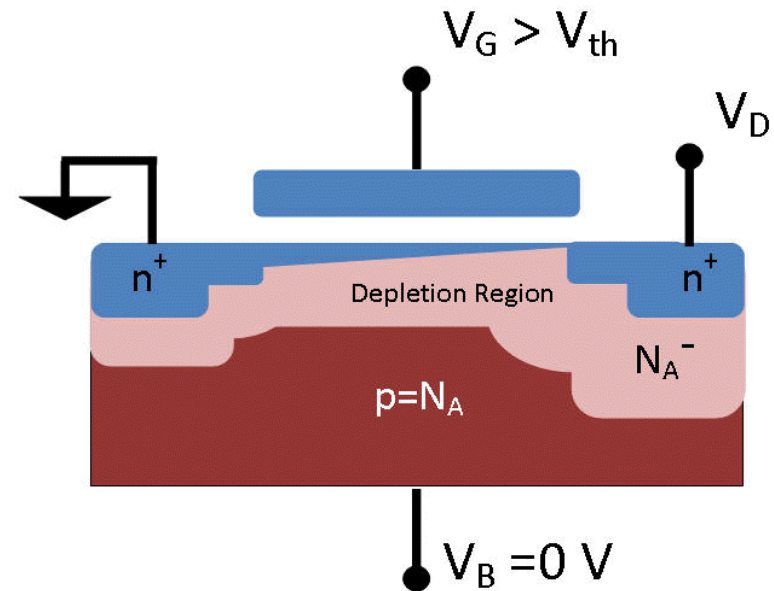
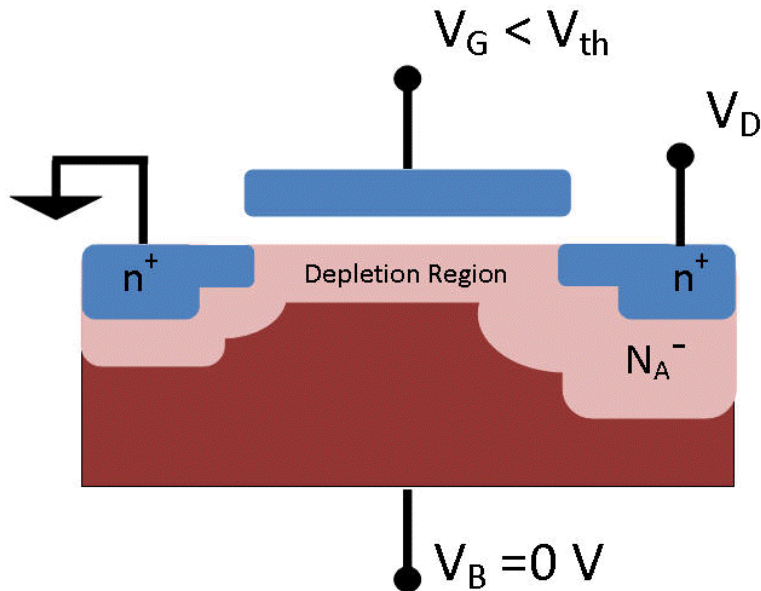


MOSFET 101



■ MOSFET (Metal Oxide Silicon Field Effect Transistor)

[A MOSFET transistor switched off] [A MOSFET transistor switched on]

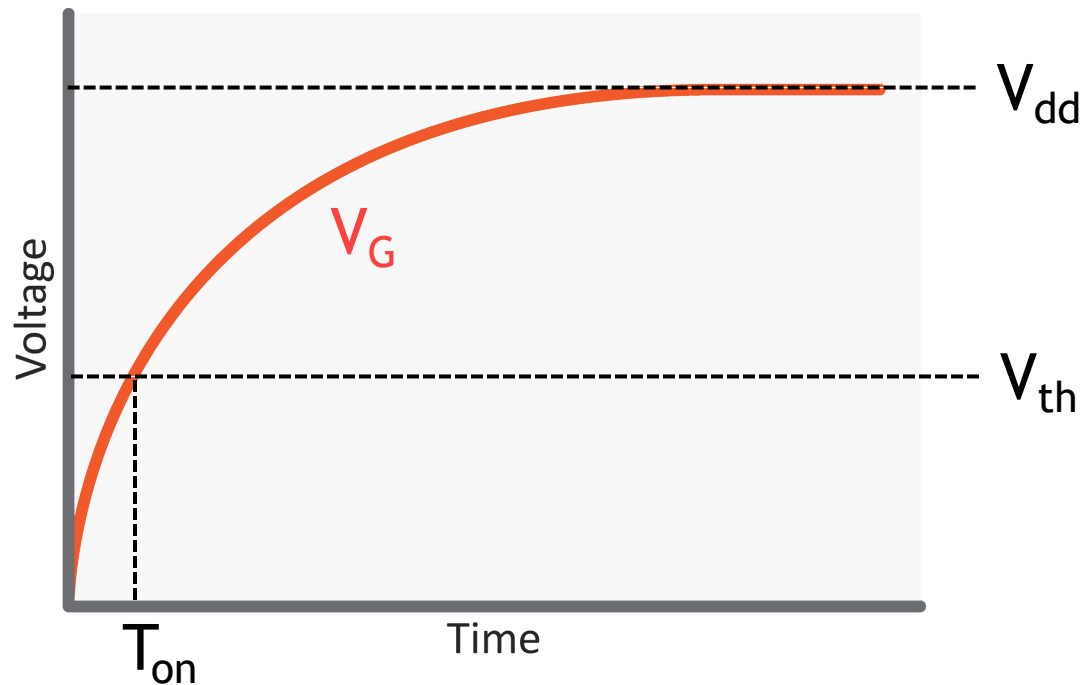


- Gate is switched on when V_G reaches a threshold V_{th}
 - By creating a channel in depletion region through field effect
 - V_{th} : threshold voltage (minimum voltage to create channel)



MOSFET 101

■ RC charging curve of V_G



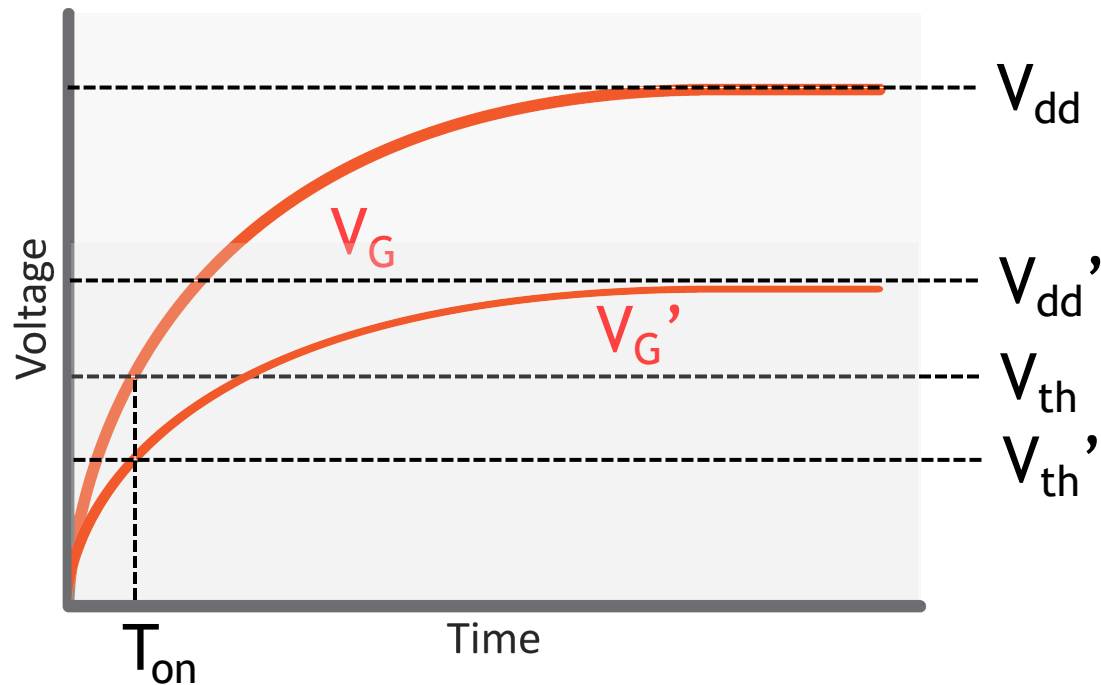
■ Speed (T_{on}) is determined by V_{dd} if V_{th} is fixed

- V_{dd} is the CPU supply voltage (the water pressure)
- If V_{dd} is lower, V_G will reach V_{th} more slowly (low pressure)



MOSFET 101

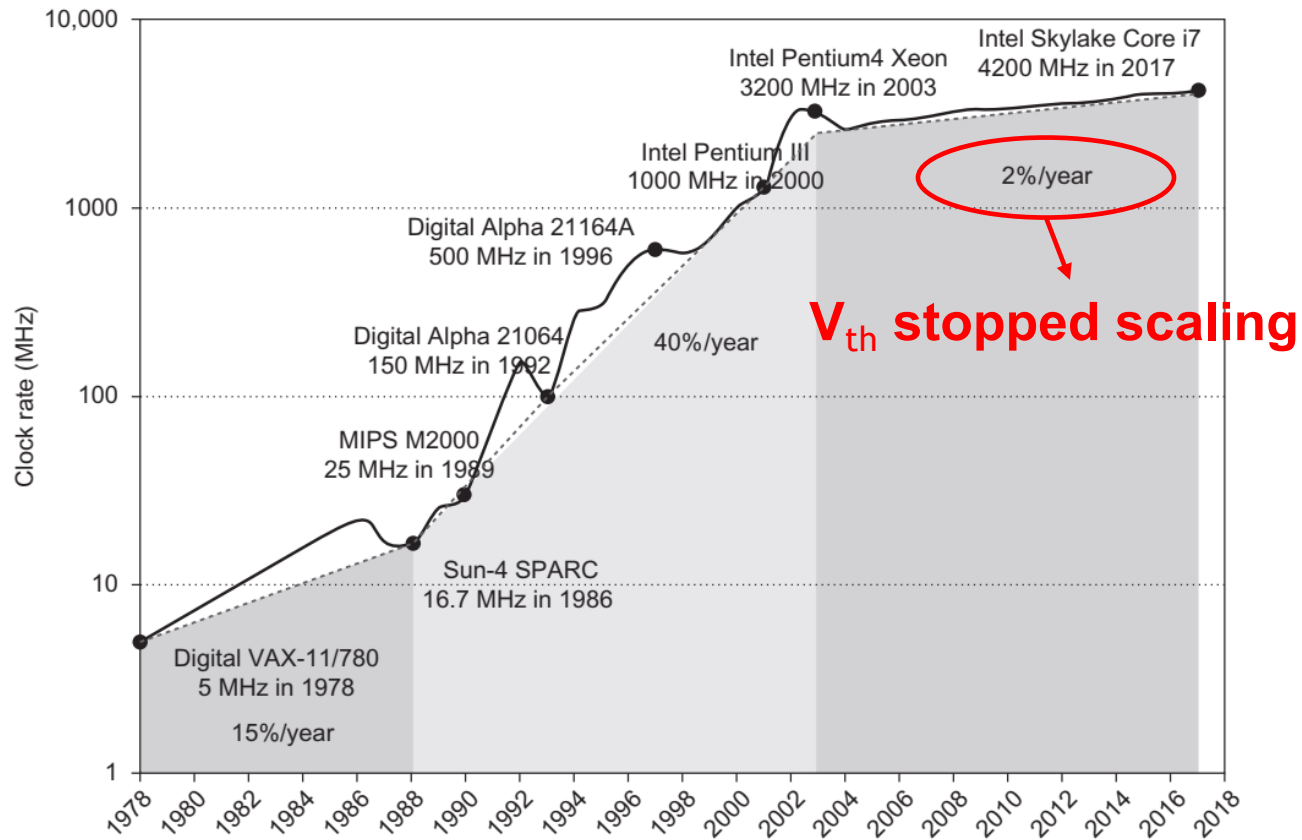
■ RC Charging Curve of V_G



- Speed (T_{on}) is maintained while reducing V_{dd} to V_{dd}' ,
only if V_{th} is also reduced to V_{th}'



End of Dennard Scaling



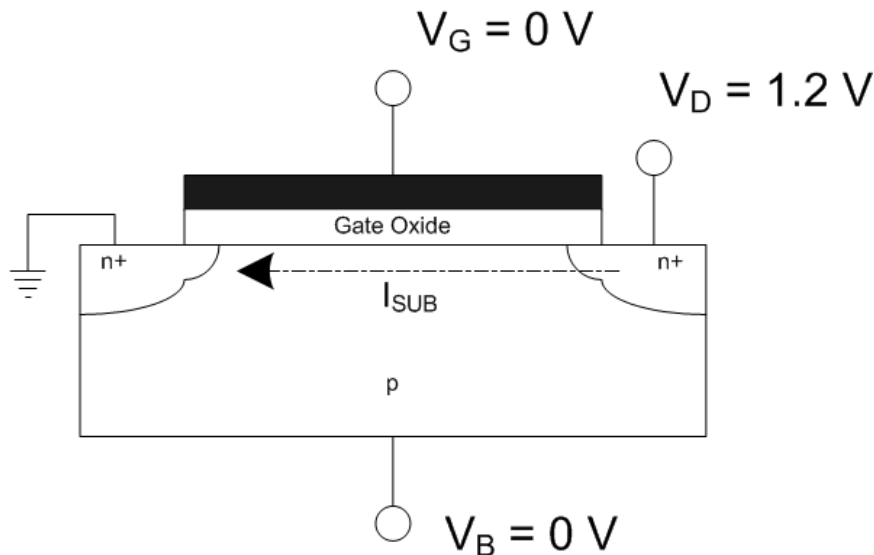
- And around 2003 is when Dennard Scaling ended



Limits to Dropping V_{th}

■ Subthreshold leakage

- Transistor leaks current even when gate is off ($V_G = 0$)

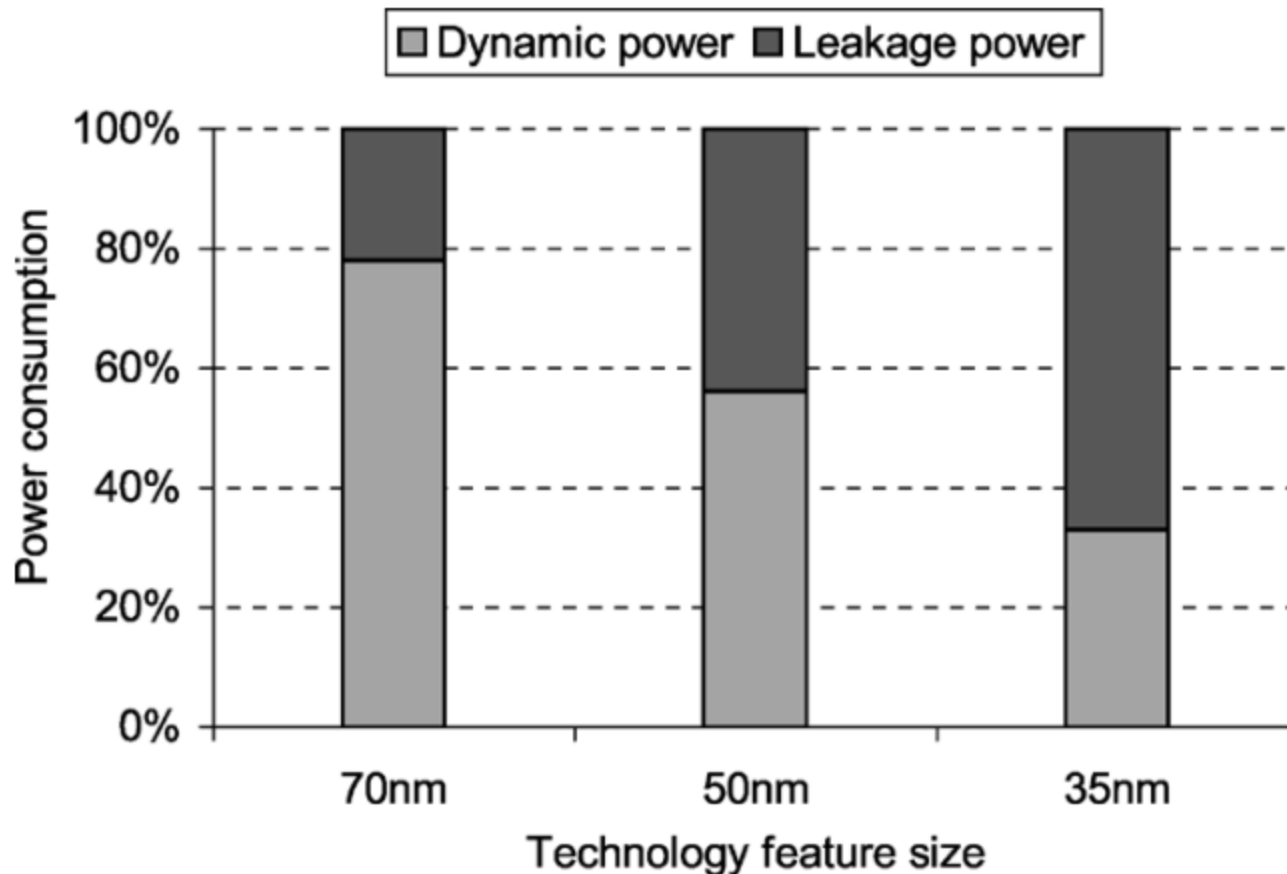


- This leakage current translates to leakage power
- Leakage worsens when V_{th} is dropped (related to oxide thickness)



Leakage Power across Generations

- Leakage power has increased across technology nodes



Source: L. Yan, Jiong Luo and N. K. Jha, "Joint dynamic voltage scaling and adaptive body biasing for heterogeneous distributed real-time embedded systems," in *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 24, no. 7, pp. 1030-1041, July 2005



End of Dennard Scaling

- Previous power calculation was incomplete
 - CPU power is the sum of both dynamic and leakage power
- $\text{Power}_{\text{CPU}} \propto \text{Power}_{\text{dynamic}} + \text{Power}_{\text{leakage}}$
 - $\text{Power}_{\text{dynamic}} \propto A * N * C F V_{\text{dd}}^2$
 - $\text{Power}_{\text{leakage}} \propto f(N, V_{\text{dd}}, V_{\text{th}}) \propto N * V_{\text{dd}} * e^{-V_{\text{th}}}$
 - Leakage worsens *exponentially* when V_{th} is dropped
 - Catch-22: when dropping V_{th} , $\text{Power}_{\text{dynamic}}$ ↓ but $\text{Power}_{\text{leakage}}$ ↑↑
- V_{th} can't be reduced further, so V_{dd} can't be reduced
- Dennard Scaling relies on reducing V_{dd} , so it's the end



“Dark Silicon” Rears its Head

- What happens to frequency without Dennard Scaling?
- $\text{Power}_{\text{dynamic}} (\propto A * N * CFV^2) + \text{Power}_{\text{leakage}} (\propto N * V * e^{-V_{th}})$
 - A = Activity factor
 - N = Number of transistors ($\propto 1/d^2$) ↑ ↑
 - C = Capacitance ($\propto d$) ↓
 - V = Supply Voltage \Leftrightarrow (Due to fixed V_{th})
 - F = Frequency ???
- To offset N, you actually have to *decrease* F
- Otherwise, if you want to maintain F, must decrease N
 - That is, you cannot power on all the transistors at any given point
 - Dark silicon: situation where chip is only partially powered



Free Ride is Over

- “Free” speed improvements from transistors is over
- Now it’s up to architects to improve performance
 - Moore’s Law is still alive (although slowing down)
 - Architects are flooded with extra transistors each generation
 - But it’s hard to even keep them powered without reducing F!
- Now is a good time to discuss technology constraints
 - Since we already mentioned a big one: TDP