



University of Pittsburgh

CS/COE 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?
 - Area of CPU chip?

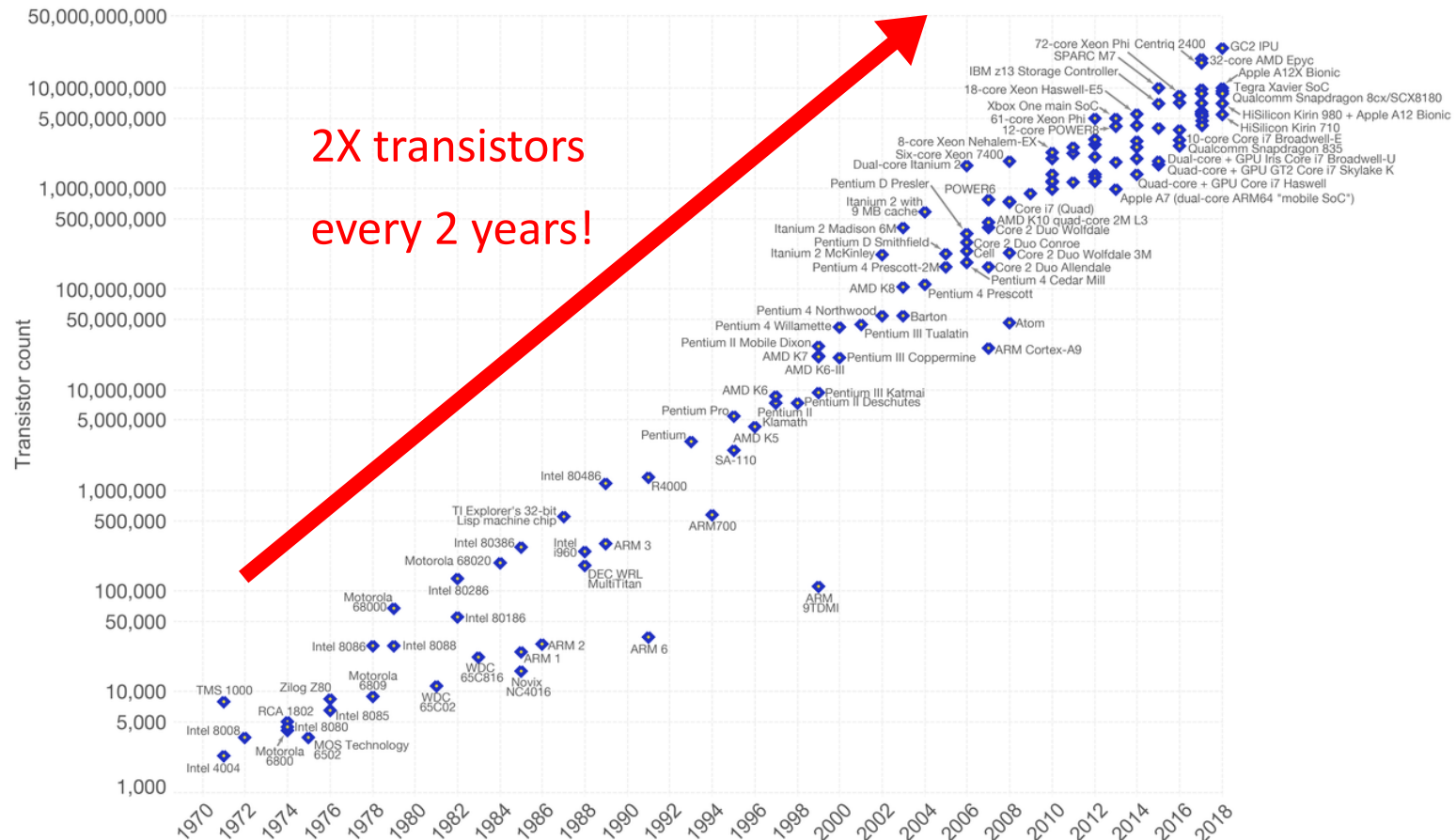


Moore's Law

Moore's Law – The number of transistors on integrated circuit chips (1971-2018)

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are linked to Moore's law.

OurWorld
in Data



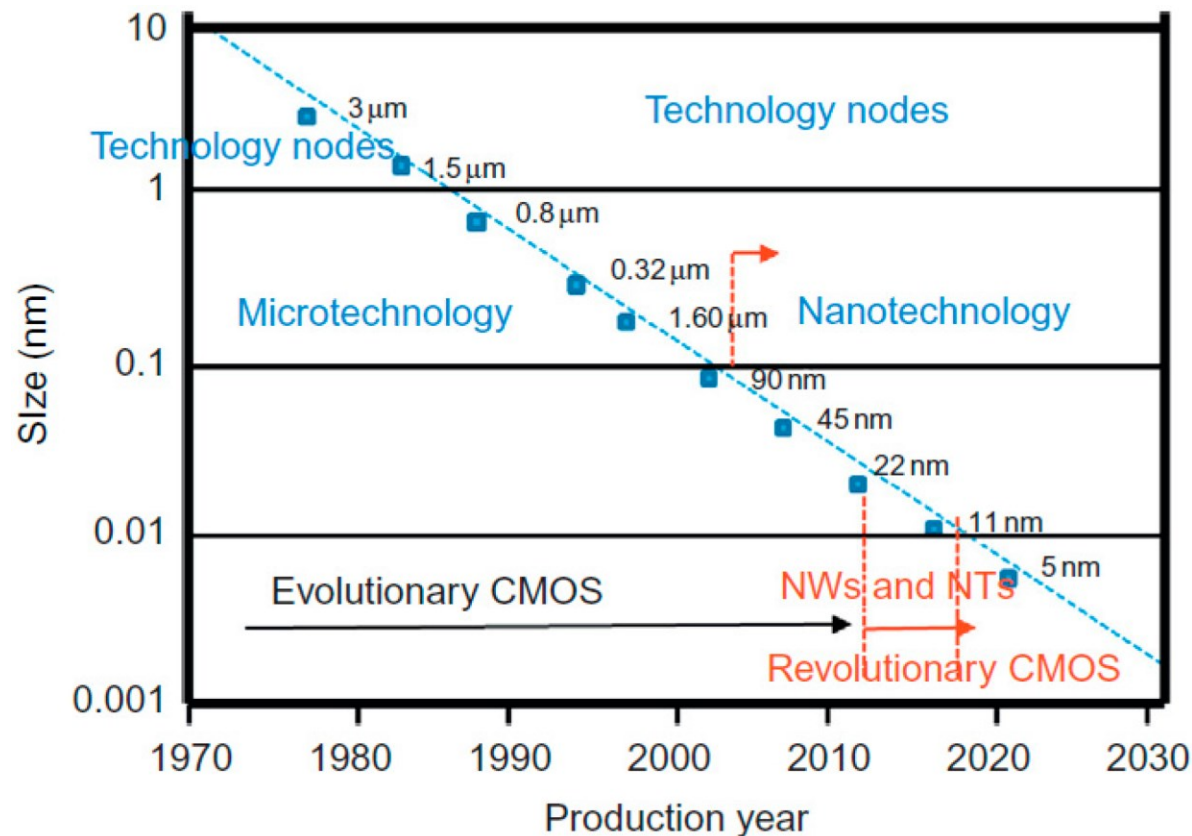
Data source: Wikipedia (https://en.wikipedia.org/wiki/Transistor_count)

The data visualization is available at [OurWorldinData.org](https://ourworldindata.org). There you find more visualizations and research on this topic.

Licensed under [CC-BY-SA](https://creativecommons.org/licenses/by-sa/4.0/) by the author Max Roser.



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)

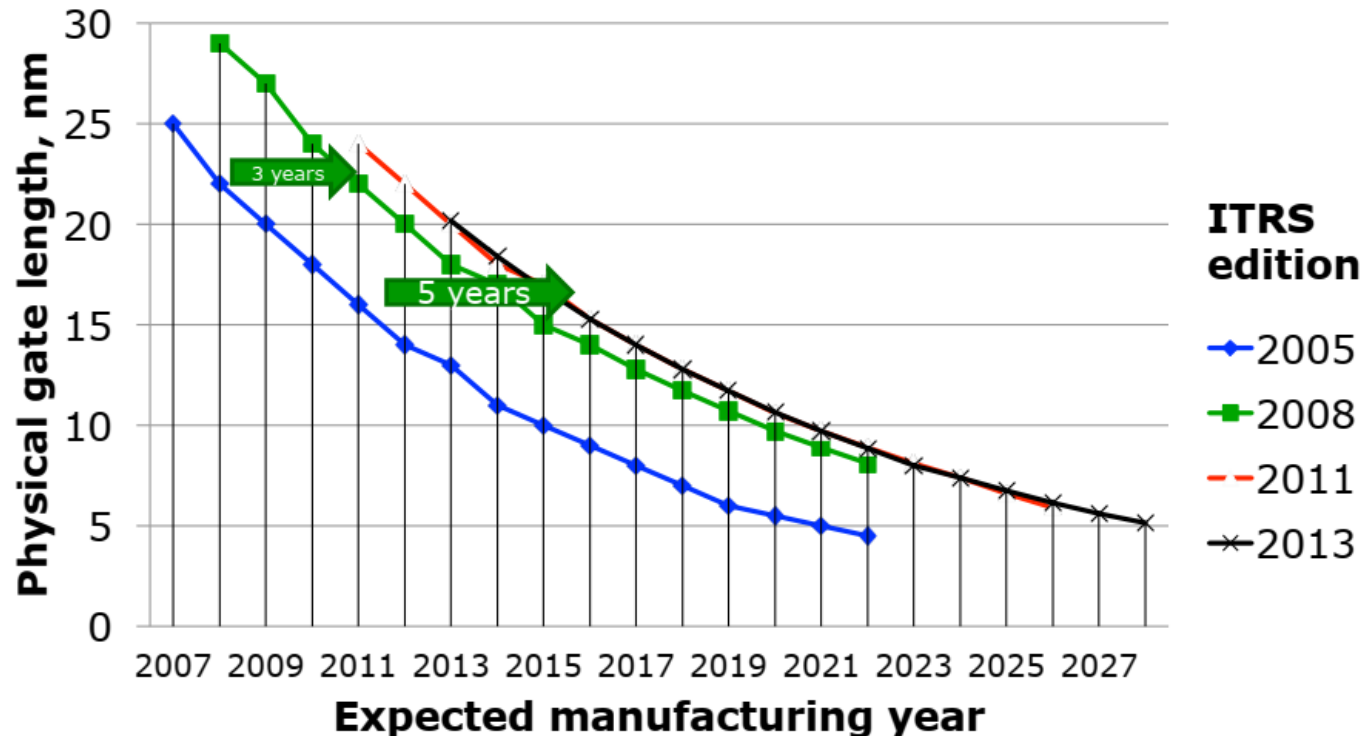
Year of Production	2018	2020	2022	2025	2028	2031	2034
Technology Node (nm)	7	5	3	2.1	1.5	1.0	0.7

- Looks like 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

- Million-dollar question:
Did Moore's Law result in higher performance CPUs?

- Please go to your respective Teams chat groups
 - But stay in the Zoom room and use only chat on Teams
 - To have chat content accessible to asynchronous students
- 1. Get to know each other
- 2. And then try to answer the following questions:
 - What do you think? Are CPUs getting faster?
 - If not, why do you think so? If yes, again why do you think so?
- 3. After 10 minutes, we will share discussions with class



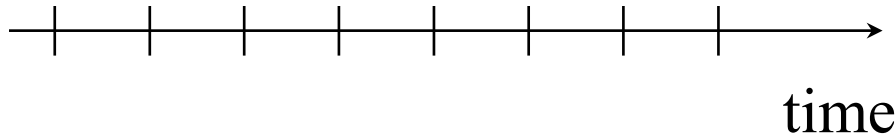
Are CPUs getting Faster?

- Yes!
- Clock speeds are increasing, power draw is decreasing - Andrew
- More cores - Jason
- Visual Studio compilation is actually faster - Nick
- No!
- Clock speeds are plateauing recently - Jason
- Seems stagnant from user's perspective e.g. Visual Studio - Josh



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 = $\frac{\text{seconds}}{\text{program}}$

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



Improving Execution Time

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

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

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

■ Improving $\frac{\text{cycles}}{\text{instructions}}$:

- 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

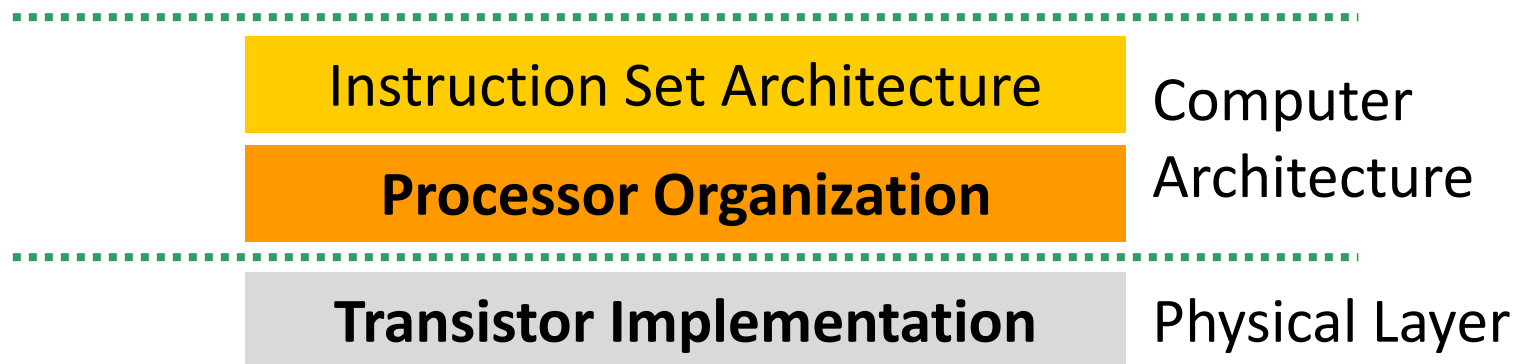
■ Improving $\frac{\text{instructions}}{\text{program}}$:

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



Moore's Law and Performance

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



- 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

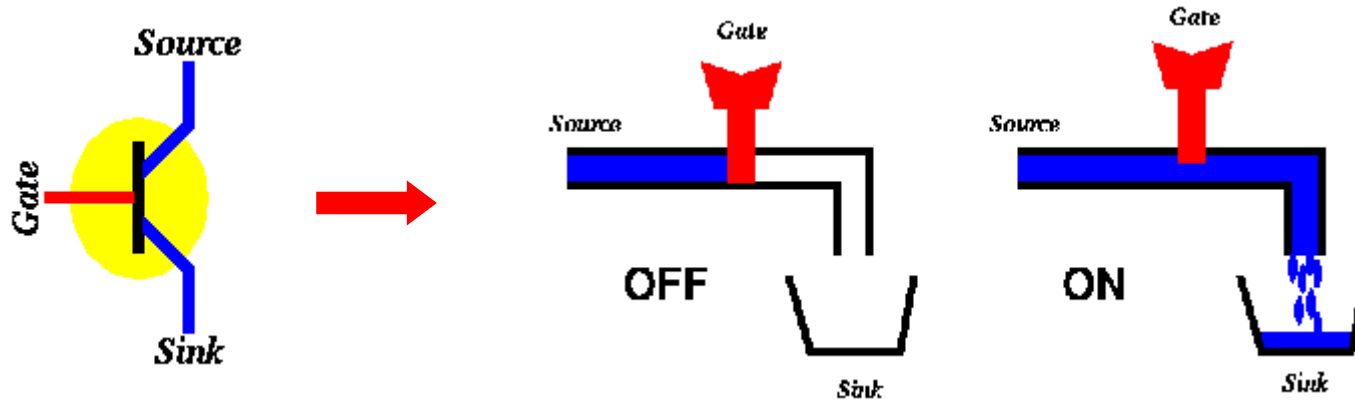
- CPU frequency is also impacted by transistor speed
 - As well as how many transistors are in between clock ticks (which is determined by processor organization)

- 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!



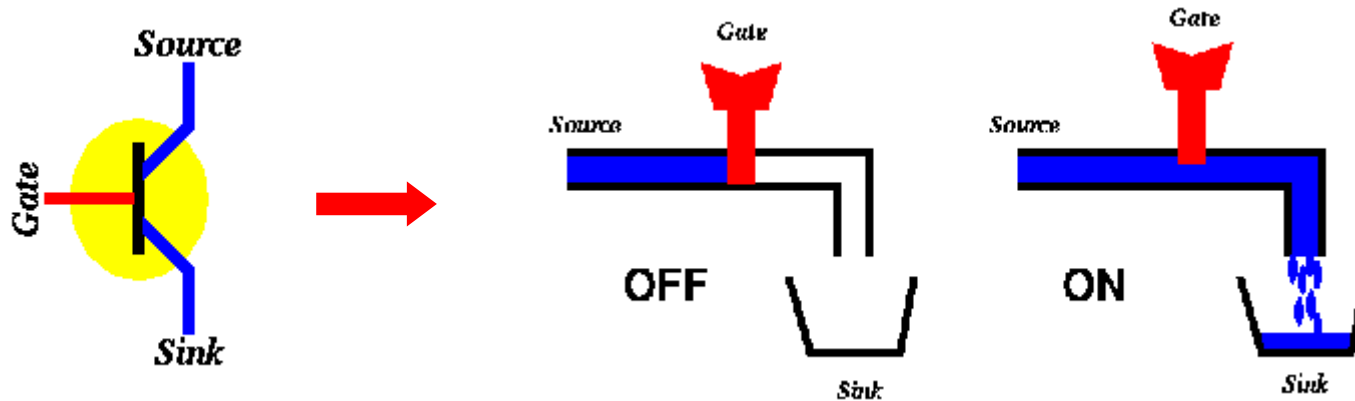
■ To make a transistor go fast, do one of the following:

- 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:

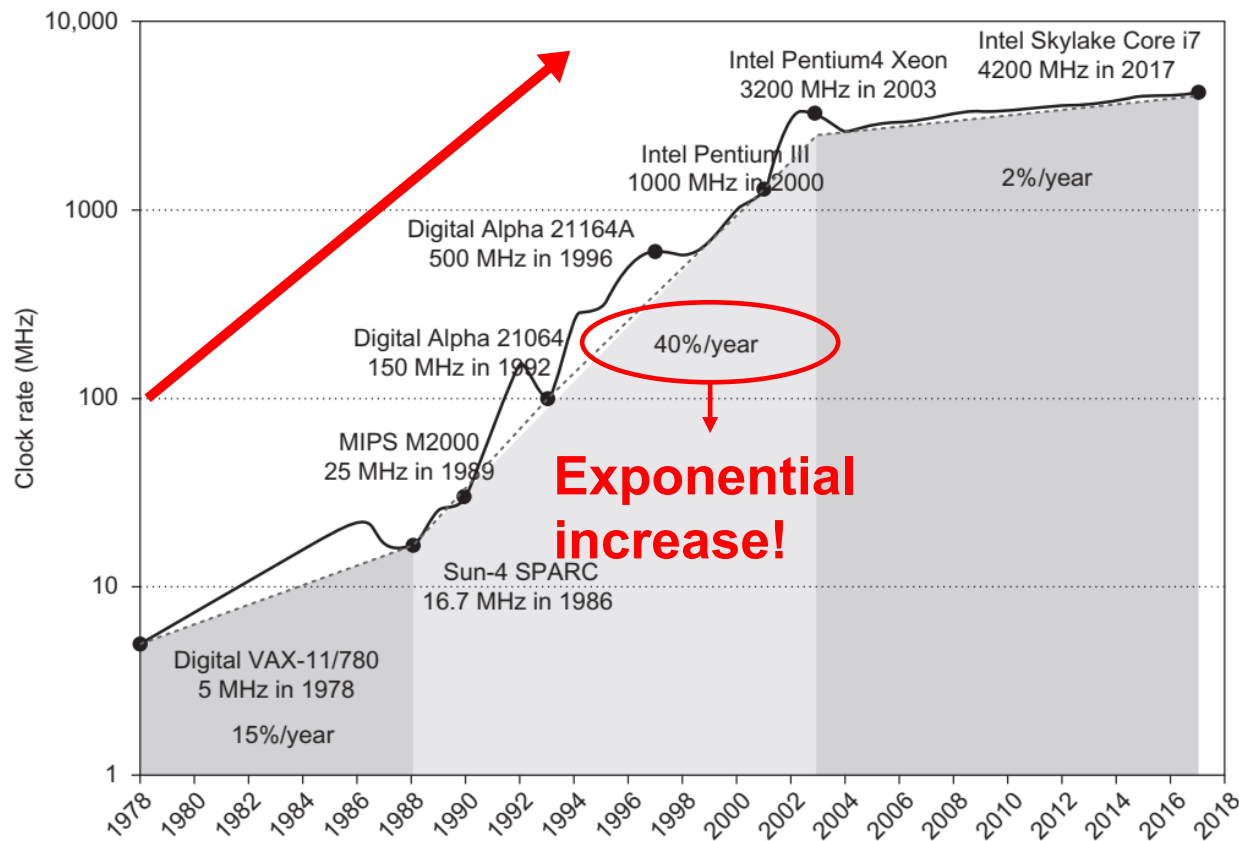
- *Channel length* (channel resistance) is reduced ↓
- *Capacitance* is reduced ↓

■ So, given the same *supply voltage*, smaller is faster!

■ So, did Moore's Law enjoy faster and faster frequencies?



Yes, for a while ...



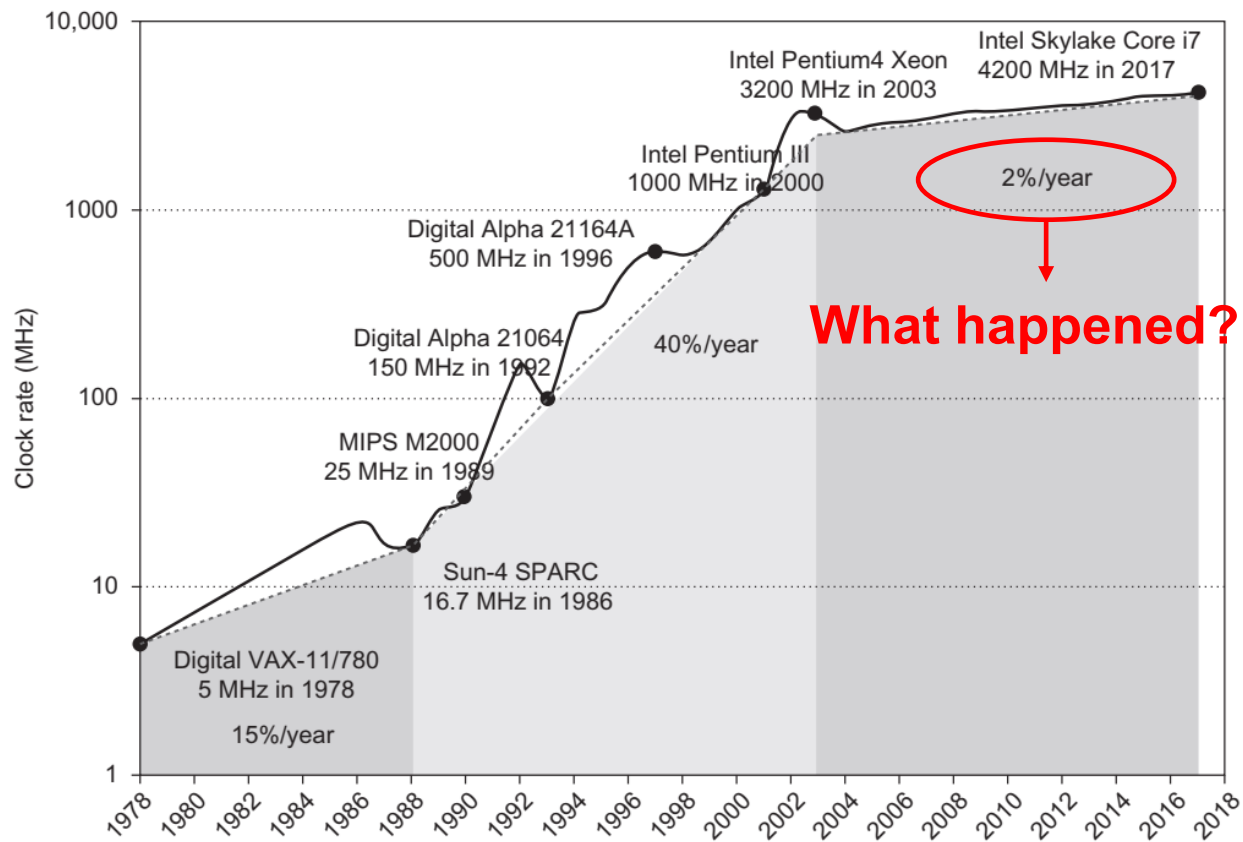
Source: Computer Architecture, A Quantitative Approach (6th ed.) by John Hennessy and David Patterson, 2017

■ Improvements in large part due to transistors

- Processor design also contributed but we'll discuss later



But not so much lately

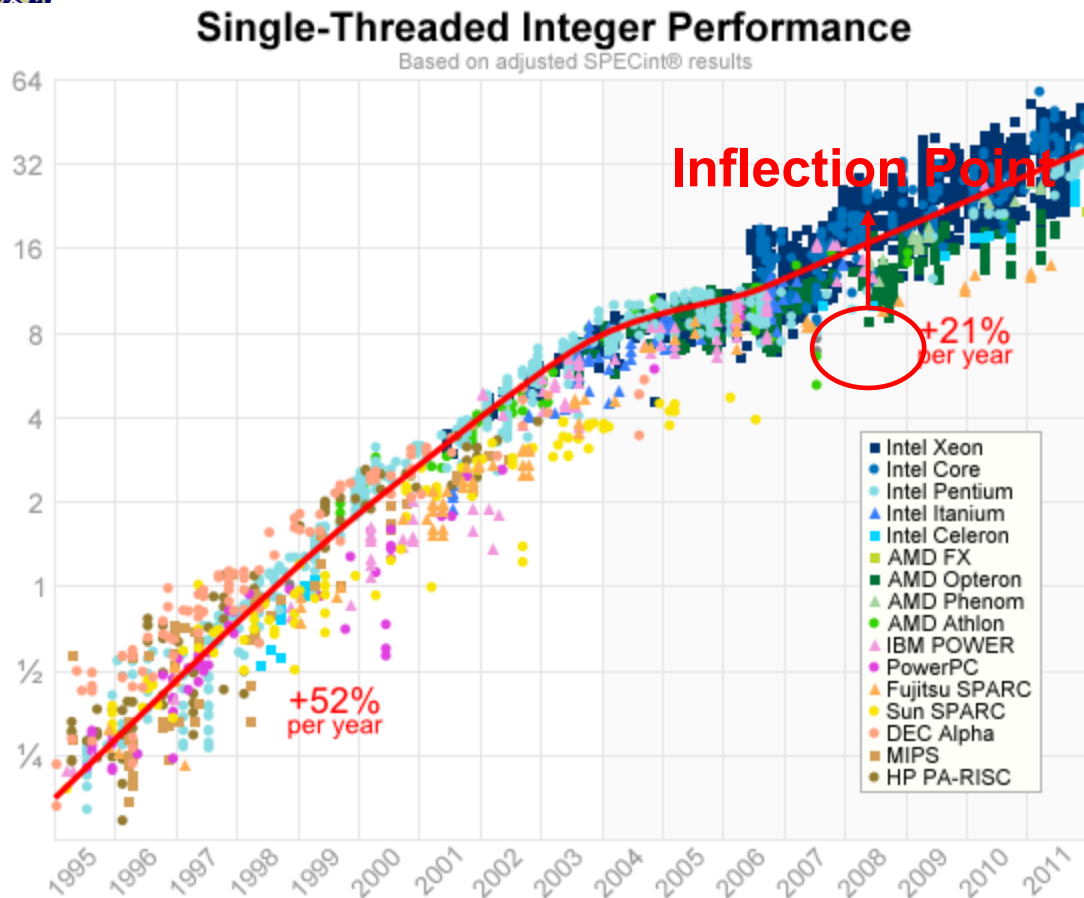


Source: Computer Architecture, A Quantitative Approach (6th ed.) by John Hennessy and David Patterson, 2017

- 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 heat (power) that cooling system can handle
- Cooling system hasn't improved much over generations (Typically a CPU cooling fan attached with thermal paste)

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

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



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



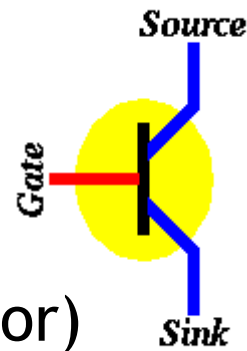
TDP and Moore's Law

- CPU Power $\propto A * N * CFV^2$ with Moore's Law
 - A = Activity factor
 - N = Number of transistors ↑ ↑
 - C = Capacitance (\propto transistor size) ↓
 - F = CPU frequency ↑ (thanks to reductions in transistor size)
 - V = CPU Supply Voltage

- Reductions in C cannot offset increases in N and F
 - Q) How did CPU frequency keep increasing up to 2003?
 - A) By maintaining power through reductions in Voltage ↓
 - Q) Wait! Voltage reduction reduces frequency! (Transistor 101)
 - A) Alright, time to do MOSFET 101

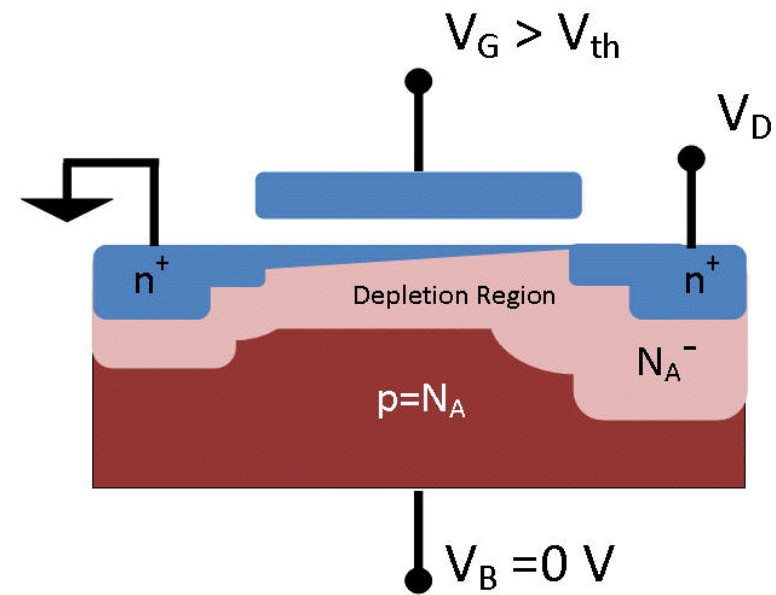
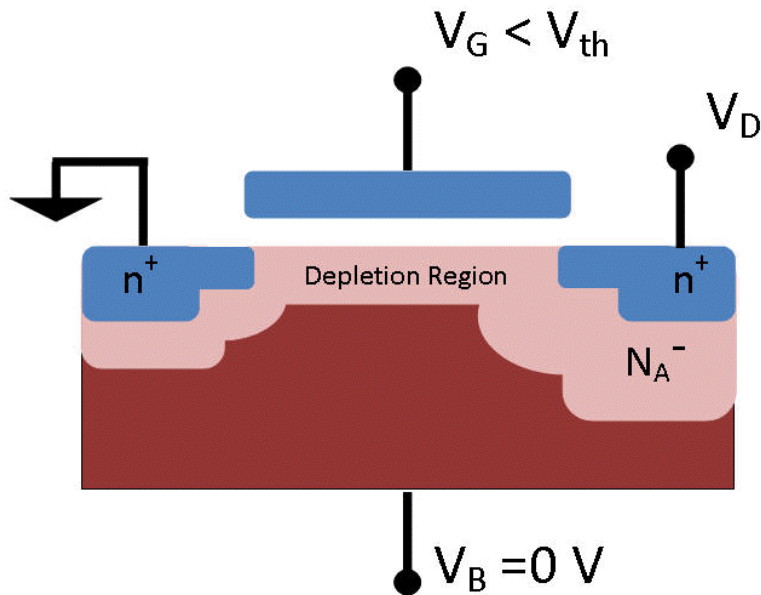


MOSFET 101



■ MOSFET (Metal Oxide Silicon Field Effect Transistor)

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

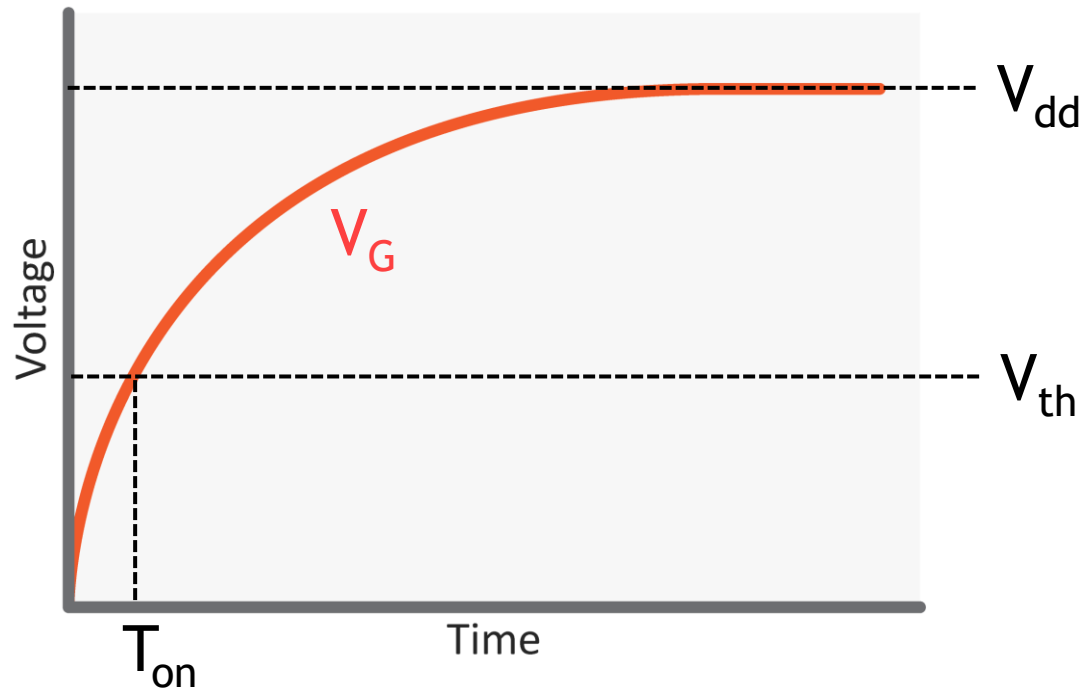


- Gate is switched on when V_G reaches a threshold V_{th}
 - By creating a channel in depletion region through field effect
 - V_{th} : threshold voltage (minimum voltage to create channel)



MOSFET 101

■ RC charging curve of V_G



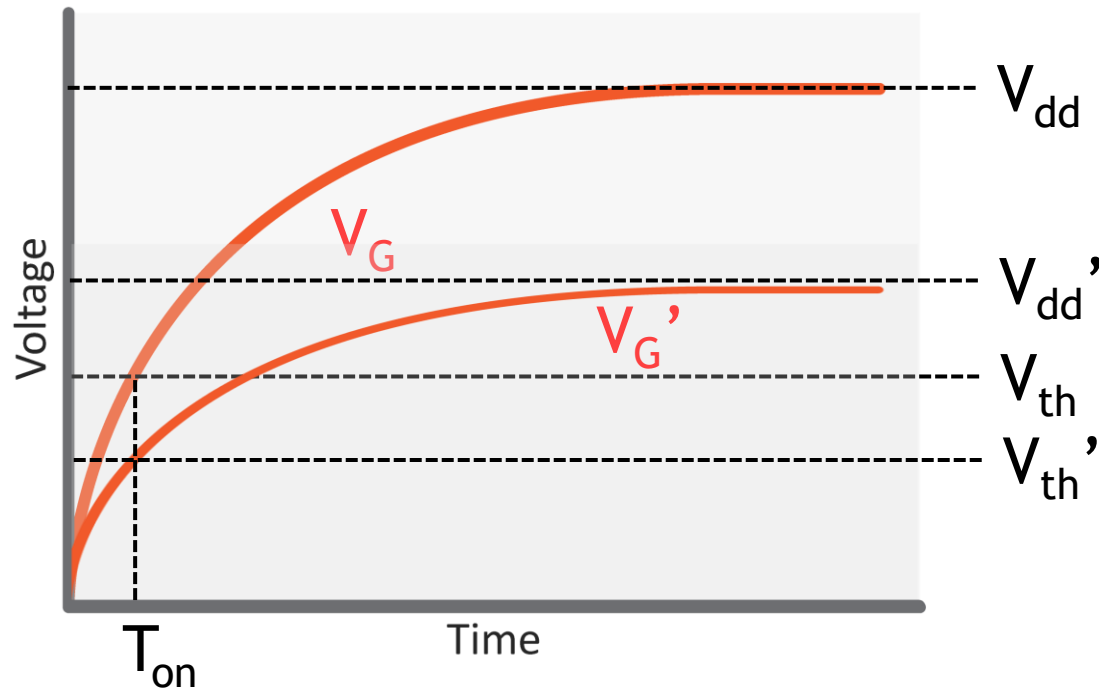
■ Speed (T_{on}) is determined by V_{dd} if V_{th} is fixed

- 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}' , if V_{th} is also reduced to V_{th}'

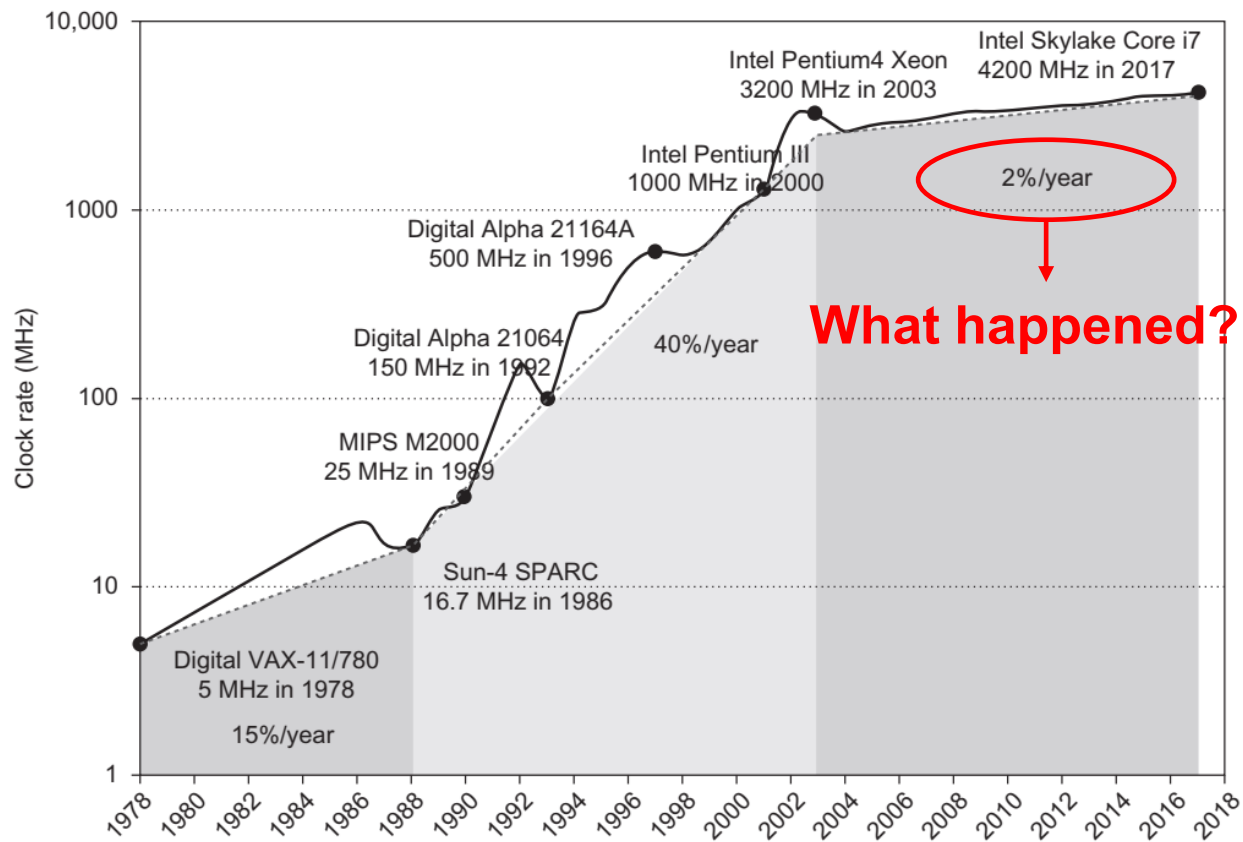


Dennard Scaling

- So in the end, this is what happens ...
- CPU Power $\propto A * N * CFV^2$ with Moore's Law
 - A = Activity factor
 - N = Number of transistors ↑ ↑
 - C = Capacitance (\propto transistor size) ↓
 - F = CPU frequency ↑ (thanks to reductions in transistor size)
 - V = CPU Supply Voltage ↓ (to reduce power)
 - V_{th} = CPU Threshold Voltage ↓ (to maintain frequency)
- Factors balance each other out to make power constant
- This recipe for scaling frequency while keeping power constant is called Dennard Scaling



End of Dennard Scaling



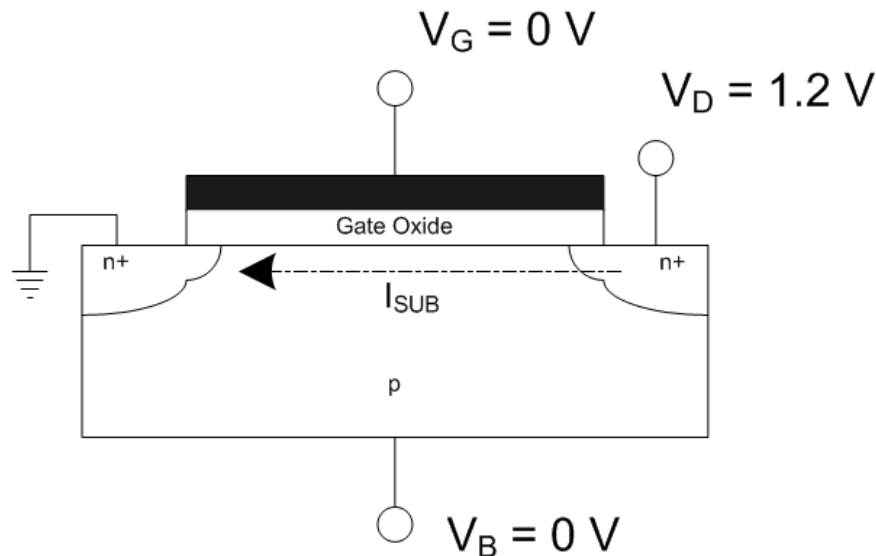
- And around 2003 is when Dennard Scaling ended



Limits to Dropping V_{th}

■ Subthreshold leakage

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

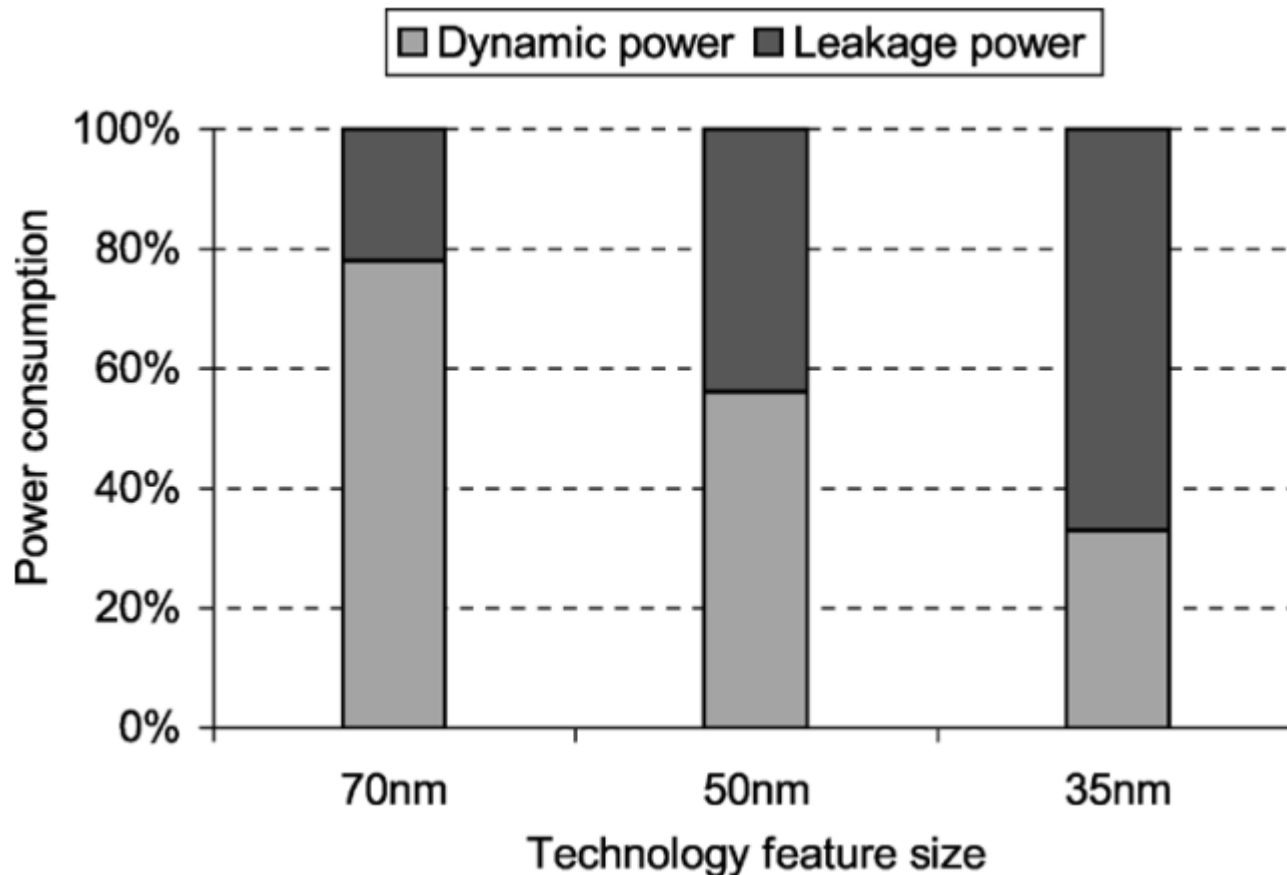


- This leakage current translates to leakage power
- Leakage worsens when V_{th} is dropped (related to oxide thickness)



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

■ $\text{Power}_{\text{CPU}} \propto \text{Power}_{\text{dynamic}} + \text{Power}_{\text{leakage}}$

- $\text{Power}_{\text{dynamic}} \propto A * N * CFV^2$
- $\text{Power}_{\text{leakage}} \propto f(N, V, V_{\text{th}}) \propto N * V * e^{-V_{\text{th}}}$
- Leakage worsens *exponentially* when V_{th} is dropped
- Catch-22: when dropping V_{th} , $\text{Power}_{\text{dynamic}}$ ↓ but $\text{Power}_{\text{leakage}}$ ↑↑
- That means V_{th} can't be reduced and V can't be reduced

■ $\text{Power}_{\text{dynamic}} (\propto A * N * CFV^2) + \text{Power}_{\text{leakage}} (\propto N * V * e^{-V_{\text{th}}})$

- A = Activity factor
- N = Number of transistors ↑↑
- C = Capacitance (\propto transistor size) ↓
- F = CPU frequency \Leftrightarrow (Can't increase without violating TDP)
- V = CPU Supply Voltage \Leftrightarrow (Due to fixed V_{th})
- V_{th} = CPU Threshold Voltage \Leftrightarrow (Due to leakage power)



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 and well (although slowing down)
 - Architects are flooded with extra transistors each generation
- Now is a good time to discuss technology constraints
 - Since we already mentioned a big one: TDP