

# **CSC 252: Computer Organization**

## **Spring 2020: Lecture 25**

Instructor: Yuhao Zhu

Department of Computer Science  
University of Rochester

# Announcements

- Lab 5: <https://www.cs.rochester.edu/courses/252/spring2020/labs/assignment5.html>

19	20	21	22	23	24	25
26	27	28 <b>Today Last Lecture</b>	29	30	May 1	2
3	4 <b>Due</b>	5	6 <b>Final</b>	7	8	9

# Announcements

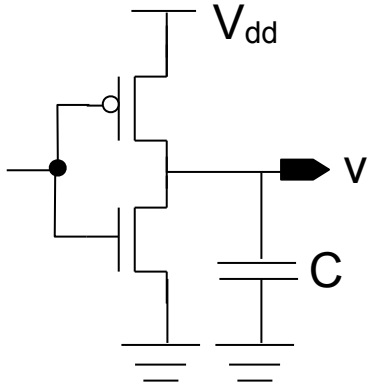
- Grades for Lab 4 are out. Let the TAs know if you have questions
- Final exam: May 6, 19:15 EST — 22:15 EST
- Let me know if you can't make this time.
- Exam will be electronic on blackboard, but we will send you an PDF version so that you can work offline in case
  - 1) you don't have Internet access at the exam time or
  - 2) you lose Internet access.
  - Write down the answers on a scratch paper, take pictures, and send us the pictures
- Same rule as before: anything on paper is fine, nothing electronic other than using the computer to take the exam
- **"I don't know"** gets 15% partial credit. Must cross/erase everything else.

# Announcements

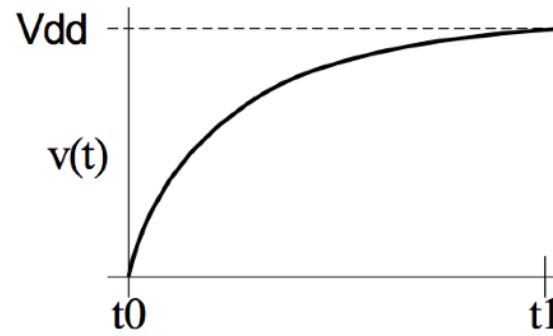
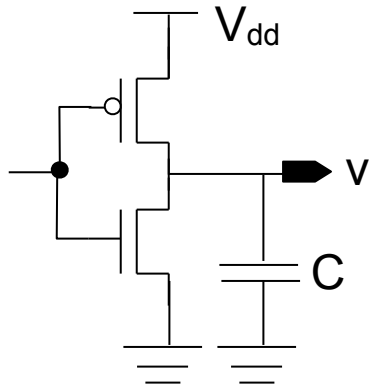
- Now the dry run on blackboard



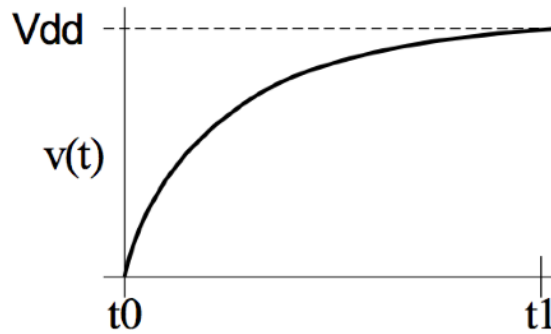
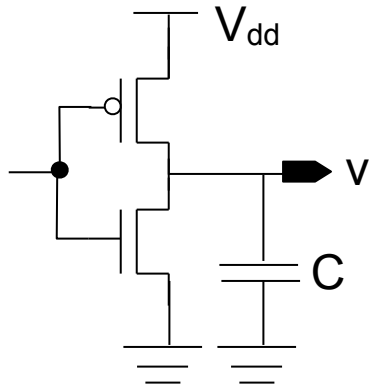
# Dynamic Power



# Dynamic Power

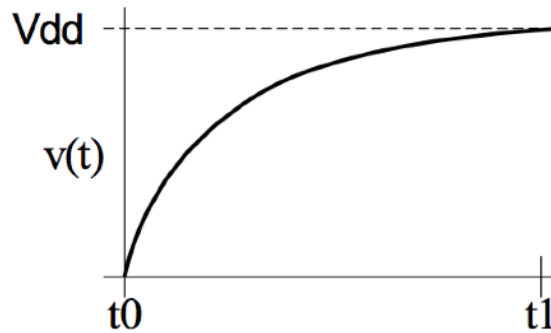
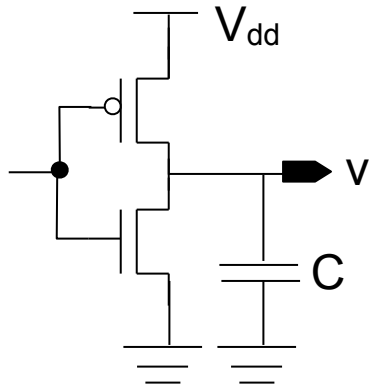


# Dynamic Power



$$\begin{aligned}
 E_{sw} &= \int_{t_0}^{t_1} P(t) dt = \int_{t_0}^{t_1} (V_{dd} - v) \cdot i(t) dt = \int_{t_0}^{t_1} (V_{dd} - v) \cdot c (dv/dt) dt = \\
 &= cV_{dd} \int_{t_0}^{t_1} dv - c \int_{t_0}^{t_1} v \cdot dv = cV_{dd}^2 - 1/2 cV_{dd}^2 = \boxed{1/2 cV_{dd}^2}
 \end{aligned}$$

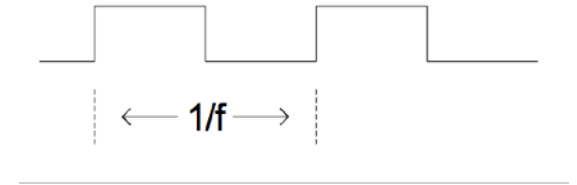
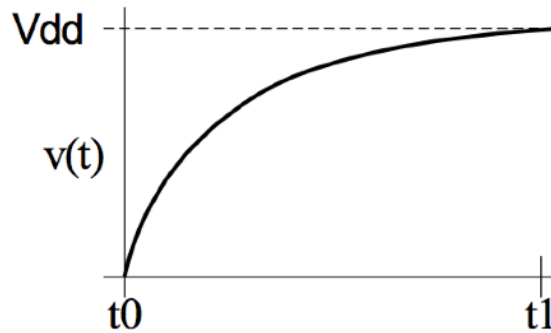
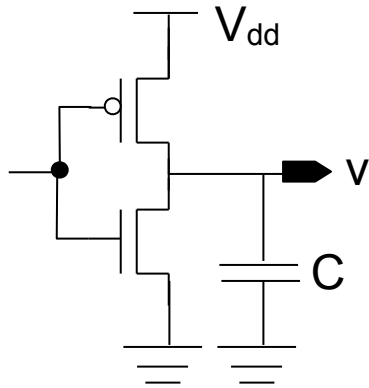
# Dynamic Power



$$\begin{aligned}
 E_{sw} &= \int_{t_0}^{t_1} P(t) dt = \int_{t_0}^{t_1} (V_{dd} - v) \cdot i(t) dt = \int_{t_0}^{t_1} (V_{dd} - v) \cdot c (dv/dt) dt = \\
 &= cV_{dd} \int_{t_0}^{t_1} dv - c \int_{t_0}^{t_1} v \cdot dv = cV_{dd}^2 - 1/2 cV_{dd}^2 = \boxed{1/2 cV_{dd}^2}
 \end{aligned}$$

Energy dissipated for every transition (0->1 or 1->0):  $\frac{1}{2} C_L V_{dd}^2$

# Dynamic Power



$$\begin{aligned}
 E_{sw} &= \int_{t_0}^{t_1} P(t) dt = \int_{t_0}^{t_1} (V_{dd} - v) \cdot i(t) dt = \int_{t_0}^{t_1} (V_{dd} - v) \cdot c (dv/dt) dt = \\
 &= cV_{dd} \int_{t_0}^{t_1} dv - c \int_{t_0}^{t_1} v \cdot dv = cV_{dd}^2 - 1/2 cV_{dd}^2 = \boxed{1/2 cV_{dd}^2}
 \end{aligned}$$

Energy dissipated for every transition (0->1 or 1->0):  $\frac{1}{2} C_L V_{dd}^2$

$$\mathbf{P = \alpha \cdot (E / T) = \alpha \cdot E f = \frac{1}{2} \alpha C V_{dd}^2 f}$$

$\alpha$ : switch activity factor. No switching, no dynamic power consumption

# Dynamic Power

$$P = k C V^2 f$$

# Dynamic Power

$$P = k C V^2 f$$

- Frequency  $f$  is proportional to Voltage  $V$

# Dynamic Power

$$P = k C V^2 f$$

- Frequency  $f$  is proportional to Voltage  $V$ 
  - Intuitively: higher voltage moves electrons faster, so the clock speed can go up also (“overclocking” just increases clock speed without increasing voltage => machine might crash)



# Dynamic Power

$$P = k C V^2 f$$

- Frequency  $f$  is proportional to Voltage  $V$ 
  - Intuitively: higher voltage moves electrons faster, so the clock speed can go up also (“overclocking” just increases clock speed without increasing voltage => machine might crash)
  - 15% reduction in voltage requires about 15% slow down in frequency

# Dynamic Power

$$P = k C V^2 f$$

- Frequency  $f$  is proportional to Voltage  $V$ 
  - Intuitively: higher voltage moves electrons faster, so the clock speed can go up also (“overclocking” just increases clock speed without increasing voltage => machine might crash)
  - 15% reduction in voltage requires about 15% slow down in frequency
  - What’s the impact on dynamic power?  $0.85^3 \approx 60\%$  -> 40% dynamic power reduction.

# Multi-core Reduces Dynamic Power

$$P = k C f^3$$

# Multi-core Reduces Dynamic Power

$$P = k C f^3$$

- Dynamic power favors parallel processing over higher clock rate
  - Take a core and replicate it 4 times: 4x speedup & **4x power**

# Multi-core Reduces Dynamic Power

$$P = k C f^3$$

- Dynamic power favors parallel processing over higher clock rate
  - Take a core and replicate it 4 times: 4x speedup & **4x power**
  - Take a core and clock it 4 times faster: 4x speedup but **64x dynamic power!**

# Multi-core Reduces Dynamic Power

$$P = k C f^3$$

- Dynamic power favors parallel processing over higher clock rate
  - Take a core and replicate it 4 times: 4x speedup & **4x power**
  - Take a core and clock it 4 times faster: 4x speedup but **64x dynamic power!**
- Another way to think about this

# Multi-core Reduces Dynamic Power

$$P = k C f^3$$

- Dynamic power favors parallel processing over higher clock rate
  - Take a core and replicate it 4 times: 4x speedup & **4x power**
  - Take a core and clock it 4 times faster: 4x speedup but **64x dynamic power!**
- Another way to think about this
  - If a task can be perfectly parallelized by 4 cores, we can reduce the clock frequency of each core to 1/4 while retaining the same performance

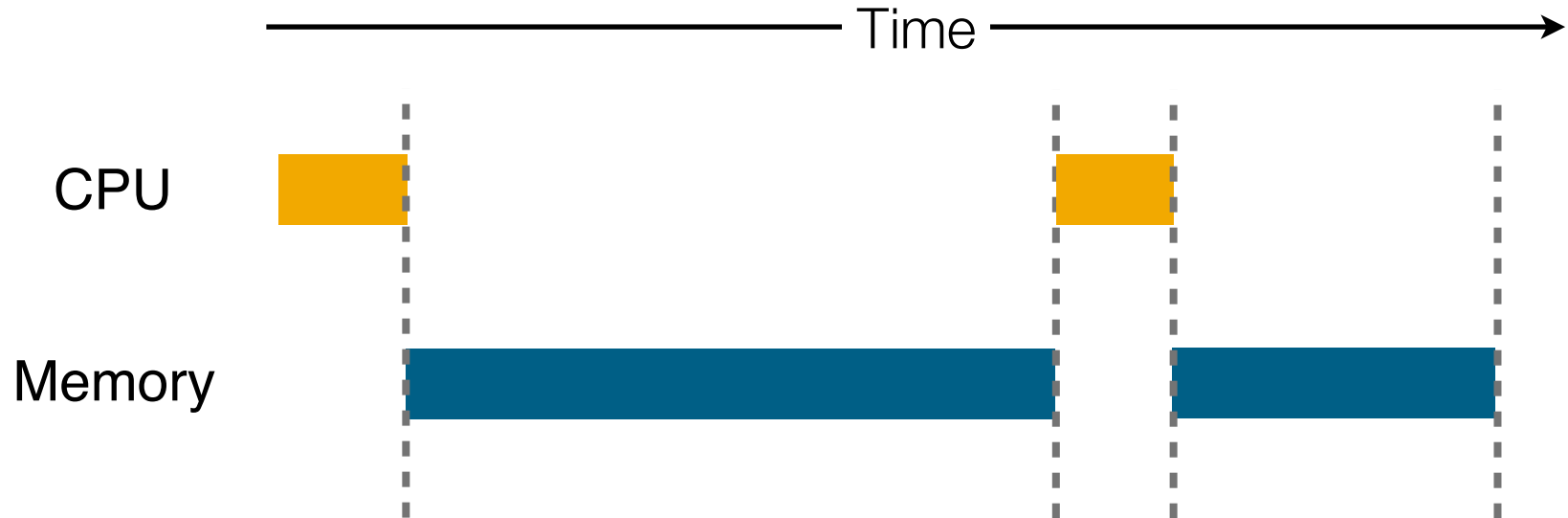
# Multi-core Reduces Dynamic Power

$$P = k C f^3$$

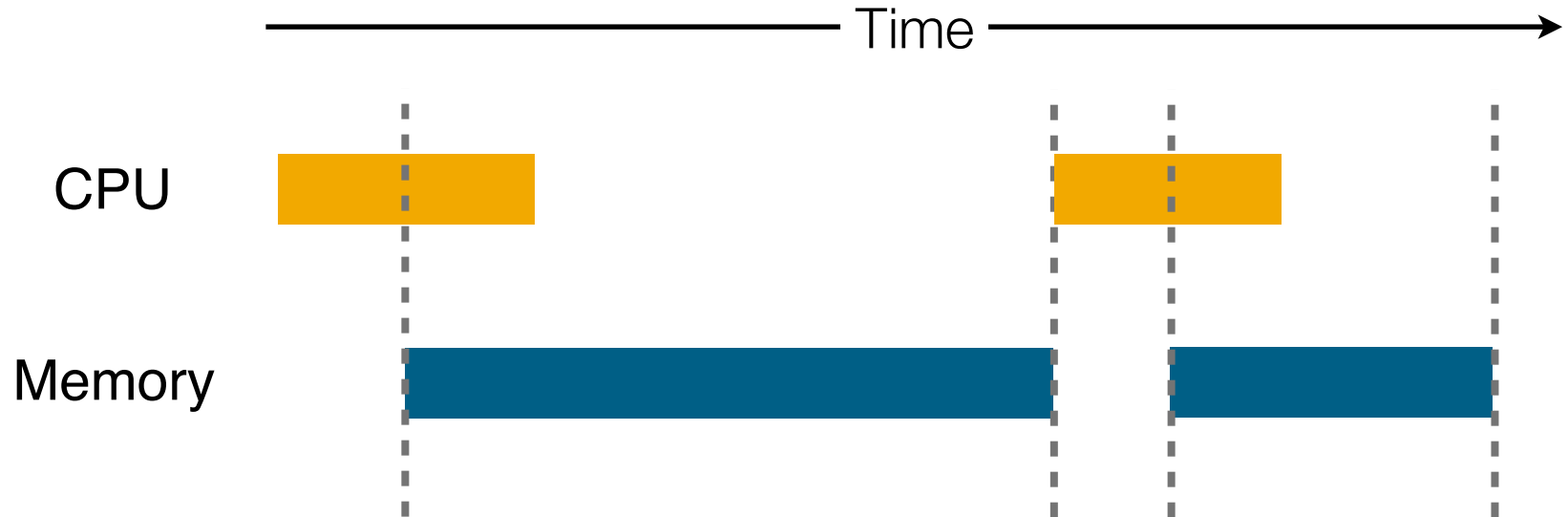
- Dynamic power favors parallel processing over higher clock rate
  - Take a core and replicate it 4 times: 4x speedup & **4x power**
  - Take a core and clock it 4 times faster: 4x speedup but **64x dynamic power!**
- Another way to think about this
  - If a task can be perfectly parallelized by 4 cores, we can reduce the clock frequency of each core to 1/4 while retaining the same performance
  - Dynamic power becomes  $4 \times (1/4)^3 = 1/16$



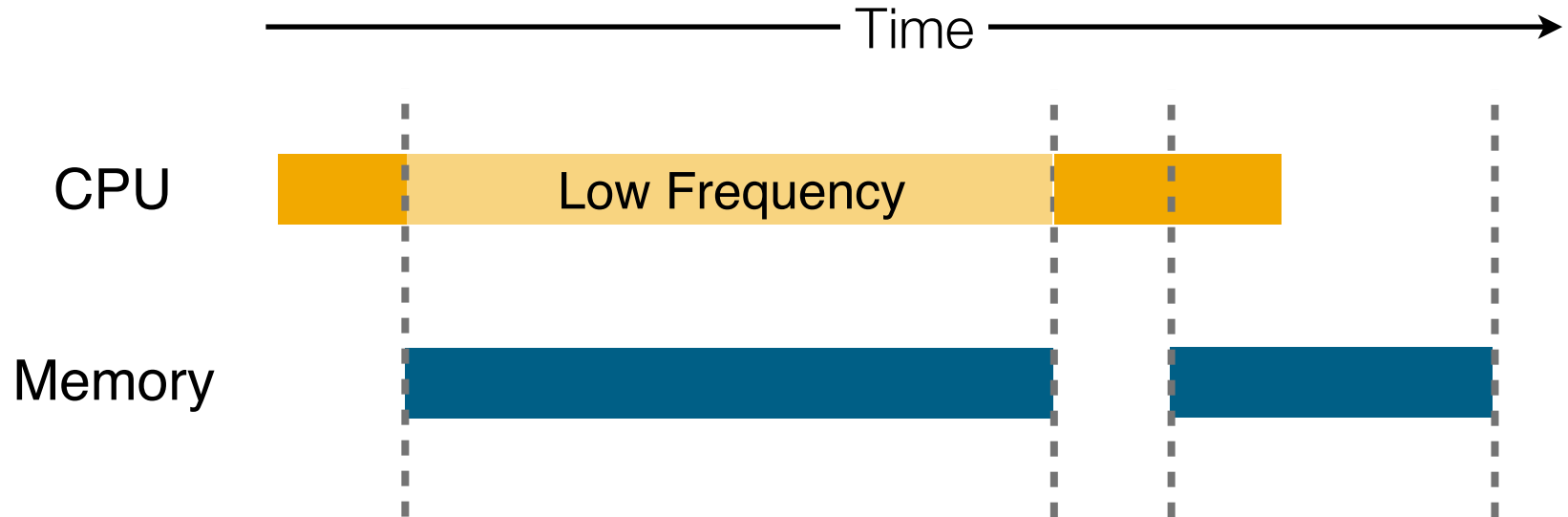
# Dynamic Voltage and Frequency Scaling



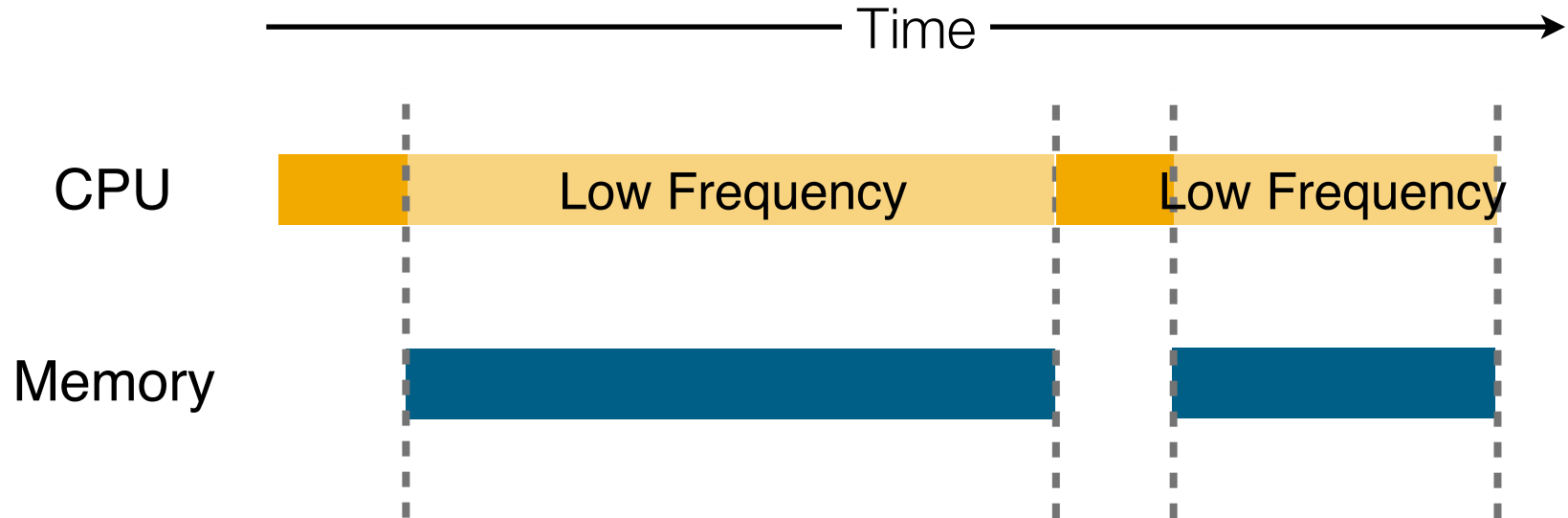
# Dynamic Voltage and Frequency Scaling



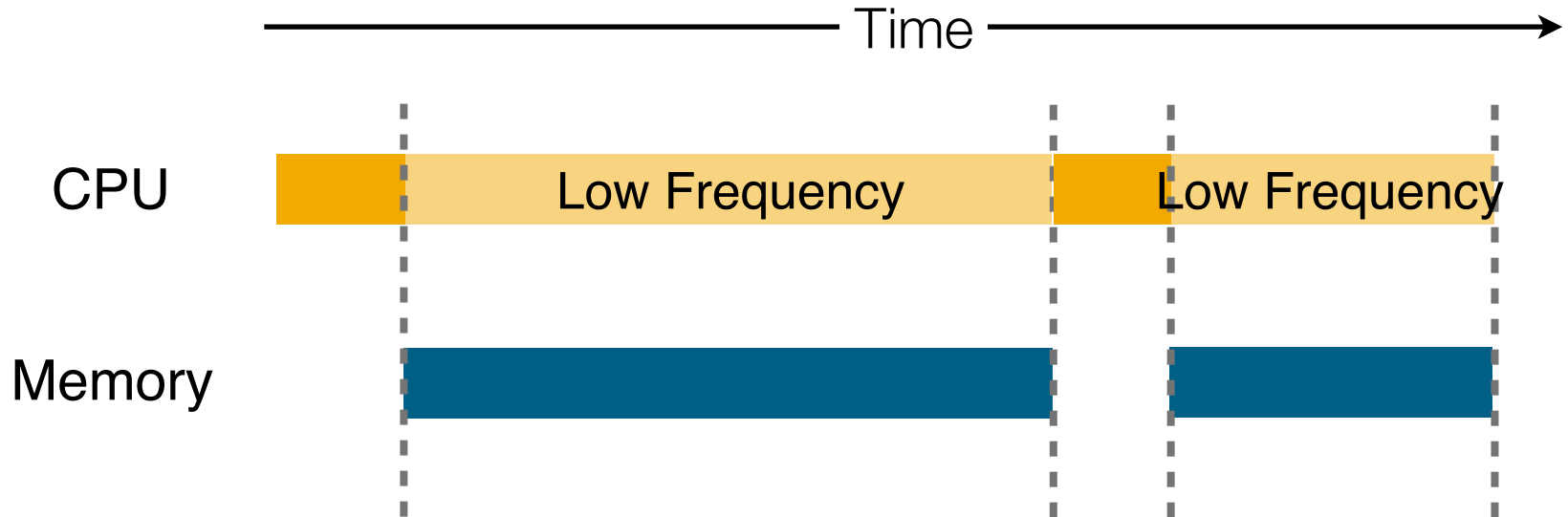
# Dynamic Voltage and Frequency Scaling



# Dynamic Voltage and Frequency Scaling

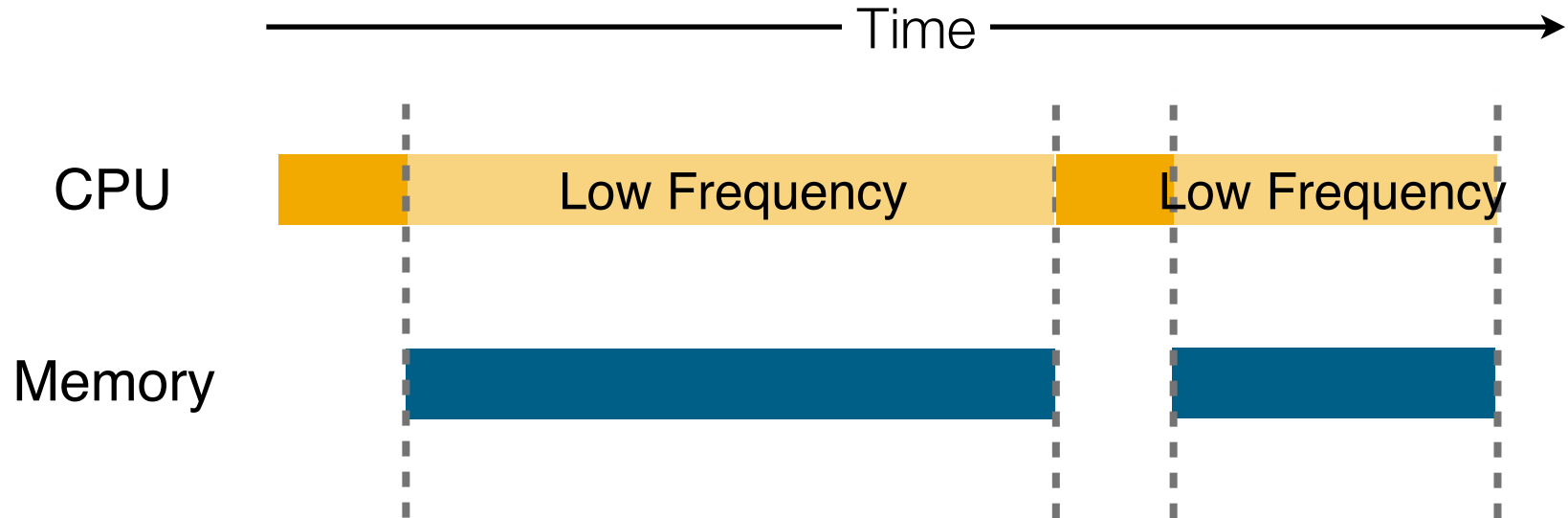


# Dynamic Voltage and Frequency Scaling



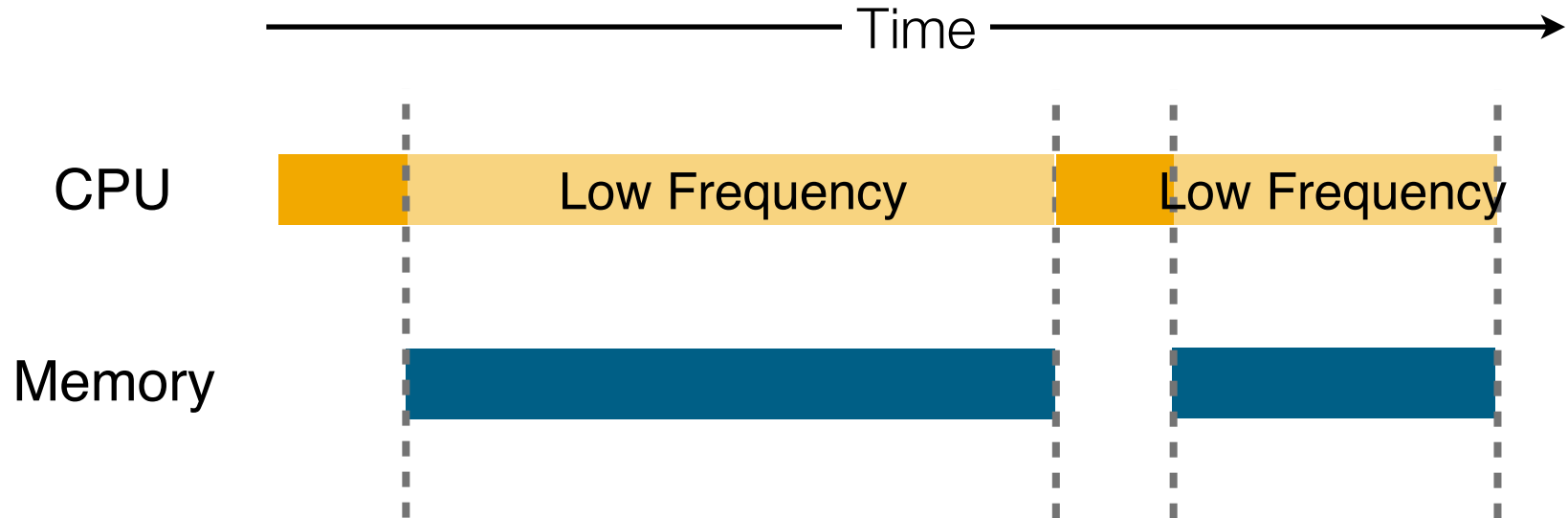
- DVFS: dynamic scales voltage and frequency depending on how much slack there is to reduce power with little performance impact

# Dynamic Voltage and Frequency Scaling



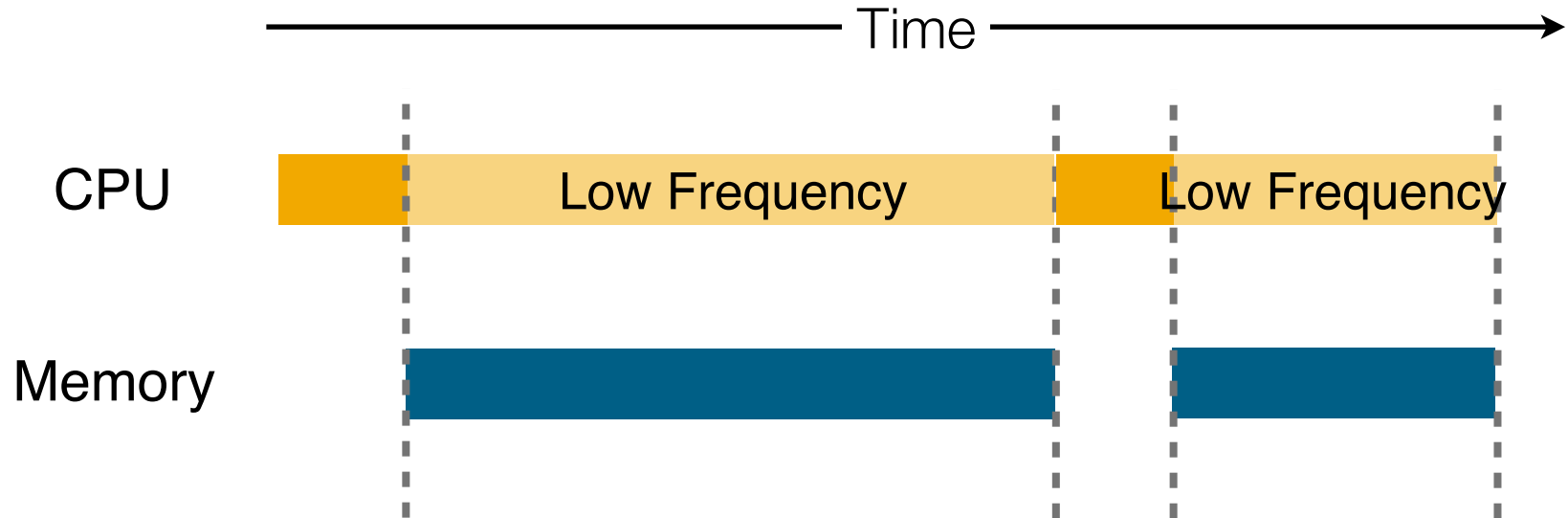
- DVFS: dynamic scales voltage and frequency depending on how much slack there is to reduce power with little performance impact
- Mostly done in OS today as the “CPU frequency governor”: based on CPU utilization, coarse-grained (e.g., ~10 ms)

# Dynamic Voltage and Frequency Scaling



- DVFS: dynamic scales voltage and frequency depending on how much slack there is to reduce power with little performance impact
- Mostly done in OS today as the “CPU frequency governor”: based on CPU utilization, coarse-grained (e.g., ~10 ms)
- You can control it. Try it yourself.

# Dynamic Voltage and Frequency Scaling



- DVFS: dynamic scales voltage and frequency depending on how much slack there is to reduce power with little performance impact
- Mostly done in OS today as the “CPU frequency governor”: based on CPU utilization, coarse-grained (e.g., ~10 ms)
- You can control it. Try it yourself.
- Hardware can do it in finer granularity, but more complex circuits

<https://www.kernel.org/doc/Documentation/cpu-freq/governors.txt>





MORGAN & CLAYPOOL PUBLISHERS

# Computer Architecture Techniques for Power Efficiency

Stefanos Kaxiras  
Margaret Martonosi

*SYNTHESIS LECTURES ON  
COMPUTER ARCHITECTURE*

Mark D. Hill, *Series Editor*

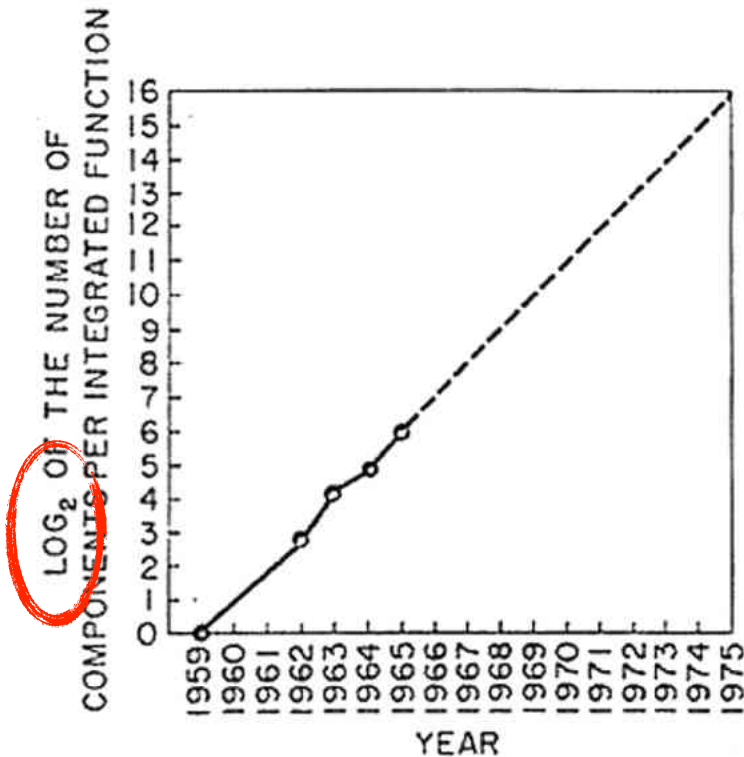
# Moore's Law

- Gordon Moore in 1965 predicted that the number of transistors doubles every year



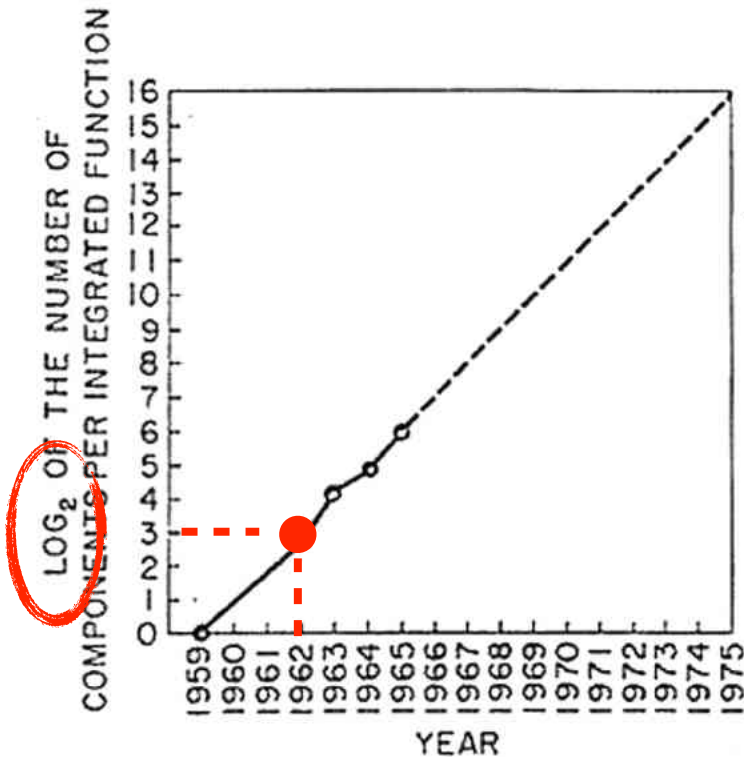
# Moore's Law

- Gordon Moore in 1965 predicted that the number of transistors doubles every year



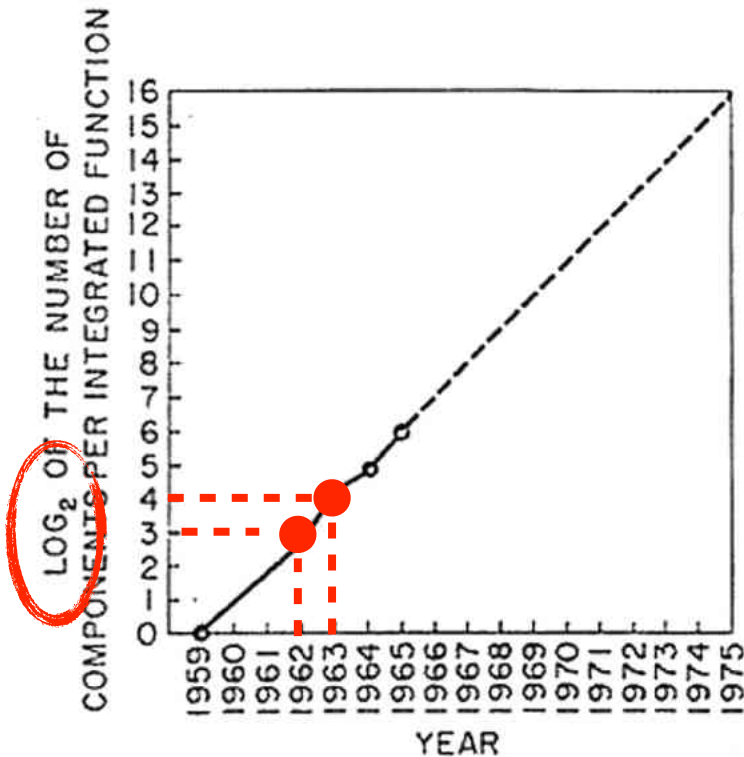
# Moore's Law

- Gordon Moore in 1965 predicted that the number of transistors doubles every year



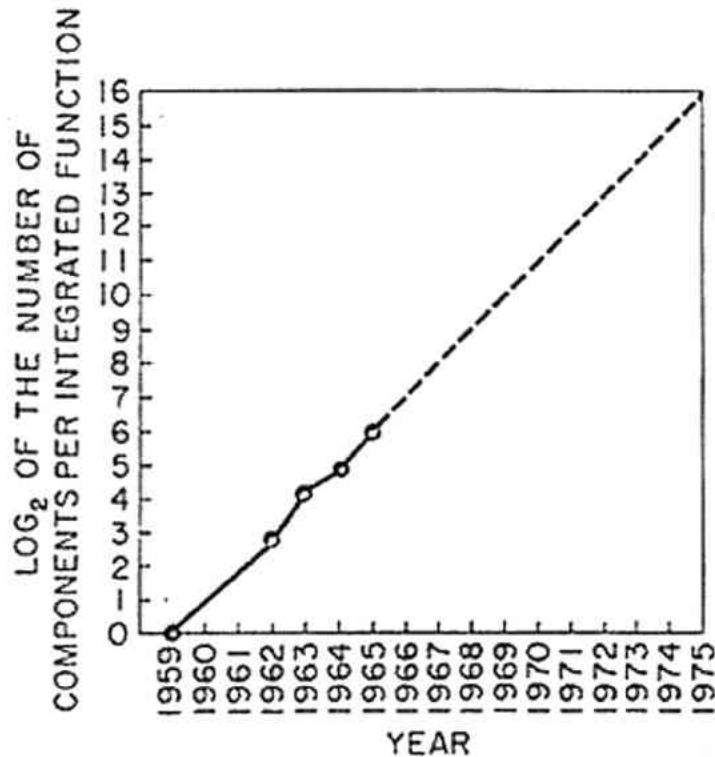
# Moore's Law

- Gordon Moore in 1965 predicted that the number of transistors doubles every year



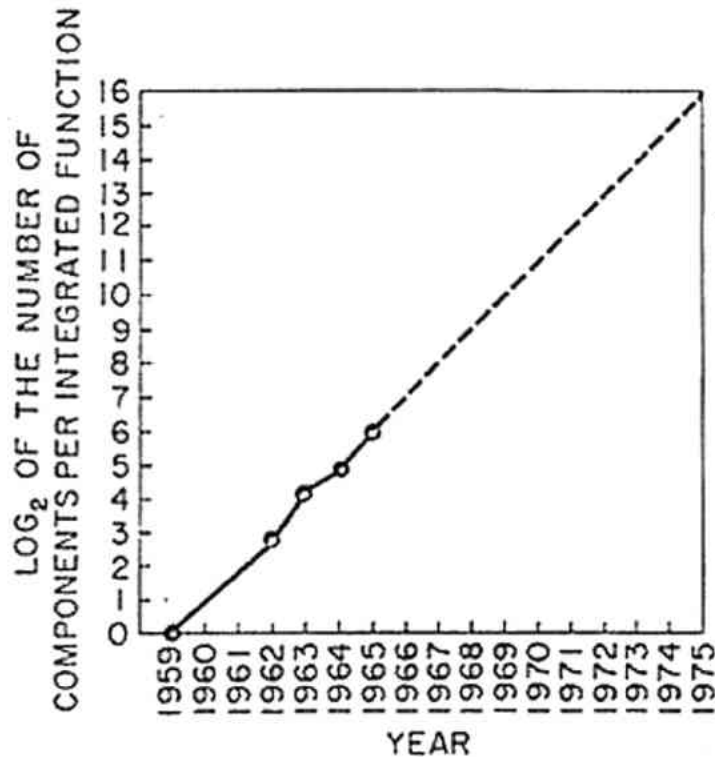
# Moore's Law

- Gordon Moore in 1965 predicted that the number of transistors doubles every year
- In 1975 he revised the prediction to doubling every 2 years



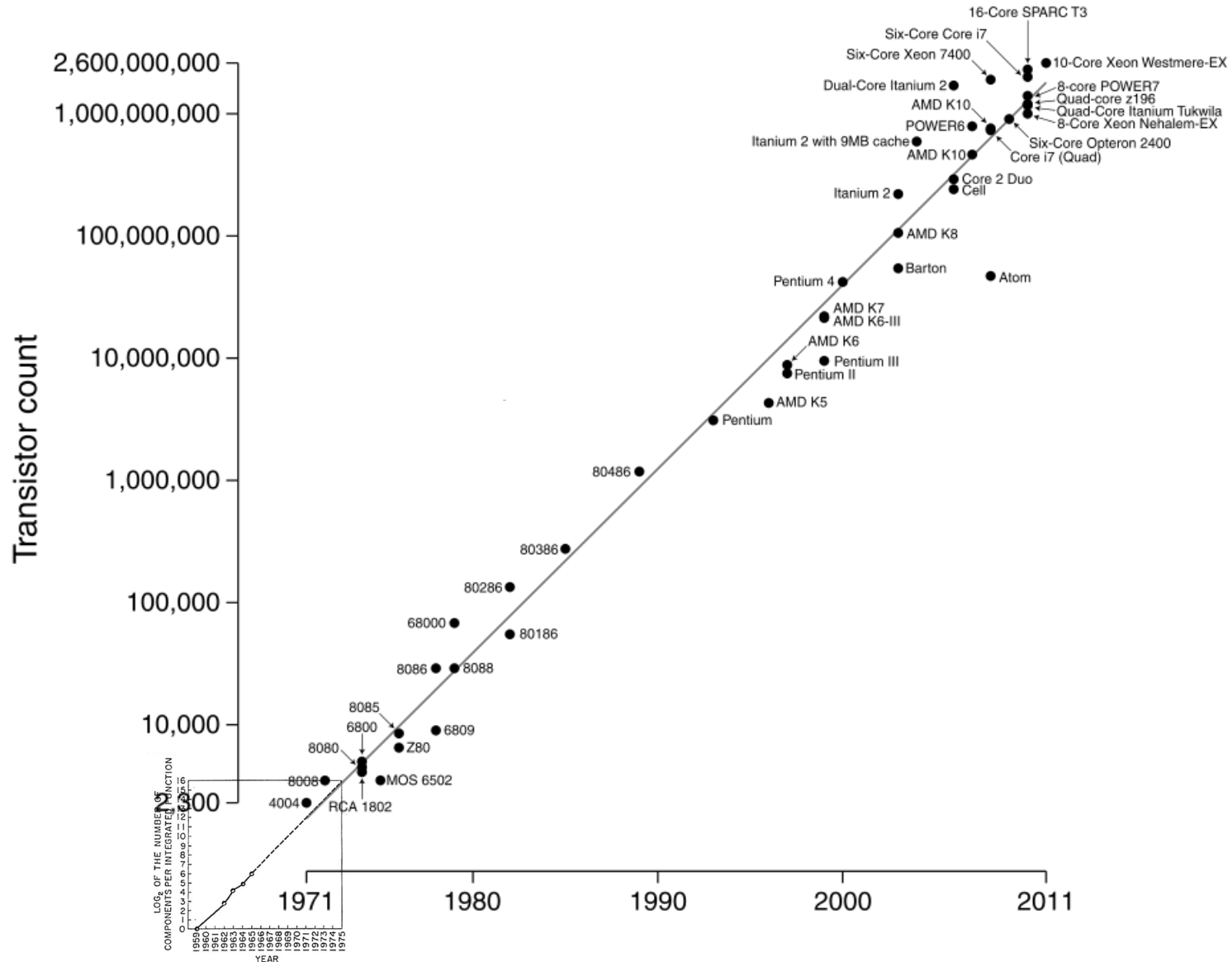
# Moore's Law

- Gordon Moore in 1965 predicted that the number of transistors doubles every year
- In 1975 he revised the prediction to doubling every 2 years
- Today's widely-known Moore's Law: number of transistors double about every 18 months (Moore never used the number 18...)



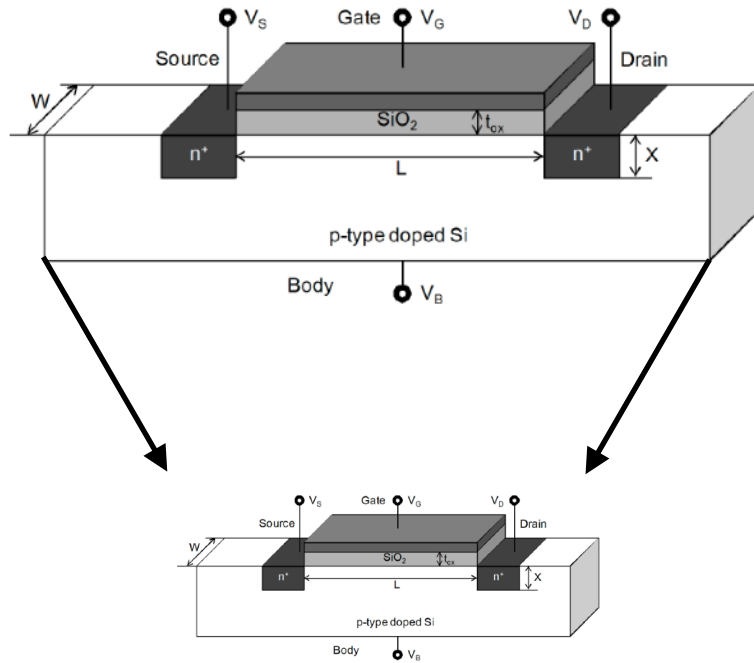


# Moore's Law





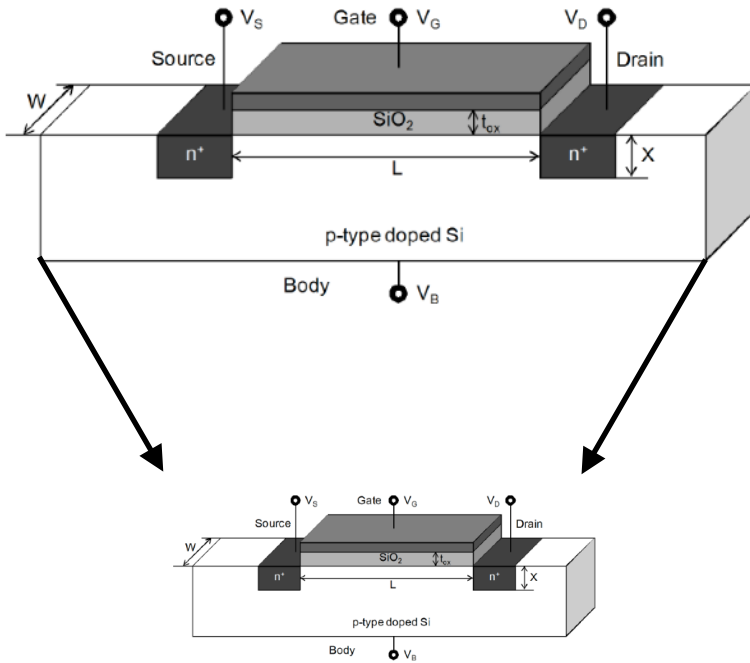
# Dennard Scaling



*Scale factor  $\alpha < 1$*

*$\alpha = 0.7 \Rightarrow 2X$  more transistors!*

# Dennard Scaling

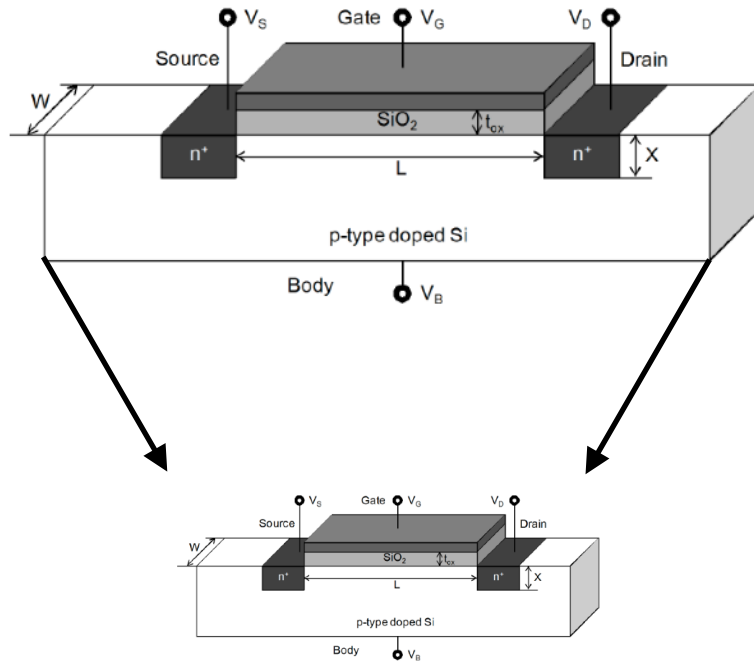


Parameter	Value	Scaled Value
Dopant concentrations	$N_a, N_d$	$N_a/\alpha, N_d/\alpha$
Dimensions	$L, W, T_{ox}$	$\alpha L, \alpha W, \alpha T_{ox}$
Field	$E$	$E$
Voltage	$V$	$\alpha V$
Capacitance	$C$	$\alpha C$
Current	$I$	$\alpha I$

*Scale factor  $\alpha < 1$*

$\alpha = 0.7 \Rightarrow 2X$  more transistors!

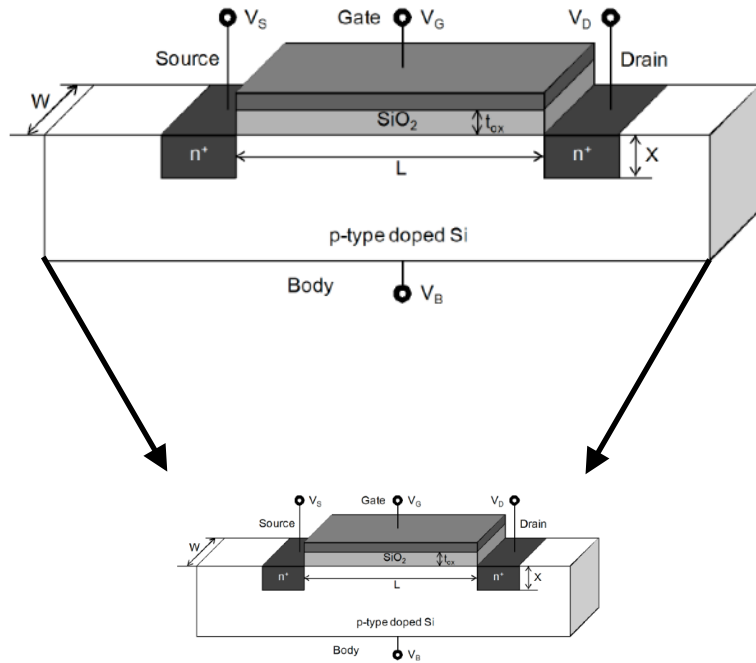
# Dennard Scaling



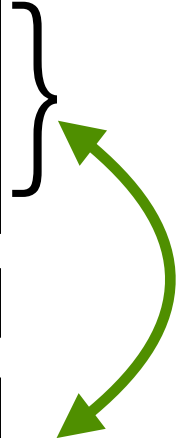
Parameter	Value	Scaled Value
Dopant concentrations	$N_a, N_d$	$N_a/\alpha, N_d/\alpha$
Dimensions	$L, W, T_{ox}$	$\alpha L, \alpha W, \alpha T_{ox}$
Field	$E$	$E$
Voltage	$V$	$\alpha V$
Capacitance	$C$	$\alpha C$
Current	$I$	$\alpha I$
Transistors/Area	$d$	$d/\alpha^2$

*Scale factor  $\alpha < 1$*   
 $\alpha = 0.7 \Rightarrow 2X$  more transistors!

# Dennard Scaling



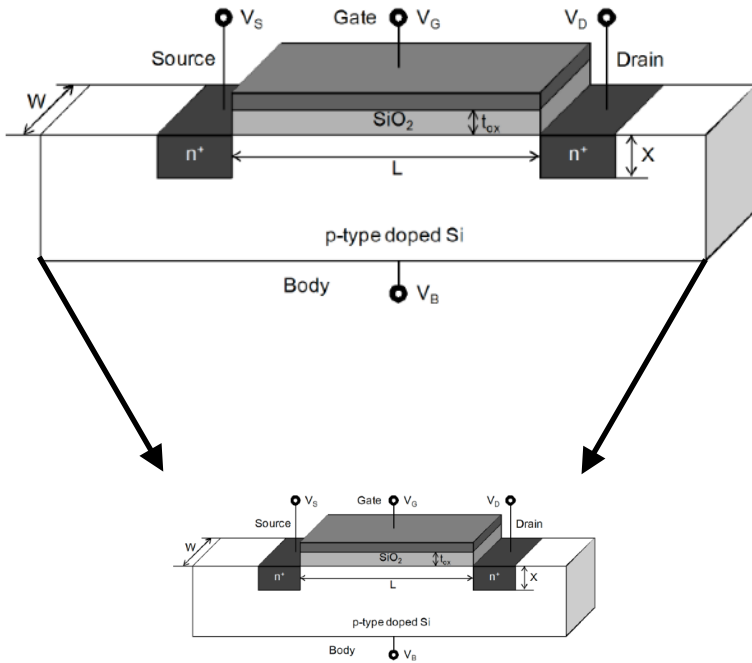
Parameter	Value	Scaled Value
Dopant concentrations	$N_a, N_d$	$N_a/\alpha, N_d/\alpha$
Dimensions	$L, W, T_{ox}$	$\alpha L, \alpha W, \alpha T_{ox}$
Field	$E$	$E$
Voltage	$V$	$\alpha V$
Capacitance	$C$	$\alpha C$
Current	$I$	$\alpha I$
Transistors/Area	$d$	$d/\alpha^2$
Propagation time ( $\sim CV/I$ )	$t$	$\alpha t$
Frequency ( $1/t$ )	$f$	$f/\alpha$



*Scale factor  $\alpha < 1$*

$\alpha = 0.7 \Rightarrow 2X$  more transistors!

# Dennard Scaling



*Scale factor  $\alpha < 1$*   
 $\alpha = 0.7 \Rightarrow 2X$  more transistors!

Parameter	Value	Scaled Value
Dopant concentrations	$N_a, N_d$	$N_a/\alpha, N_d/\alpha$
Dimensions	$L, W, T_{ox}$	$\alpha L, \alpha W, \alpha T_{ox}$
Field	$E$	$E$
Voltage	$V$	$\alpha V$
Capacitance	$C$	$\alpha C$
Current	$I$	$\alpha I$

Transistors/Area	$d$	$d/\alpha^2$
------------------	-----	--------------

Propagation time ( $\sim CV/I$ )	$t$	$\alpha t$
Frequency ( $1/t$ )	$f$	$f/\alpha$

Power ( $CV^2f$ )	$P$	$\alpha^2 P$
<b>Power/area (Power density)</b>	<b><math>P_d</math></b>	<b><math>P_d</math></b>

# Implications of Dennard Scaling

# Implications of Dennard Scaling

- Each new processor generation can have more transistors (Moore's Law), and will run at higher frequency but won't consume more power given the same area budget

# Implications of Dennard Scaling

- Each new processor generation can have more transistors (Moore's Law), and will run at higher frequency but won't consume more power given the same area budget
- More transistors means better microarchitecture, which leads to better performance even under the same frequency



# Implications of Dennard Scaling

- Each new processor generation can have more transistors (Moore's Law), and will run at higher frequency but won't consume more power given the same area budget
- More transistors means better microarchitecture, which leads to better performance even under the same frequency
- Higher frequency means better performance even under the same microarchitecture

# Implications of Dennard Scaling

- Each new processor generation can have more transistors (Moore's Law), and will run at higher frequency but won't consume more power given the same area budget
- More transistors means better microarchitecture, which leads to better performance even under the same frequency
- Higher frequency means better performance even under the same microarchitecture
- Overall, software gets a free ride: wait for the next generation of hardware and performance will naturally increase without consuming more power

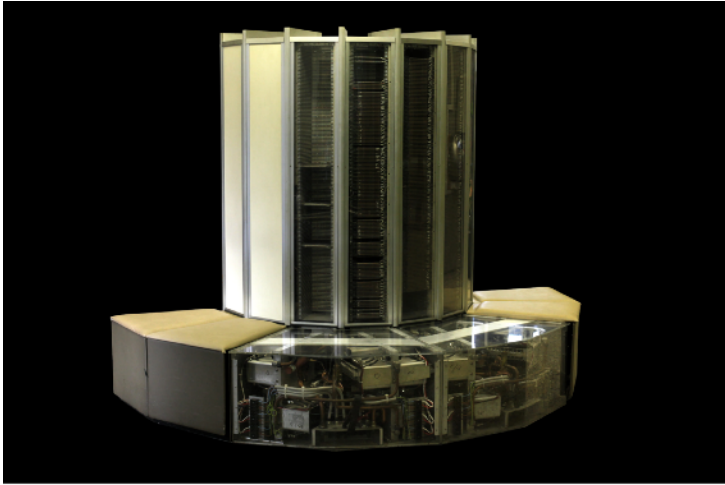
# Implications of Dennard Scaling

- Each new processor generation can have more transistors (Moore's Law), and will run at higher frequency but won't consume more power given the same area budget
- More transistors means better microarchitecture, which leads to better performance even under the same frequency
- Higher frequency means better performance even under the same microarchitecture
- Overall, software gets a free ride: wait for the next generation of hardware and performance will naturally increase without consuming more power

Moore's law gave us more transistors;  
Dennard scaling made them useful.

Bob Colwell, DAC 2013, June 4, 2013

# Moore's Law + Dennard Scaling Have Meant



- 1976 Cray 1
  - Frequency: 80 MHz
  - Speed: 250 M Ops/second
  - Scale: 2.5 Million transistors
  - 5,000 kg, 115 KW
  - Price: \$9M
  - 80 manufactured



- 2014 iPhone 6
  - Frequency: 1.4 GHz
  - Speed: > 4 B Ops/second
  - Scale: > 3 Billion transistors
  - 120 g, < 5 W
  - Price: \$649
  - 10 million sold in first 3 days

# 2005: End of Dennard Scaling

- What Happened?
  - Supply voltage  $V_{dd}$  stops scaling (Can't drop voltage below  $\sim 1$  V)
  - Remember Power =  $CV^2f$

# 2005: End of Dennard Scaling

- What Happened?

- Supply voltage  $V_{dd}$  stops scaling (Can't drop voltage below  $\sim 1$  V)
- Remember Power =  $CV^2f$

- Why?

- There is a fundamental limit as to how much voltage we need to switch a transistor, called threshold voltage ( $V_{th}$ ).
- $V_{th}$  stopped scaling because leakage power/reliability/variation becomes huge issues, and accordingly  $V_{dd}$  stops scaling

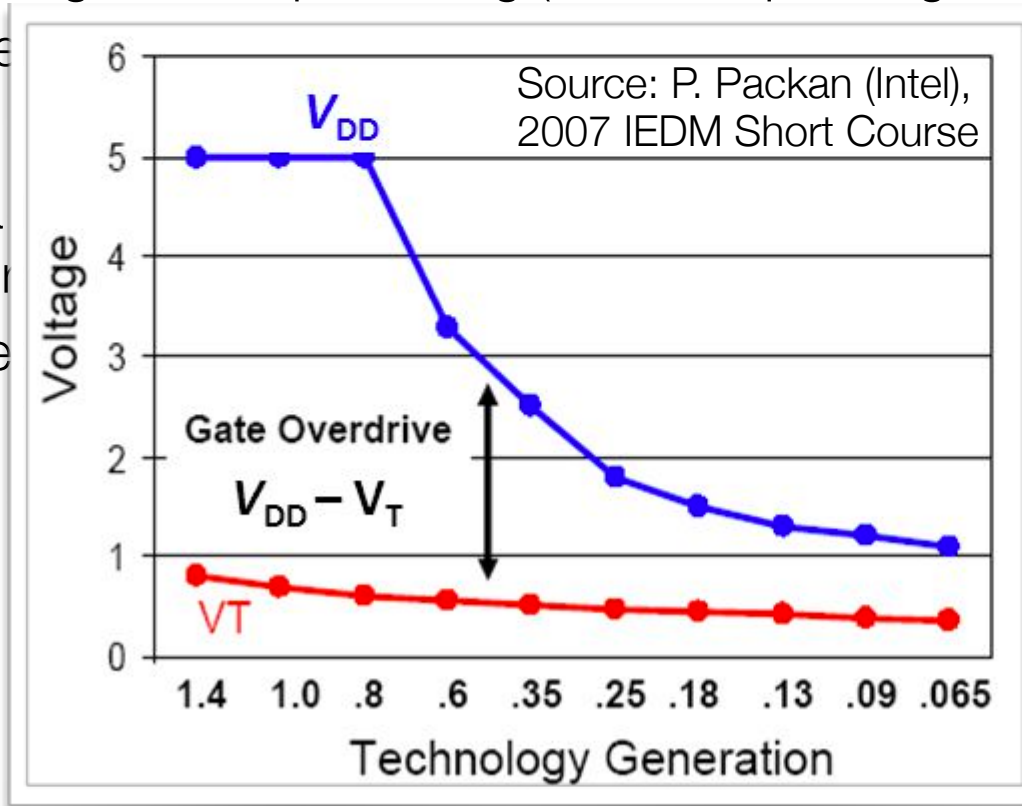
# 2005: End of Dennard Scaling

- What Happened?

- Supply voltage  $V_{dd}$  stops scaling (Can't drop voltage below  $\sim 1$  V)
- Remember

- Why?

- There is a switch a tr
- $V_{th}$  stoppe
- becomes



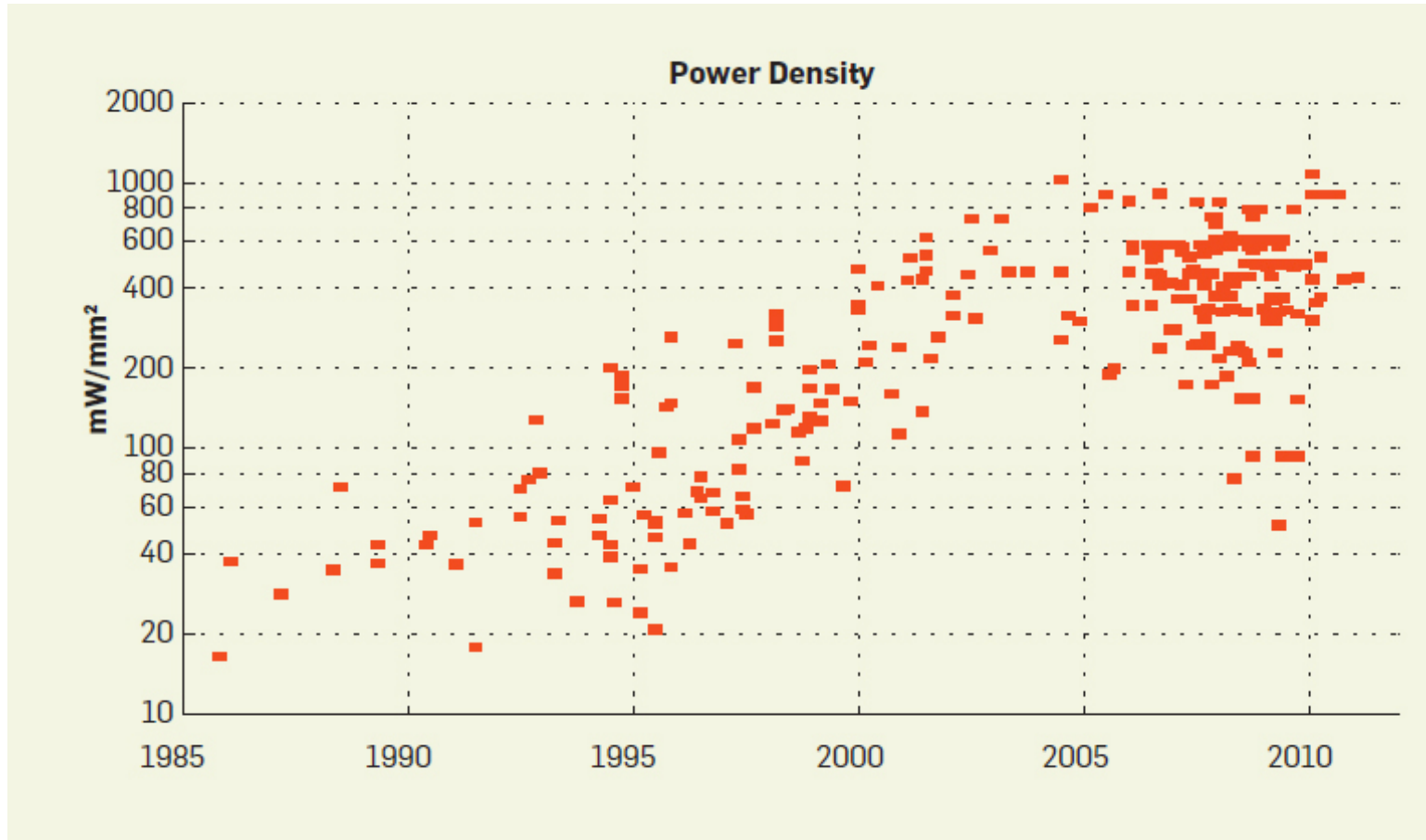
the need to  
variation  
g

# 2005: End of Dennard Scaling

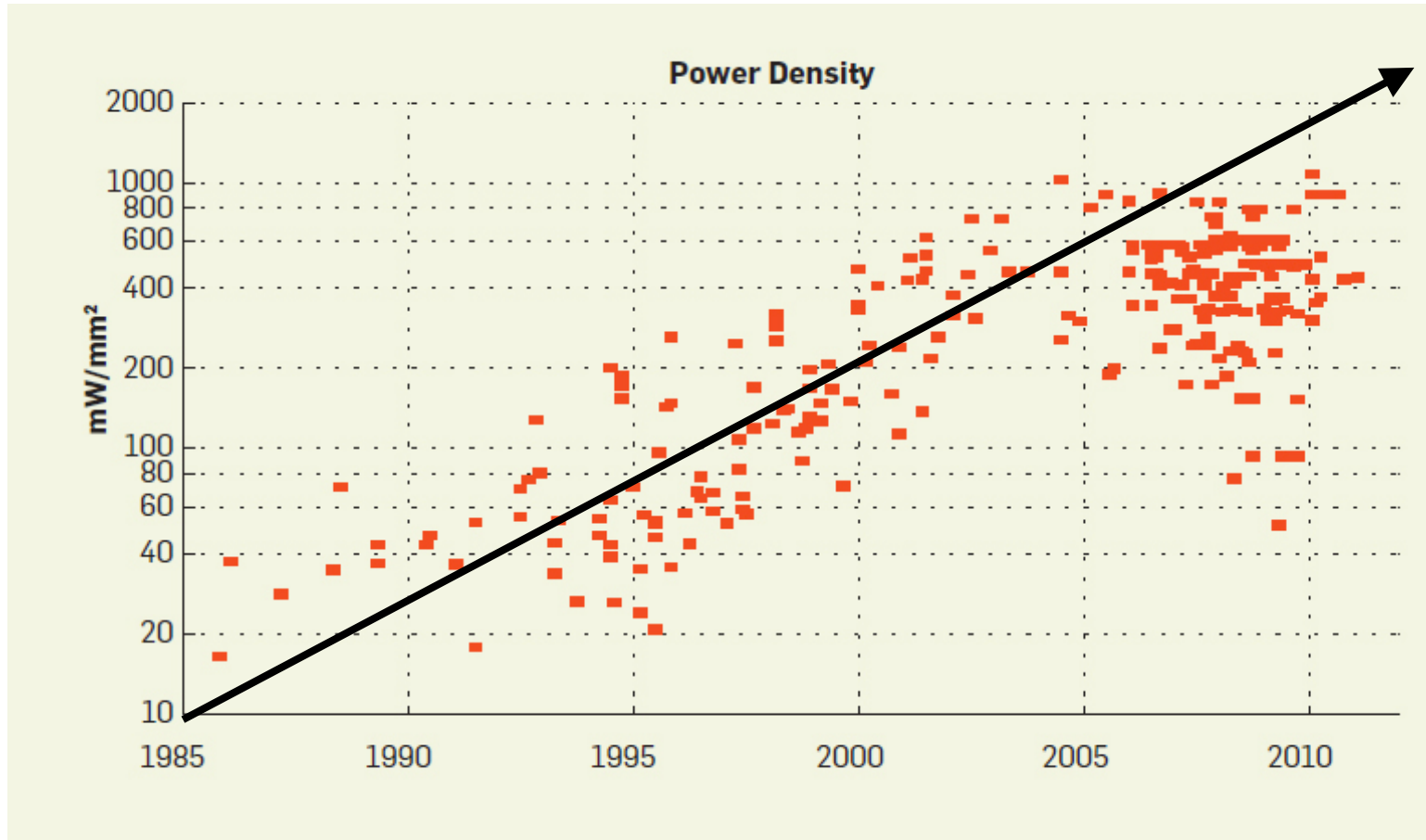
- What Happened?
  - Supply voltage  $V_{dd}$  stops scaling (Can't drop voltage below  $\sim 1$  V)
  - Remember  $\text{Power} = CV^2f$
- Why?
  - There is a fundamental limit as to how much voltage we need to switch a transistor, called threshold voltage ( $V_{th}$ ).
  - $V_{th}$  stopped scaling because leakage power/reliability/variation becomes huge issues, and accordingly  $V_{dd}$  stops scaling
- The demise of Dennard Scaling means the power density (power consumption per unit area) will increase rather than staying stable.



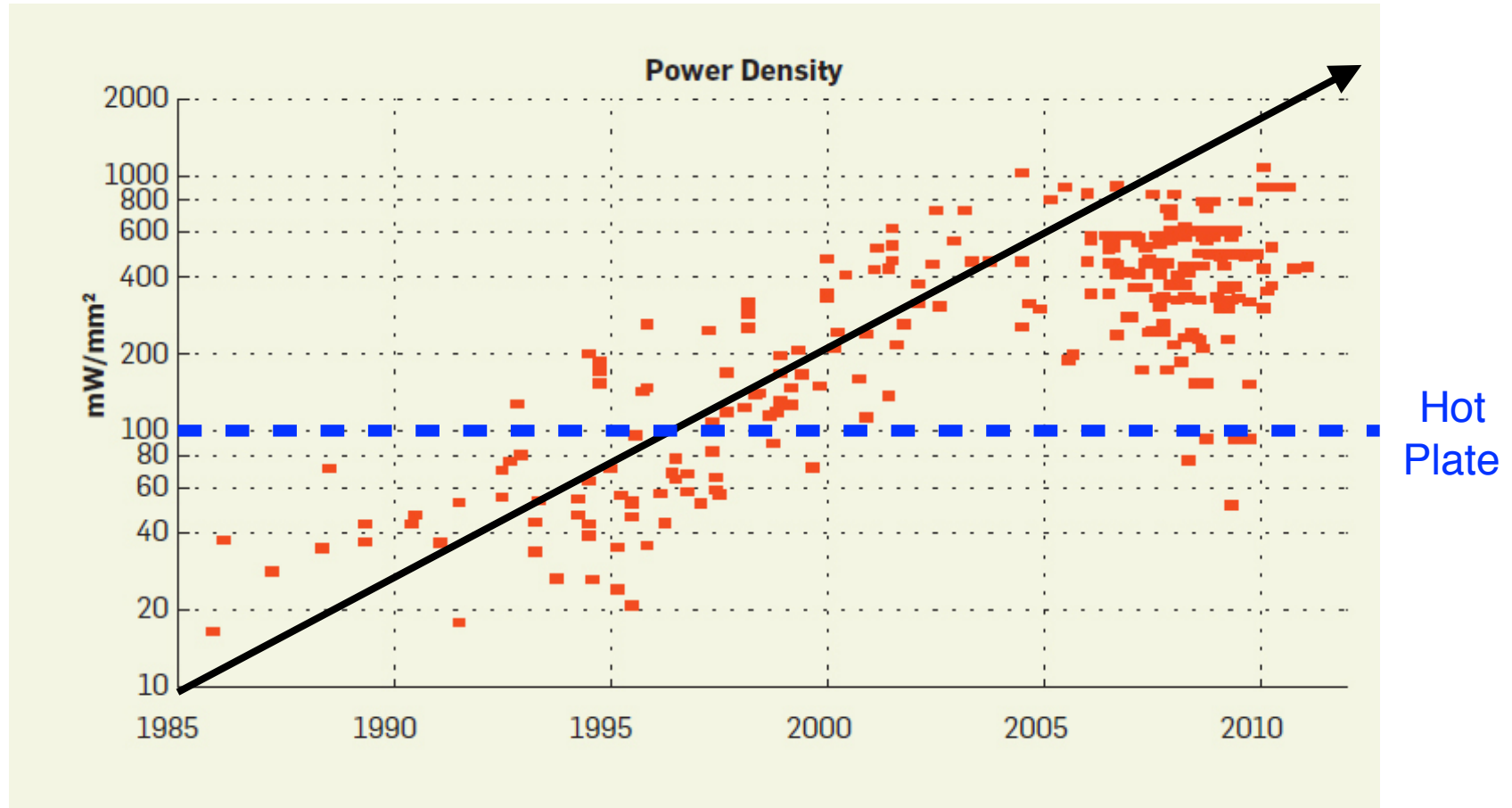
# 2005: End of Dennard Scaling



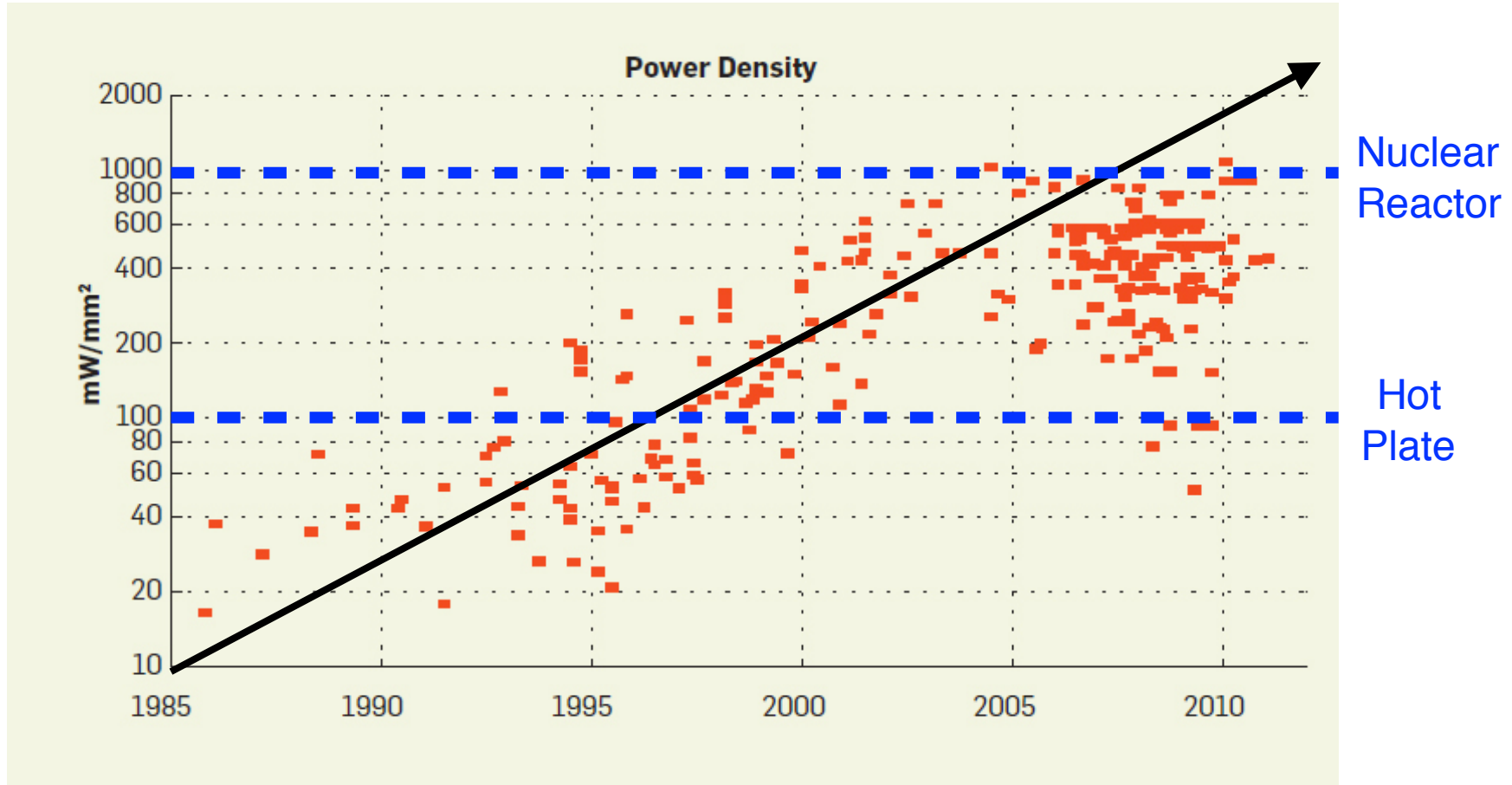
# 2005: End of Dennard Scaling



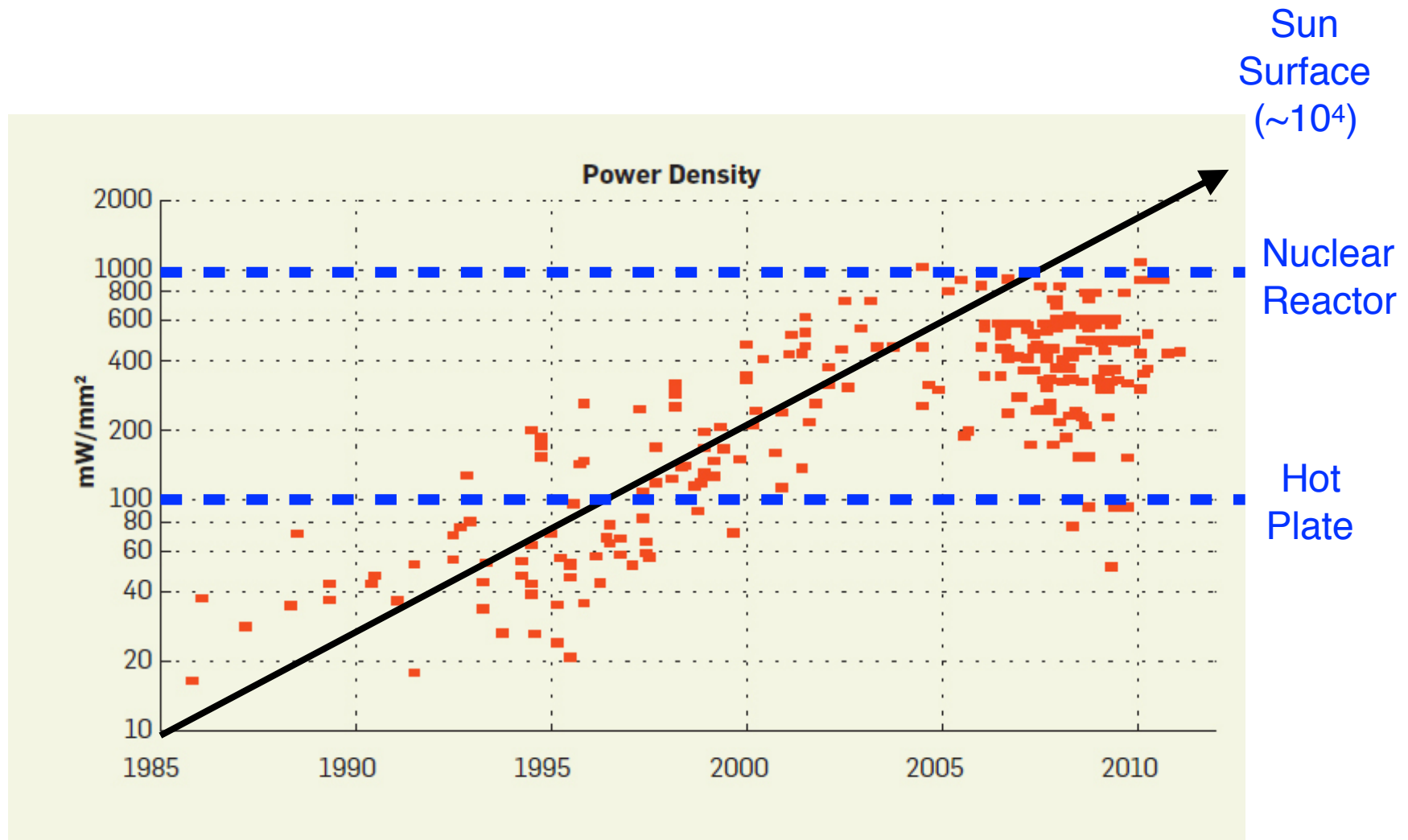
# 2005: End of Dennard Scaling



# 2005: End of Dennard Scaling



# 2005: End of Dennard Scaling



# Dark Silicon

*n.* [därk, sɪl'ɪ-kən, -kɒn']

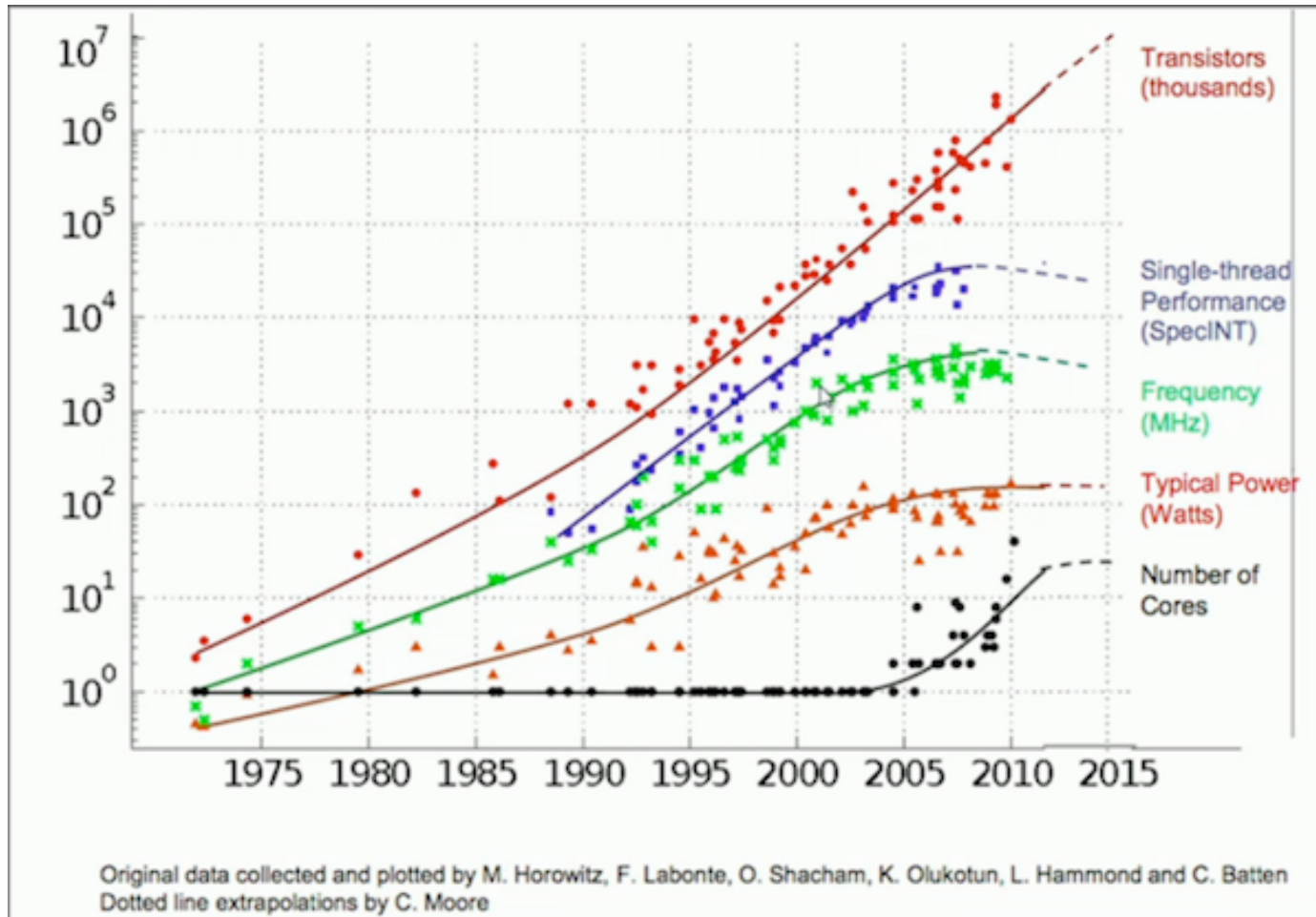
More transistors on chip (Moore's Law), but a growing fraction cannot actually be used due to power limits.

# 2005: End of Dennard Scaling

- Initial response has been to lower frequency and increase cores / chip

# 2005: End of Dennard Scaling

- Initial response has been to lower frequency and increase cores / chip





# Dark Silicon and the End of Multicore Scaling

Hadi Esmaeilzadeh<sup>†</sup> Emily Blem<sup>‡</sup> Renée St. Amant<sup>§</sup> Karthikeyan Sankaralingam<sup>¶</sup> Doug Burger<sup>¶</sup>

<sup>†</sup>University of Washington <sup>‡</sup>University of Wisconsin-Madison

<sup>§</sup>The University of Texas at Austin <sup>¶</sup>Microsoft Research

hadianeh@cs.washington.edu blem@cs.wisc.edu stamant@cs.utexas.edu karu@cs.wisc.edu dburger@microsoft.com

## ABSTRACT

Since 2005, processor designers have increased core counts to exploit Moore's Law scaling, rather than focusing on single-core performance. The failure of Dennard scaling, to which the shift to multicore parts is partially a response, may soon limit multicore scaling just as single-core scaling has been curtailed. This paper models multicore scaling limits by combining device scaling, single-core scaling, and multicore scaling to measure the speedup potential for a set of parallel workloads for the next five technology generations. For device scaling, we use both the ITRS projections and a set of more conservative device scaling parameters. To model single-core scaling, we combine measurements from over 150 processors to derive Pareto-optimal frontiers for area/performance and power/performance. Finally, to model multicore scaling, we build a detailed performance model of upper-bound performance and lower-bound core power. The multicore designs we study include single-threaded CPU-like and massively threaded GPU-like multicore chip organizations with symmetric, asymmetric, dynamic, and composed topologies. The study shows that regardless of chip organization and topology, multicore scaling is power limited to a degree not widely appreciated by the computing community. Even at 22 nm (just one year from now), 21% of a fixed-size chip must be powered off, and at 8 nm, this number grows to more than 50%. Through 2024, only 7.9× average speedup is possible across commonly used parallel workloads, leaving a nearly 24-fold gap from a target of doubled performance per generation.

**Categories and Subject Descriptors:** C.0 [Computer Systems Organization] General — Modeling of computer architecture; C.0 [Computer Systems Organization] General — System architectures

**General Terms:** Design, Measurement, Performance

**Keywords:** Dark Silicon, Modeling, Power, Technology Scaling, Multicore

## 1. INTRODUCTION

Moore's Law [24] (the doubling of *transistors* on chip every 18 months) has been a fundamental driver of computing. For the past three decades, through device, circuit, microarchitecture, architec-

ture, and compiler advances, Moore's Law, coupled with Dennard scaling [11], has resulted in commensurate exponential performance increases. The recent shift to multicore designs has aimed to increase the number of cores along with transistor count increases, and continue the proportional scaling of performance. As a result, architecture researchers have started focusing on 100-core and 1000-core chips and related research topics and called for changes to the undergraduate curriculum to solve the parallel programming challenge for multicore designs at these scales.

With the failure of Dennard scaling—and thus slowed supply voltage scaling—core count scaling may be in jeopardy, which would leave the community with no clear scaling path to exploit continued transistor count increases. Since future designs will be power limited, higher core counts must provide performance gains despite the worsening energy and speed scaling of transistors, and given the available parallelism in applications. By studying these characteristics together, it is possible to predict for how many additional technology generations multicore scaling will provide a clear benefit. Since the energy efficiency of devices is not scaling along with integration capacity, and since few applications (even from emerging domains such as recognition, mining, and synthesis [5]) have parallelism levels that can efficiently use a 100-core or 1000-core chip, it is critical to understand how good multicore performance will be in the long term. In 2024, will processors have 32 times the performance of processors from 2008, exploiting five generations of core doubling?

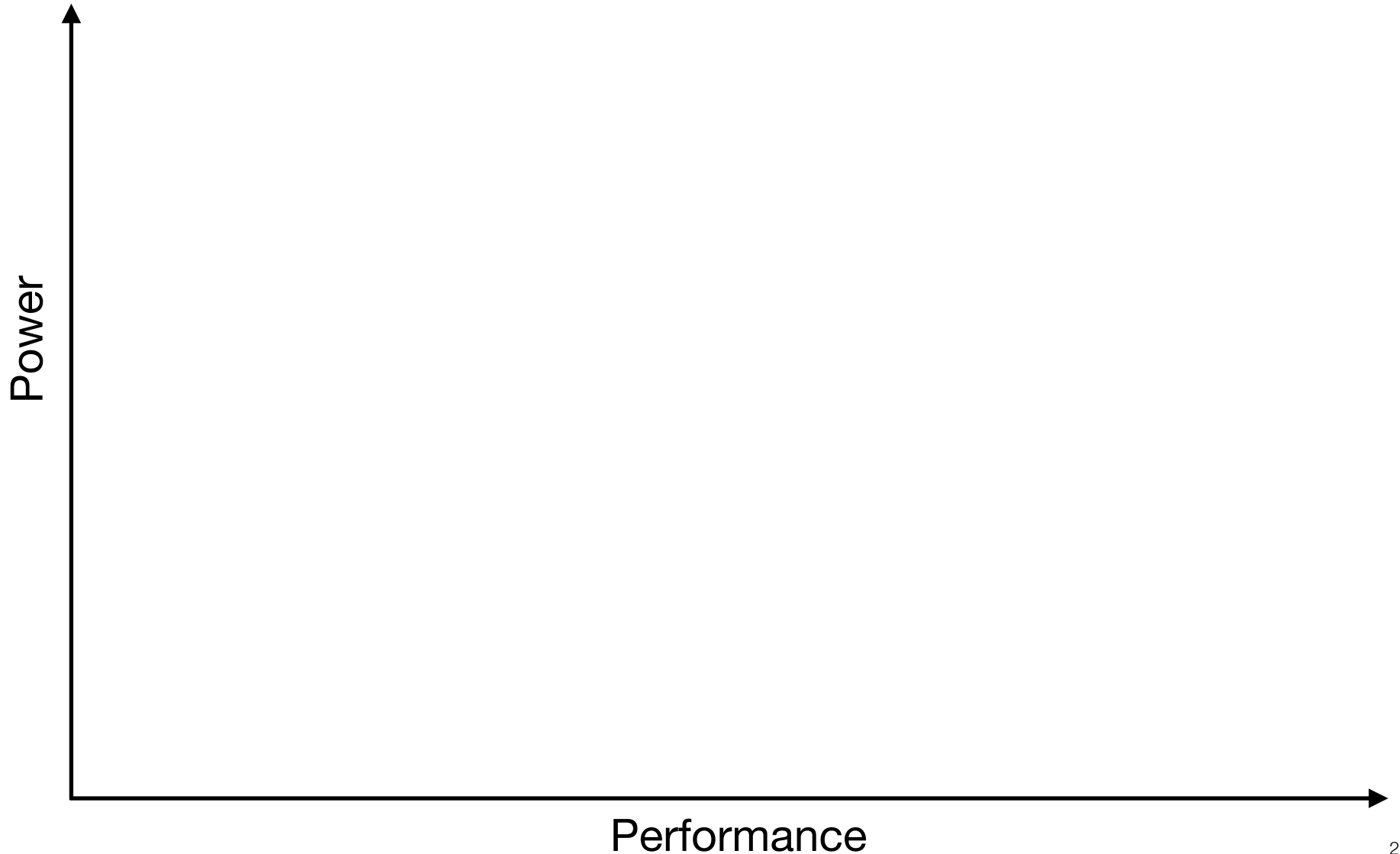
Such a study must consider devices, core microarchitectures, chip organizations, and benchmark characteristics, applying area and power limits at each technology node. This paper considers all those factors together, projecting upper-bound performance achievable through multicore scaling, and measuring the effects of *non-ideal* device scaling, including the percentage of “dark silicon” (transistor under-utilization) on future multicore chips. Additional projections include best core organization, best chip-level topology, and optimal number of cores.

We consider technology scaling projections, single-core design scaling, multicore design choices, actual application behavior, and microarchitectural features together. Previous studies have also analyzed these features in various combinations, but not all together [8, 9, 10, 15, 16, 21, 23, 28, 29]. This study builds and

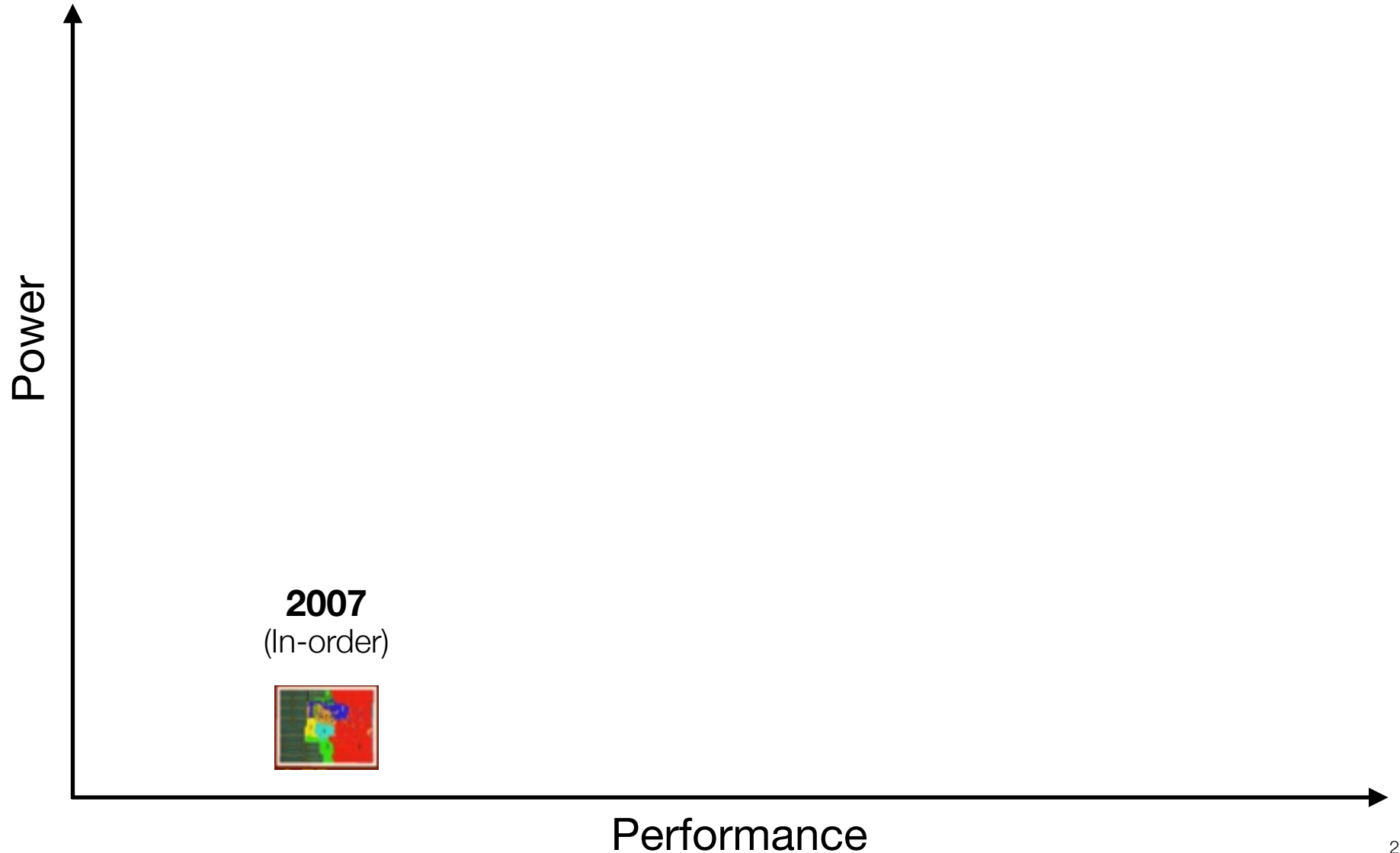
# 2007: A Revolutionary New Computer



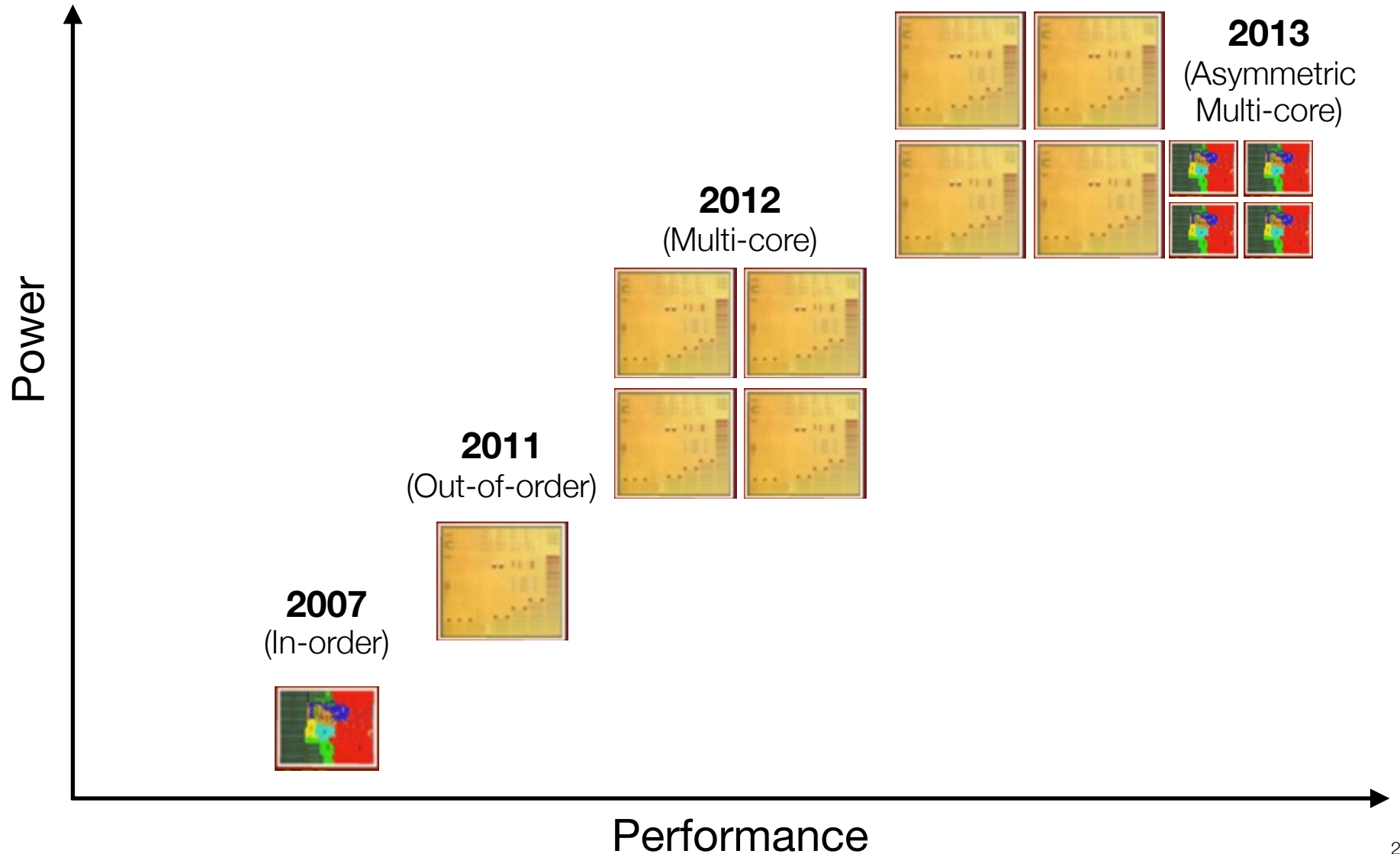
# Mobile Processor Design “Strategy”



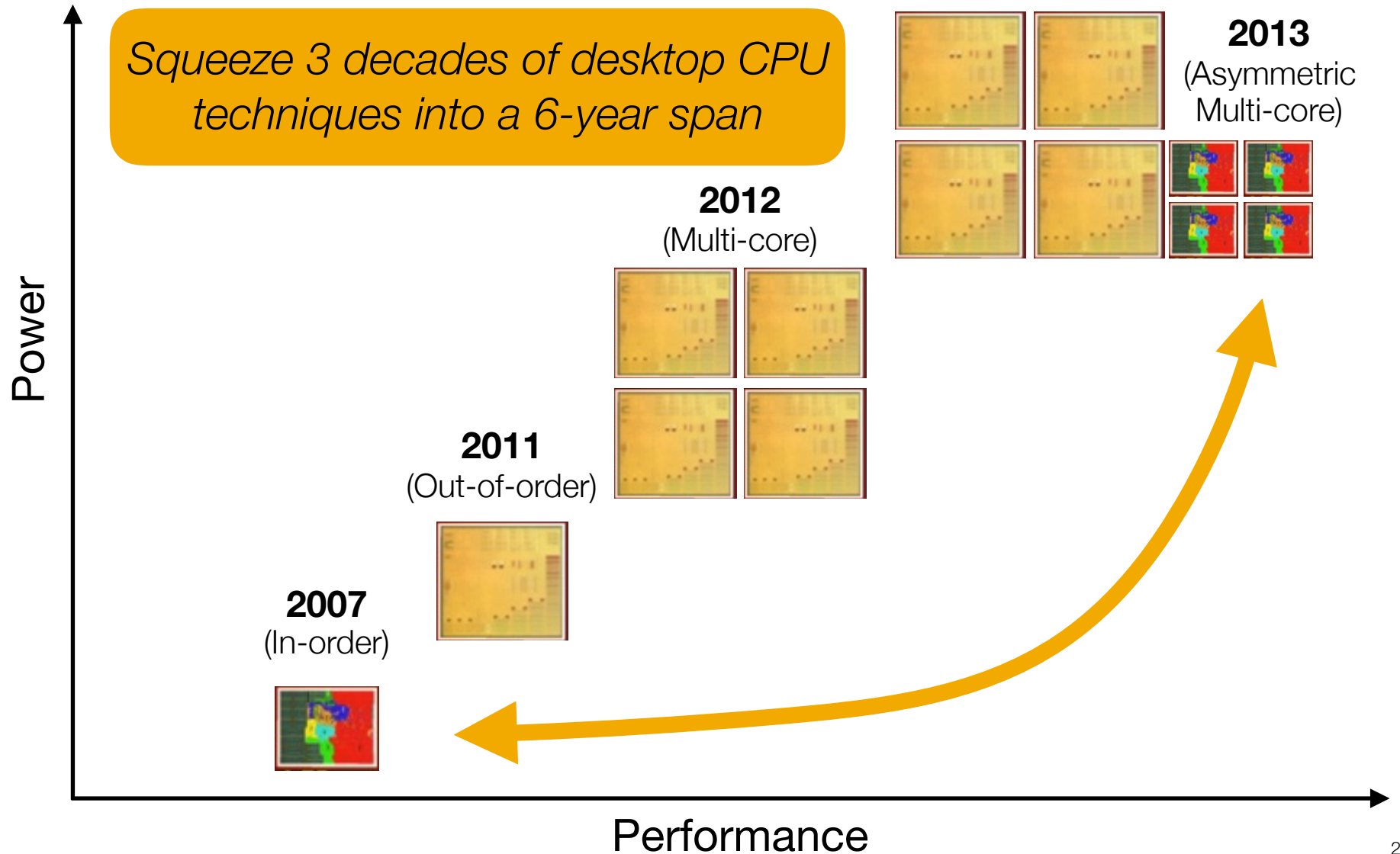
# Mobile Processor Design “Strategy”



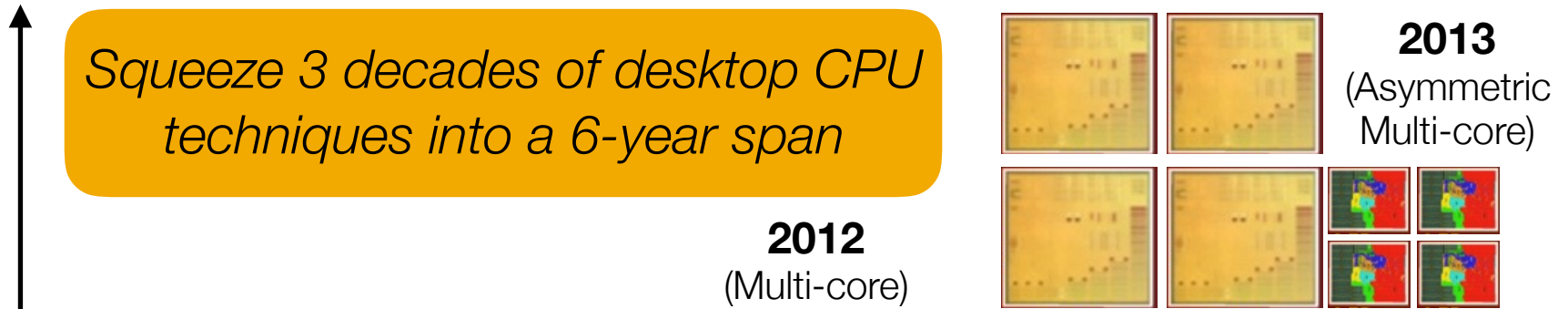
# Mobile Processor Design “Strategy”



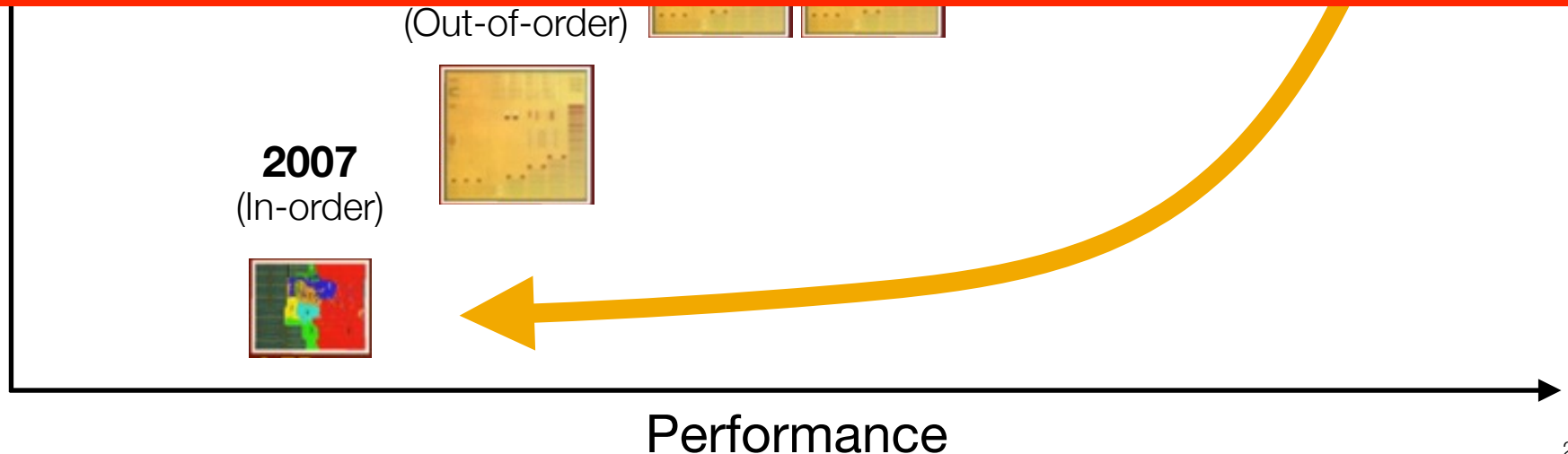
# Mobile Processor Design “Strategy”



# Mobile Processor Design “Strategy”



Ran Into Dark Silicon Issue in 6 Years





# No Moore's Law for batteries

**Fred Schlachter<sup>1</sup>**

*American Physical Society, Washington, DC 20045*

The public has become accustomed to rapid progress in mobile phone technology, computers, and access to information; tablet computers, smart phones, and other powerful new devices are familiar to most people on the planet.

These developments are due in part to the ongoing exponential increase in computer processing power, doubling approximately every 2 years for the past several decades. This pattern is usually called Moore's Law and is named for Gordon Moore, a co-founder of Intel. The law is not a law like that for gravity; it is an empirical observation, which has become a self-fulfilling prophecy. Unfortunately, much of the public has come to expect that all technology does, will, or should follow such a law, which is not consistent with our everyday observations: For example, the maximum

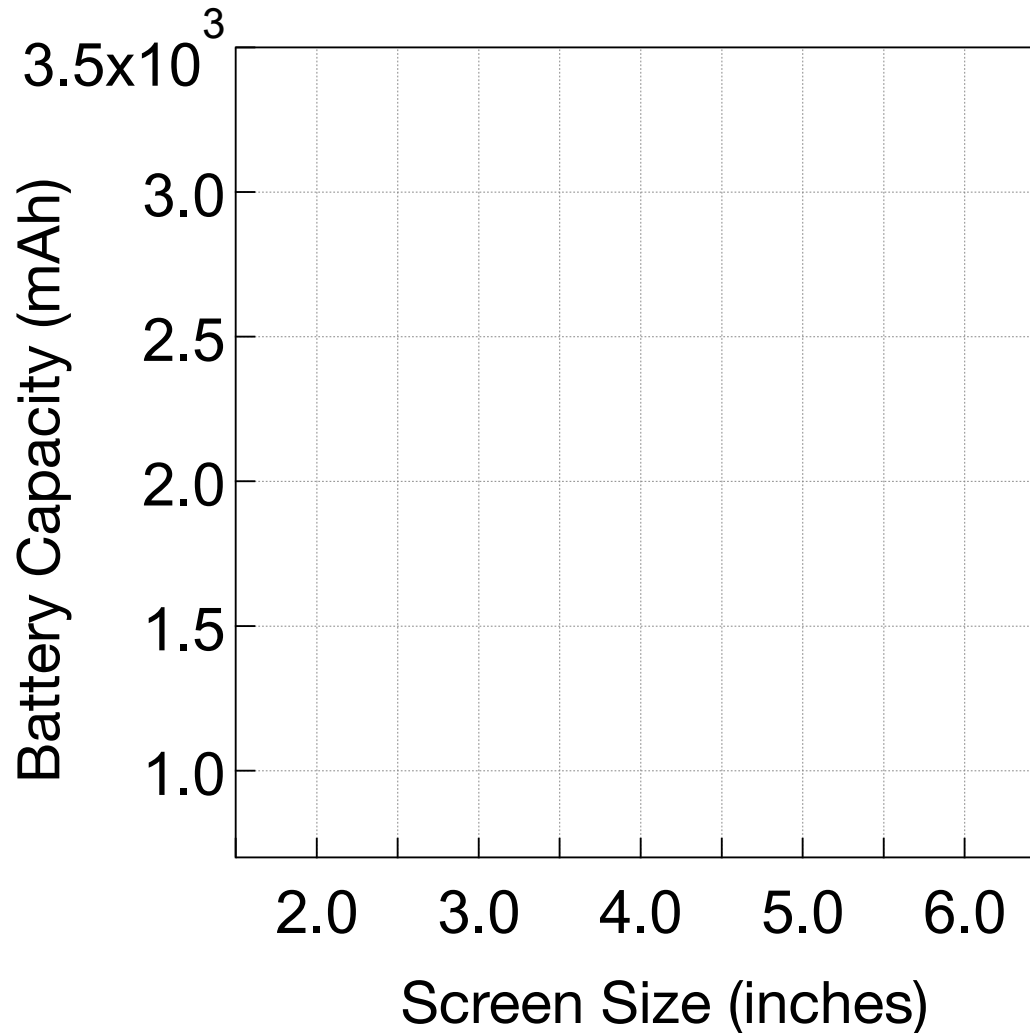
there is a Moore's Law for computer processors is that electrons are small and they do not take up space on a chip. Chip performance is limited by the lithography technology used to fabricate the chips; as lithography improves ever smaller features can be made on processors. Batteries are not like this. Ions, which transfer charge in batteries, are large, and they take up space, as do anodes, cathodes, and electrolytes. A D-cell battery stores more energy than an AA-cell. Potentials in a battery are dictated by the relevant chemical reactions, thus limiting eventual battery performance. Significant improvement in battery capacity can only be made by changing to a different chemistry.

Scientists and battery experts, who have been optimistic in the recent past about improving lithium-ion batteries and about de-



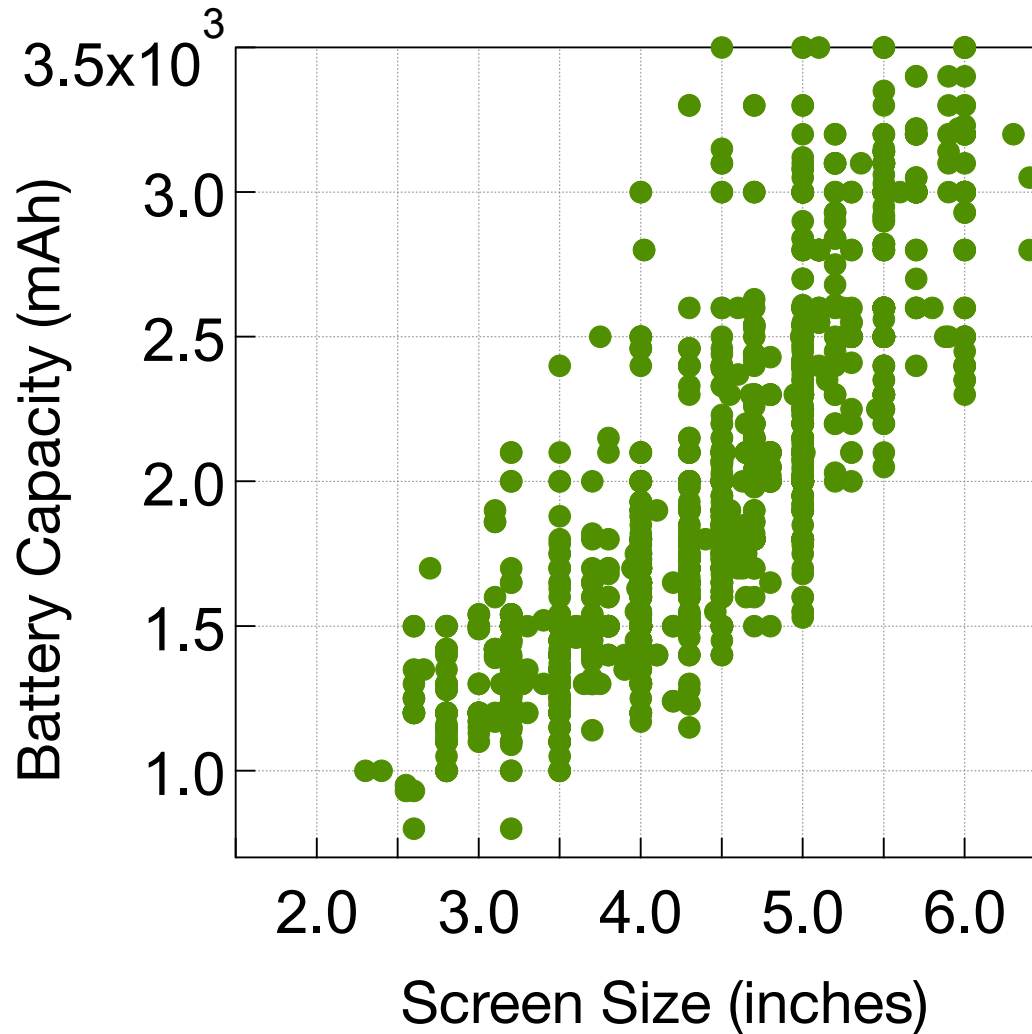
# “Improving” Energy Capacity

# “Improving” Energy Capacity



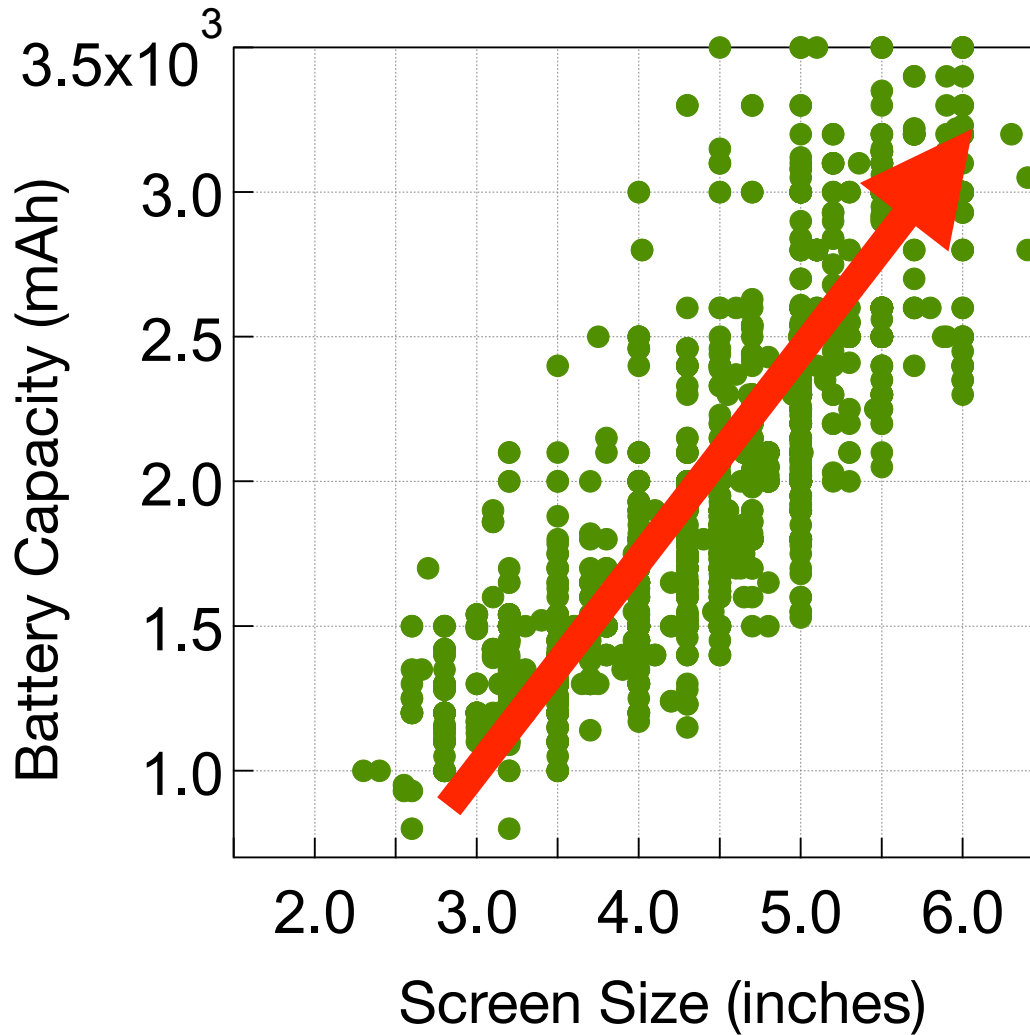
600 smartphone from 2006 to 2014 on <http://www.gsmarena.com/makers.php3>

# “Improving” Energy Capacity



600 smartphone from 2006 to 2014 on <http://www.gsmarena.com/makers.php3>

# “Improving” Energy Capacity



600 smartphone from 2006 to 2014 on <http://www.gsmarena.com/makers.php3>

# **“Improving” Energy Capacity**

# “Improving” Energy Capacity



# “Improving” Energy Capacity



SMARTPHONE



PHABLET



TABLET

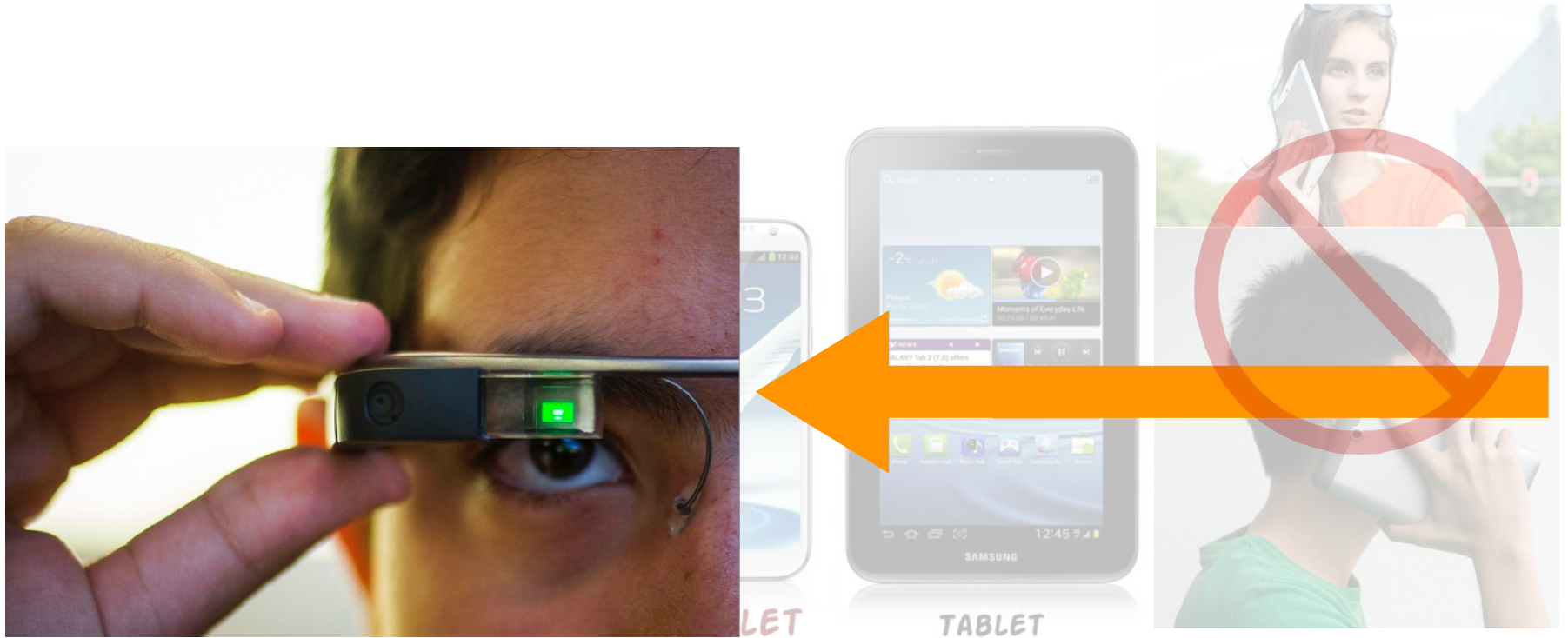


# “Improving” Energy Capacity





# “Improving” Energy Capacity



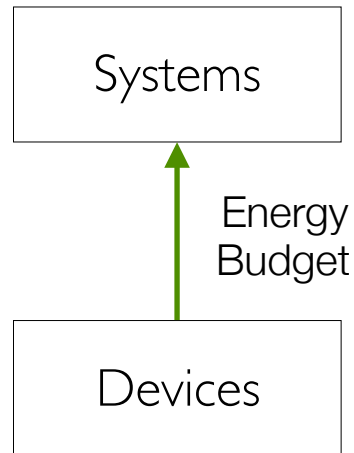
# The (Mobile) Computing Virtuous Cycle

# The (Mobile) Computing Virtuous Cycle

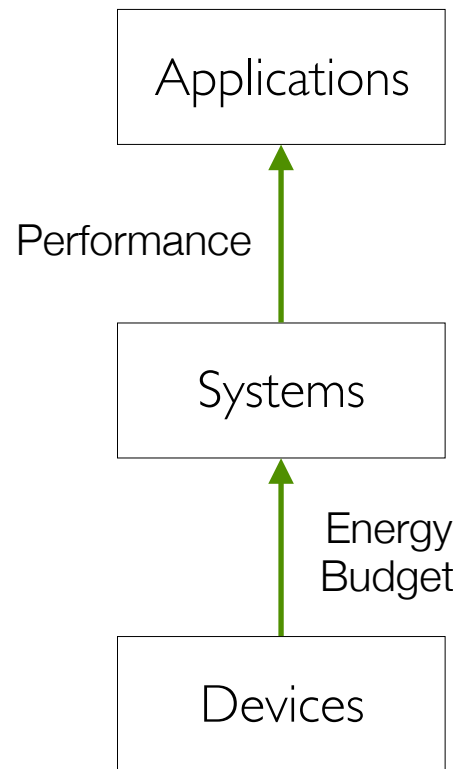


Devices

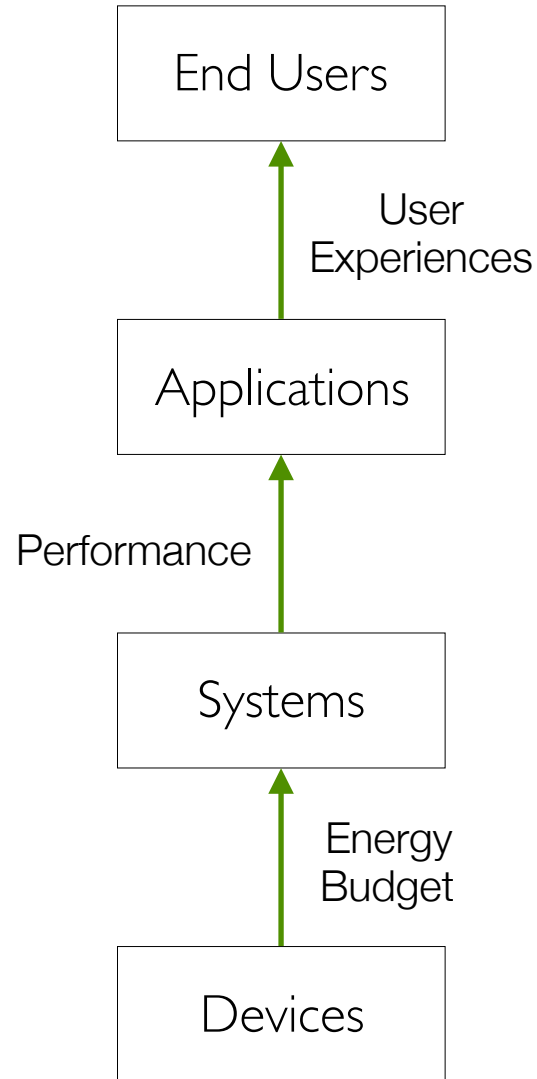
# The (Mobile) Computing Virtuous Cycle



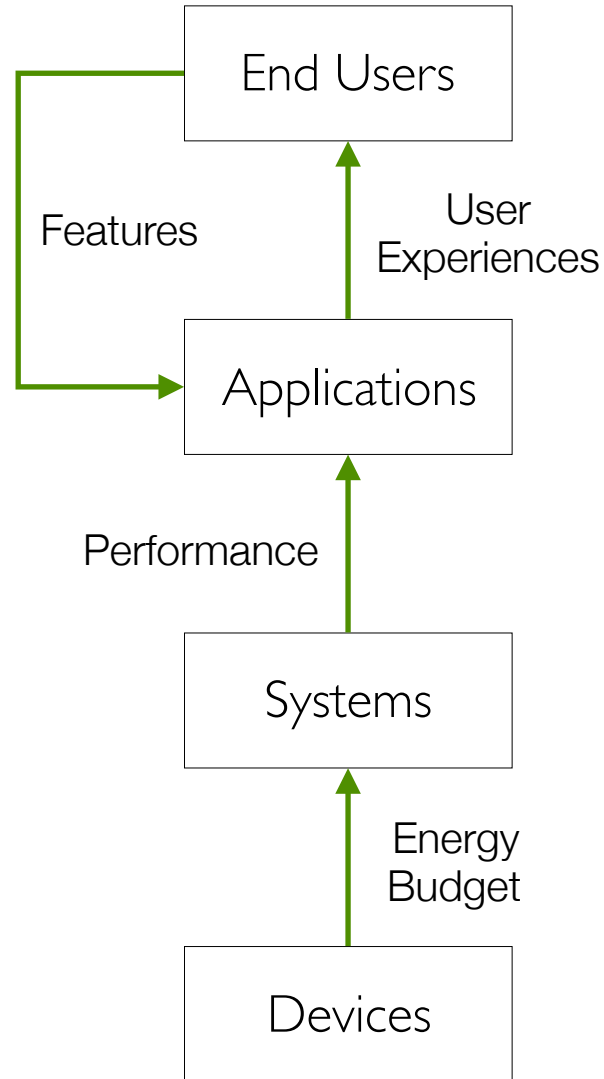
# The (Mobile) Computing Virtuous Cycle



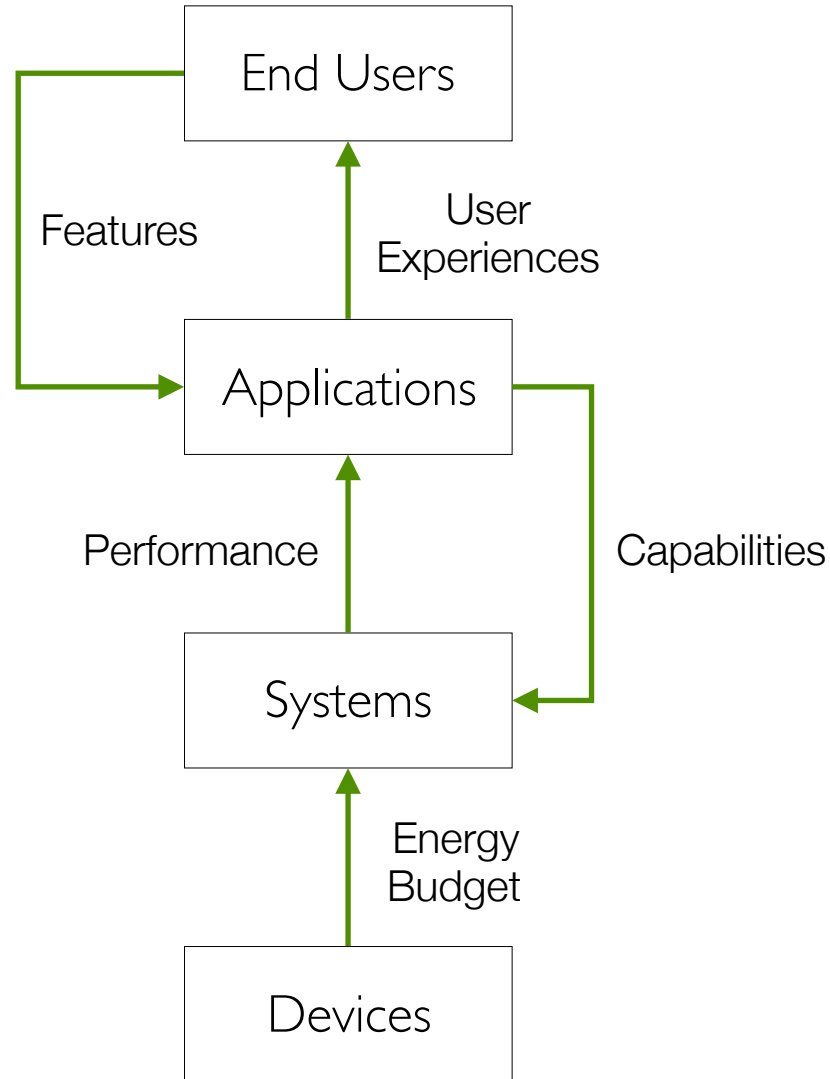
# The (Mobile) Computing Virtuous Cycle



# The (Mobile) Computing Virtuous Cycle

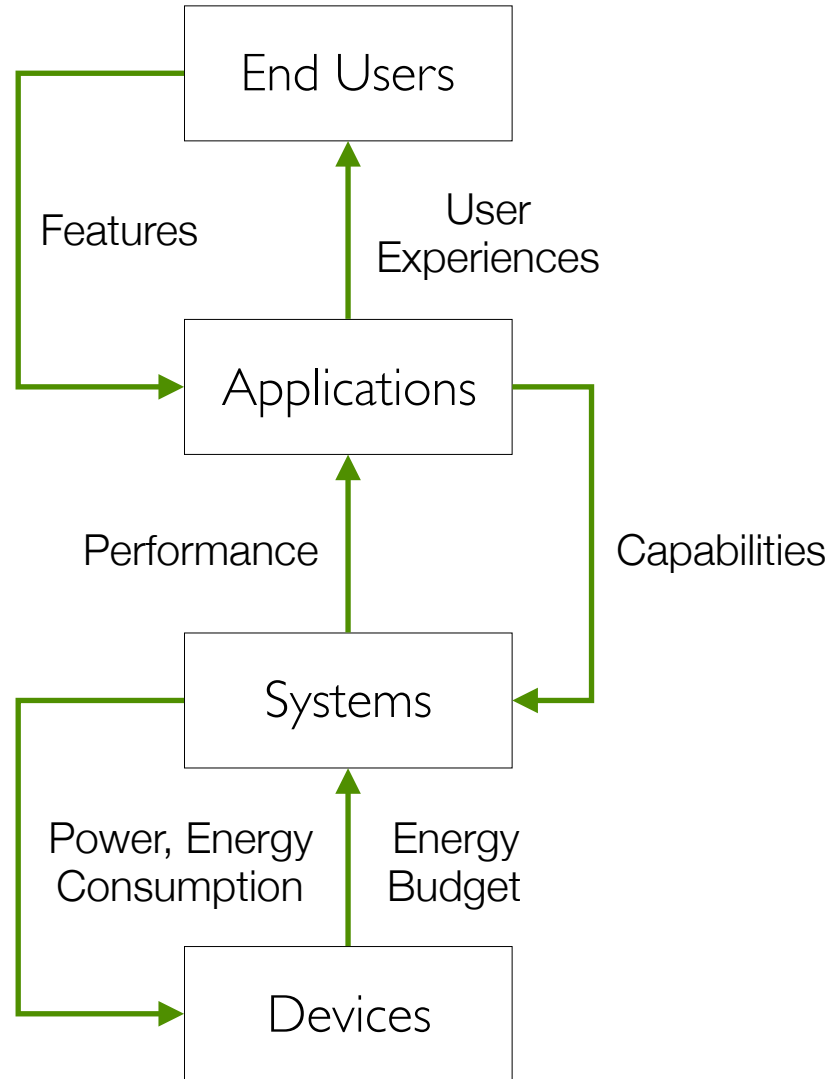


# The (Mobile) Computing Virtuous Cycle

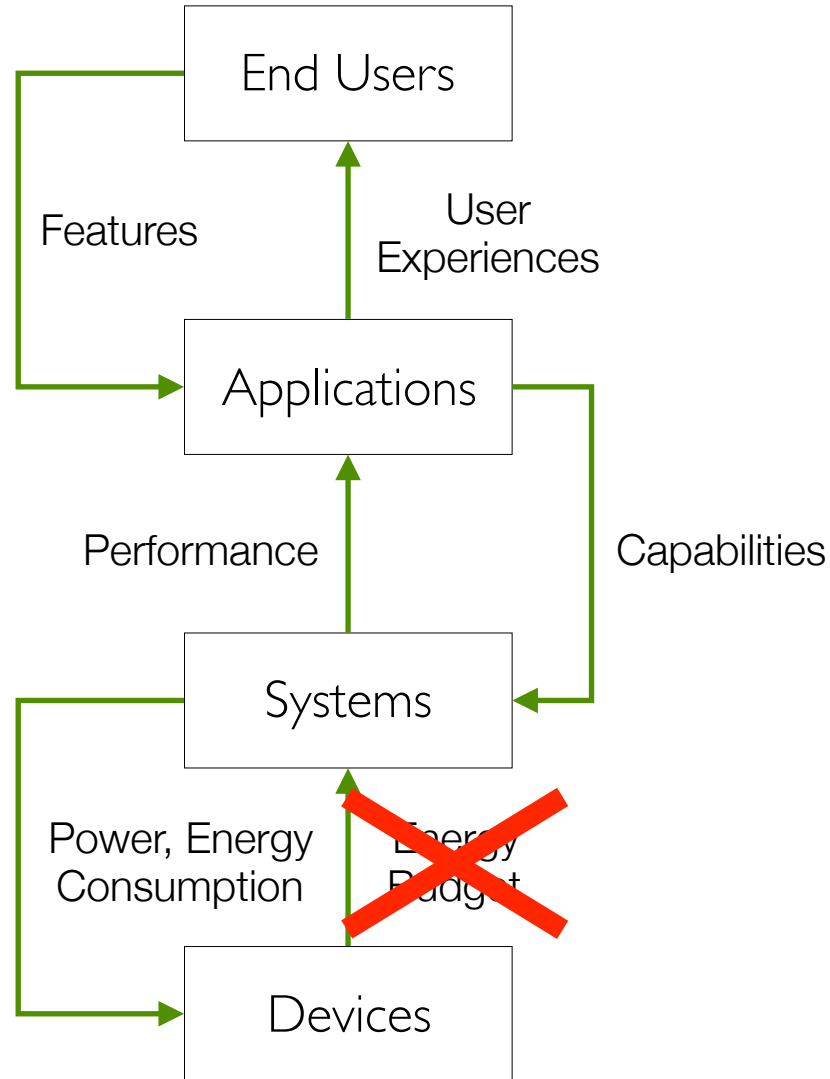




# The (Mobile) Computing Virtuous Cycle



# The (Mobile) Computing Virtuous Cycle



# Mobile CPU's Rise to Power: Quantifying the Impact of Generational Mobile CPU Design Trends on Performance, Energy, and User Satisfaction

Matthew Halpern Yuhao Zhu Vijay Janapa Reddi

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

{matthalp, yzhu}@utexas.edu, vj@ece.utexas.edu

## Abstract

*In this paper, we assess the past, present, and future of mobile CPU design. We study how mobile CPU designs trends have impacted the end-user, hardware design, and the holistic mobile device. We analyze the evolution of ten cutting-edge mobile CPU designs released over the past seven years. Specifically, we report measured performance, power, energy and user satisfaction trends across mobile CPU generations.*

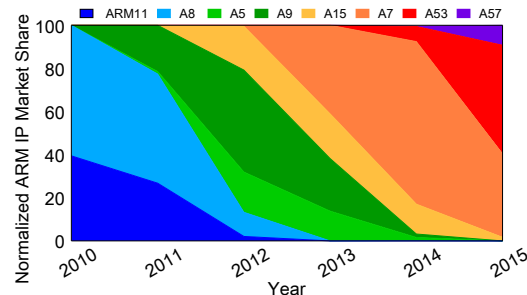
*A key contribution of our work is that we contextualize the mobile CPU's evolution in terms of user satisfaction, which has largely been absent from prior mobile hardware studies. To bridge the gap between mobile CPU design and user satisfaction, we construct and conduct a novel crowdsourcing study that spans over 25,000 survey participants using the Amazon Mechanical Turk service. Our methodology allows us to identify what mobile CPU design techniques provide the most benefit to the end-user's quality of user experience.*

*Our results quantitatively demonstrate that CPUs play a crucial role in modern mobile system-on-chips (SoCs). Over the last seven years, both single- and multicore performance improvements have contributed to end-user satisfaction by reducing user-critical application response latencies. Mobile CPUs aggressively adopted many power-hungry desktop-oriented design techniques to reach these performance levels. Unlike other smartphone components (e.g. display and radio) whose peak power consumption has decreased over time, the mobile CPU's peak power consumption has steadily increased.*

*As the limits of technology scaling restrict the ability of desktop-like scaling to continue for mobile CPUs, specialized accelerators appear to be a promising alternative that can help sustain the power, performance, and energy improvements that mobile computing necessitates. Such a paradigm shift will redefine the role of the CPU within future SoCs, which merit several design considerations based on our findings.*

## 1. Introduction

Mobile hardware design is driven by ambitious user requirements. Users demand that each generation compute faster, last longer, and fit more components into increasingly thin form factors. The fast pace at which new application use cases, wireless technologies, and sensor capabilities emerge implies that mobile system-on-chip (SoC) designs must quickly adopt



**Fig. 1: Breakdown of yearly ARM Cortex-A CPU design market share. Mobile CPU core designs have rapid design iteration and innovation. At least one new core design is released each year and newer designs overshadow the older ones.**

rate to keep pace with end-user demands. Fig. 1, based on data mined from over 1700 Android smartphone specifications, conveys the fast pace at which mobile CPU designs have evolved. Considering the ARM-based Cortex-A series alone, the most dominant mobile CPU design in smartphones and tablets to date [1], at least one new CPU core design has been released each year for the last six years – each significantly more advanced than the last. In comparison, x86-based desktop CPU designs did not exhibit as dramatic changes. Intel-based desktop processors only exhibited four significant core design changes throughout the same time span.

The rapid design innovation, pervasiveness in society, and power-constrained nature of mobile hardware necessitate the need to understand the implications of their current design trends on future designs. Mobile CPUs have evolved from embedded processors to desktop-like single-chip multiprocessors to provide application responsiveness to end-users. However, mobile CPUs, like any embedded processor, operate under a stricter set of power, thermal and energy constraints than their desktop counterparts. Therefore, there is a need to understand the effectiveness of these designs trends, both in terms of the end-users' satisfaction and hardware efficiency.

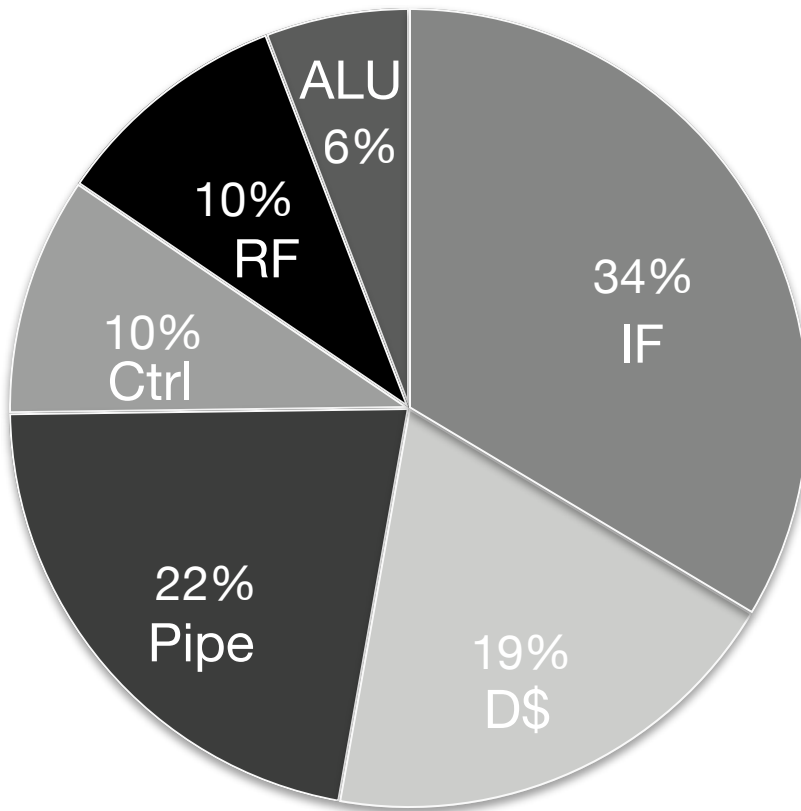
In this paper, we take the first steps towards understanding the mobile hardware evolution by studying the mobile CPU in conjunction with end-user experience. We measure and quantify the performance, power, energy, and user satisfac-

# Sources of Energy-Inefficiencies

General-Purpose CPU = Instruction Delivery + Data Feeding + Execution + Control, where instruction delivery, data feeding & control are pure overhead

# Sources of Energy-Inefficiencies

General-Purpose CPU = Instruction Delivery + Data Feeding + Execution + Control, where instruction delivery, data feeding & control are pure overhead



**IF:** Instruction fetch

**Ctrl:** Other control logics

**Pipe:** Pipeline reg, bus, clock

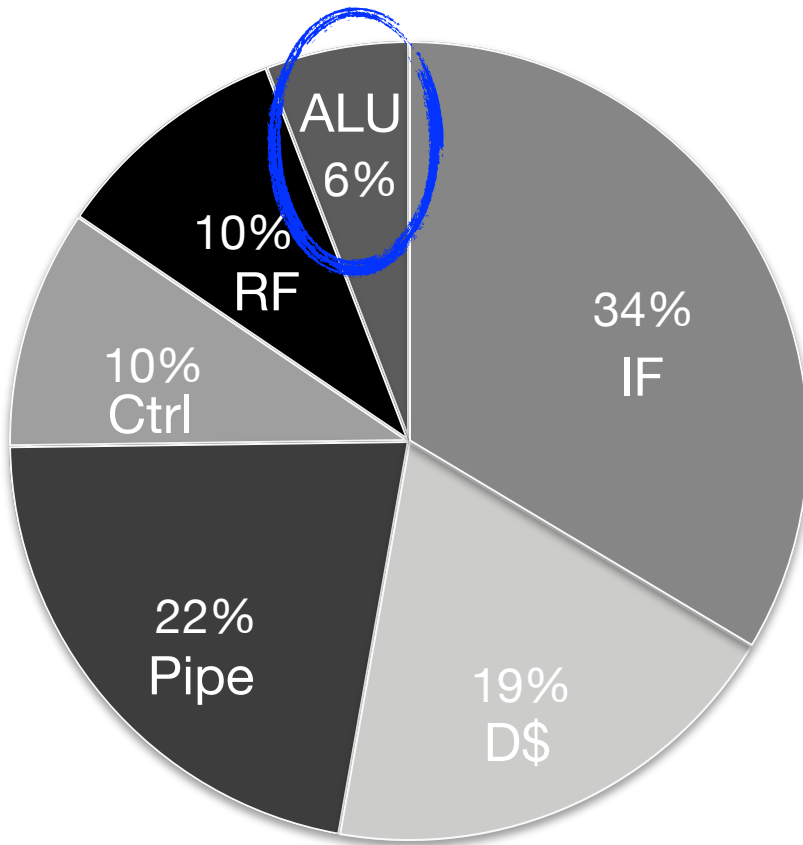
**D\$:** Data cache

**RF:** Register file

**ALU:** Functional units

# Sources of Energy-Inefficiencies

General-Purpose CPU = Instruction Delivery + Data Feeding + Execution + Control, where instruction delivery, data feeding & control are pure overhead



**IF:** Instruction fetch

**Ctrl:** Other control logics

**Pipe:** Pipeline reg, bus, clock

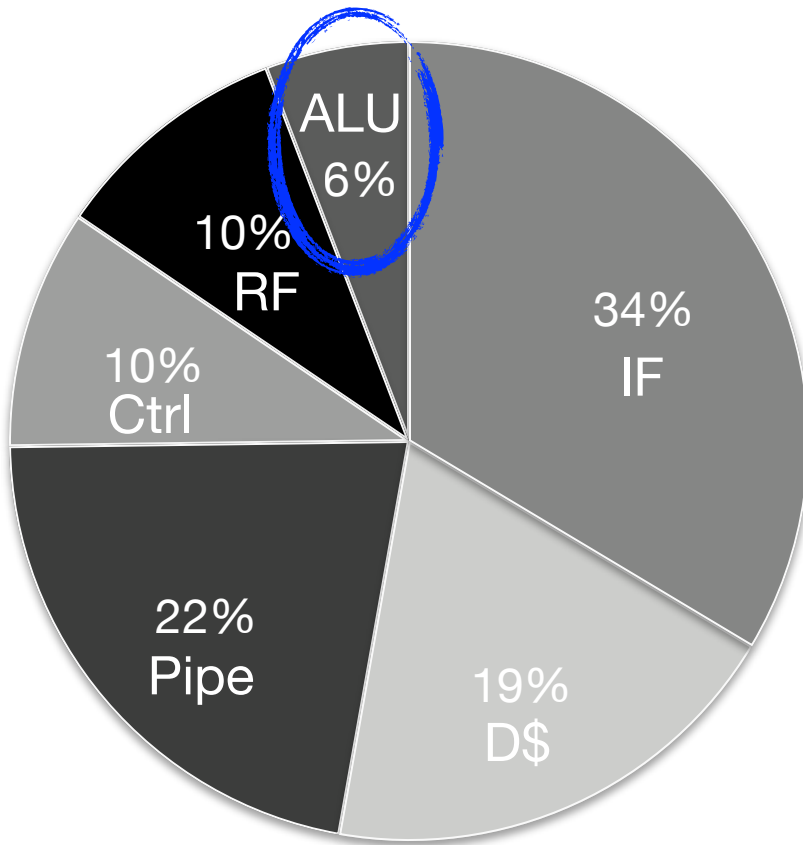
**D\$:** Data cache

**RF:** Register file

**ALU:** Functional units

# Sources of Energy-Inefficiencies

General-Purpose CPU = Instruction Delivery + Data Feeding + Execution + Control, where instruction delivery, data feeding & control are pure overhead



## Pure Overhead

**IF:** Instruction fetch

**Ctrl:** Other control logics

**Pipe:** Pipeline reg, bus, clock

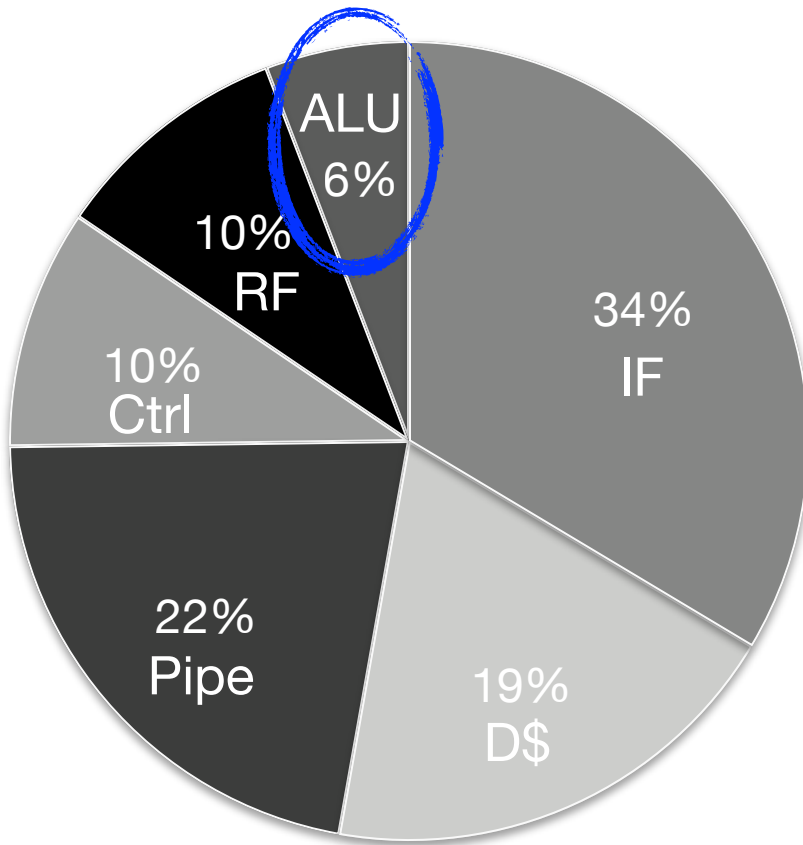
**D\$:** Data cache

**RF:** Register file

**ALU:** Functional units

# Sources of Energy-Inefficiencies

General-Purpose CPU = Instruction Delivery + Data Feeding + Execution + Control, where instruction delivery, data feeding & control are pure overhead



## Pure Overhead

**IF:** Instruction fetch

**Ctrl:** Other control logics

**Pipe:** Pipeline reg, bus, clock

**D\$:** Data cache

**RF:** Register file

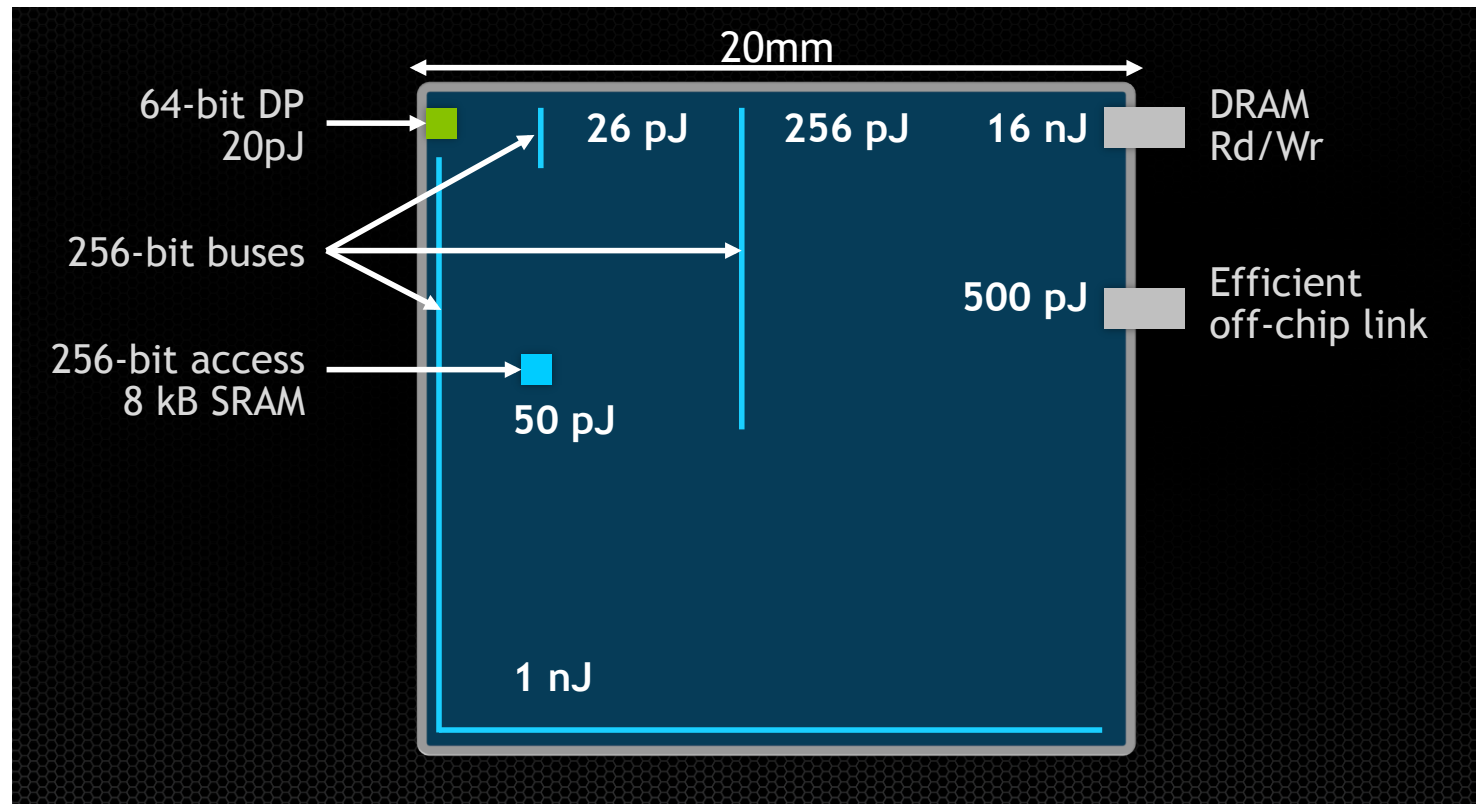
**ALU:** Functional units

**Doing Actual Work**



# Sources of Energy-Inefficiencies

General-Purpose CPU = Instruction Delivery + Data Feeding + Execution + Control, where instruction delivery, data feeding & control are pure overhead



Now What?

# Entering the Era of Specialization

# Entering the Era of Specialization

- Instead of building general-purpose processors that can do everything, but inefficiently, let's build specialized processors that can only do limited things, but extremely efficiently.

# Entering the Era of Specialization

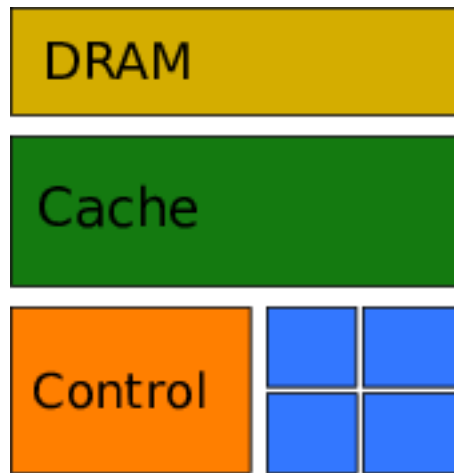
- Instead of building general-purpose processors that can do everything, but inefficiently, let's build specialized processors that can only do limited things, but extremely efficiently.
- Extreme example: Application-Specific Integrated Circuit (ASIC)

# Entering the Era of Specialization

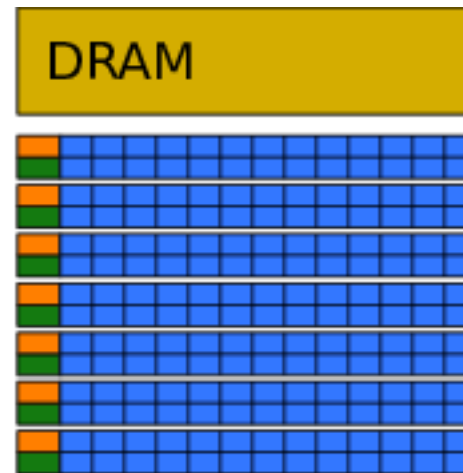
- Instead of building general-purpose processors that can do everything, but inefficiently, let's build specialized processors that can only do limited things, but extremely efficiently.
- Extreme example: Application-Specific Integrated Circuit (ASIC)
- Imagine the entire processor is an FP adder with 2 registers
  - Instruction delivery: none. Instructions are implicit (add).
  - Data feeding: simple.
  - Control: little: clock the adder, simple pipeline, etc.
  - Execution: will be the major power consumer.
  - Flexibility/programmability: very limited: addition only, nothing else!

# Entering the Era of Specialization

- Another example: Graphics Processing Unit (GPU)
  - SIMT: Single instruction multiple thread
  - SIMT amortizes control/instruction delivery overhead
  - Data feeding is still very complex
  - More efficient execution for applications that are massively parallel

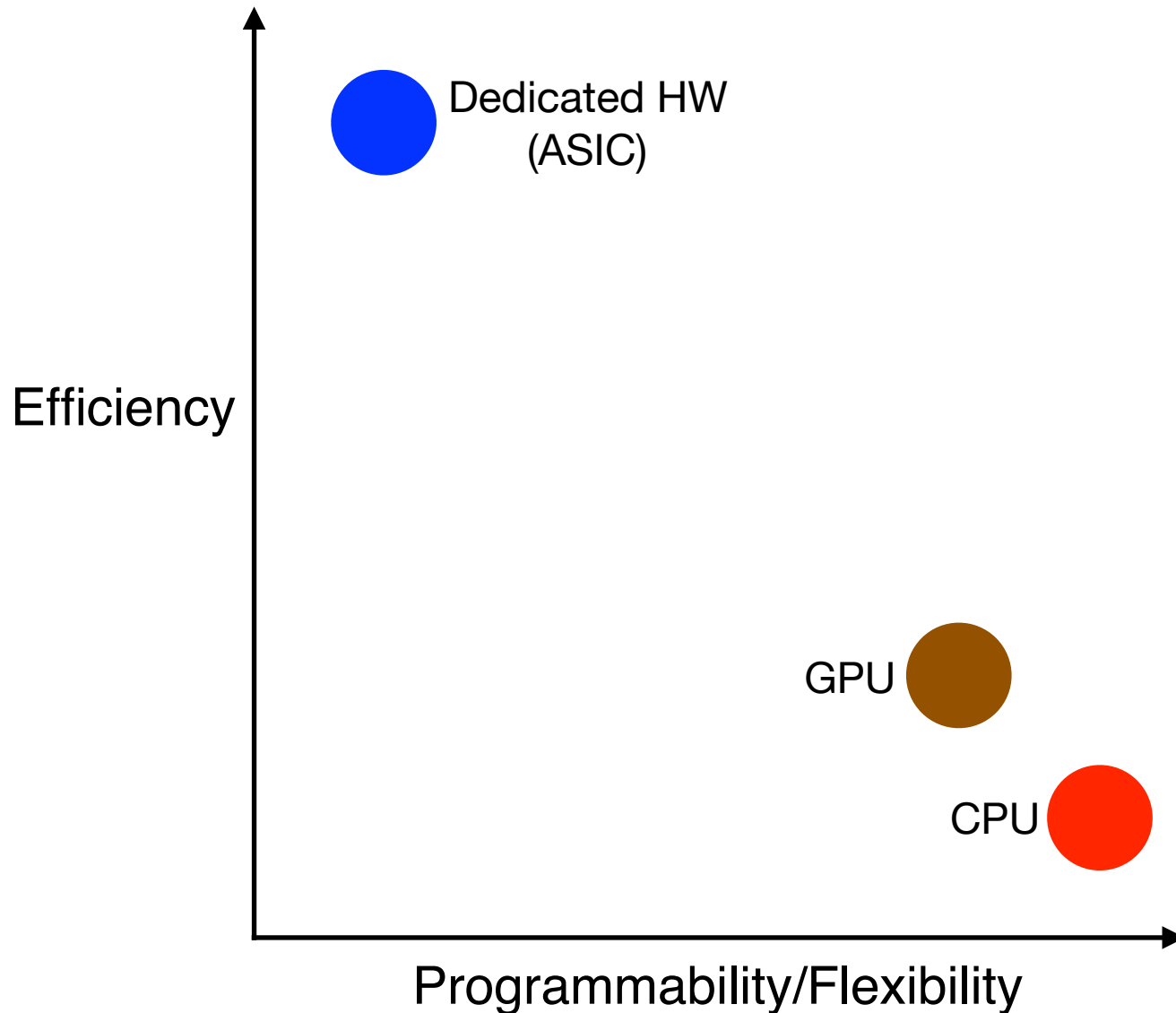


CPU



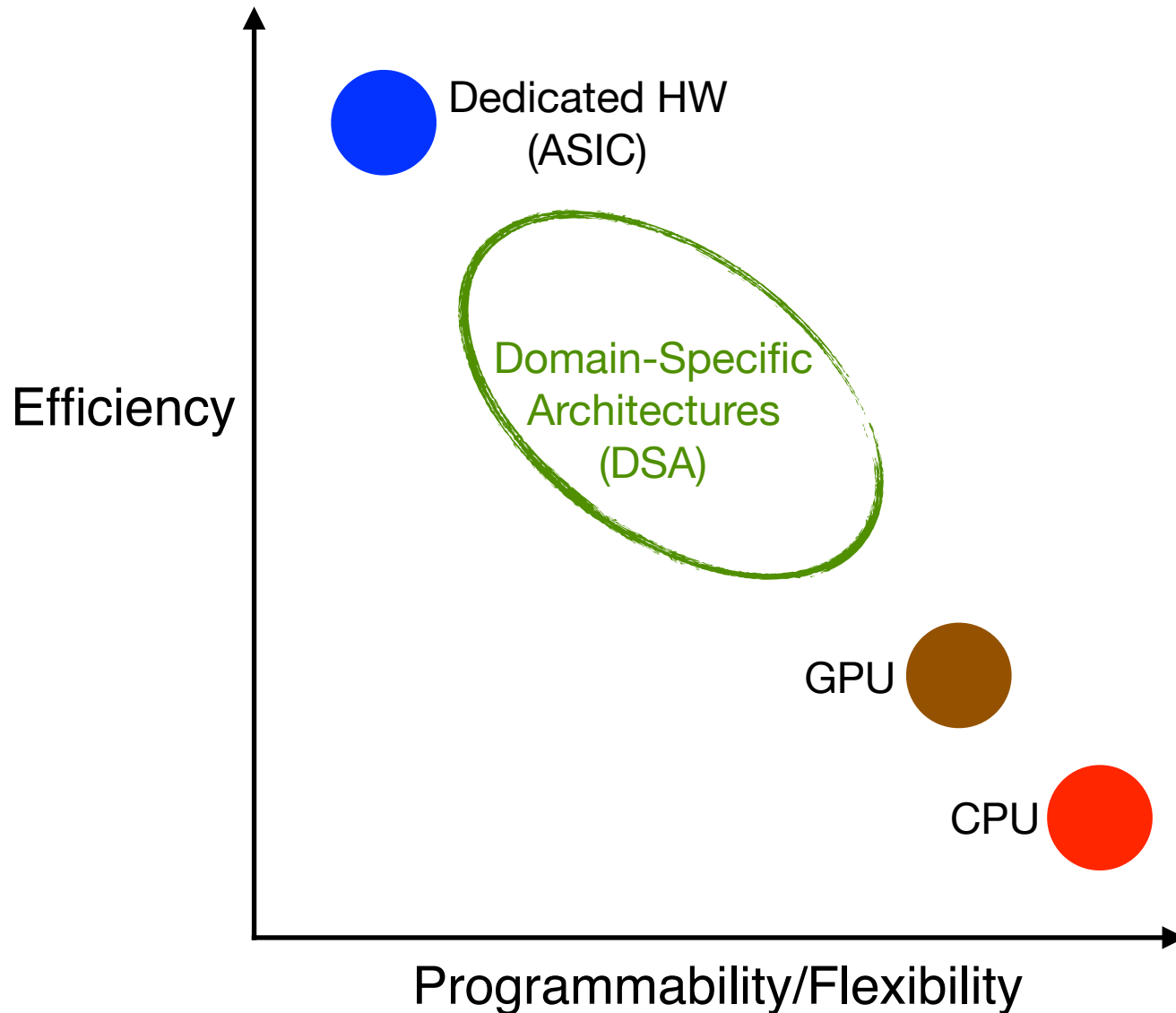
GPU

# Entering the Era of Specialization





# Entering the Era of Specialization



# Example: Video (De)-Compression





# Example: Video (De)-Compression

30 second video: 1920 x 1080, @ 30fps

8-bits per color  $\rightarrow$  24 bits/pixel  $\rightarrow$  6.2MB/frame ( $6.2 \text{ MB} \times 30 \text{ sec} \times 30 \text{ fps} = 5.2 \text{ GB}$ )

Actual H.264 video file size: 65.4 MB (80-to-1 compression ratio)

Compression/encoding performed in real time on iPhone 5s



# Example: Computational Photography

- High-quality cameras are key product differentiators for today's consumer smartphones.



# Example: Computational Photography

- High-quality cameras are key product differentiators for today's consumer smartphones.
- But smartphone cameras aren't as good as high-end digital single-lens reflex cameras.
  - Aperture, focal length, pixels, dynamic range, sensors, resolution, etc.

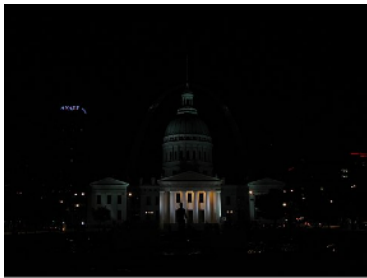


# Example: Computational Photography

- High-quality cameras are key product differentiators for today's consumer smartphones.
- But smartphone cameras aren't as good as high-end digital single-lens reflex cameras.
  - Aperture, focal length, pixels, dynamic range, sensors, resolution, etc.
- How to achieve DSLR-like quality in a smartphone form factor?

# Example: Computational Photography

- High-quality cameras are key product differentiators for today's consumer smartphones.
- But smartphone cameras aren't as good as high-end digital single-lens reflex cameras.
  - Aperture, focal length, pixels, dynamic range, sensors, resolution, etc.
- How to achieve DSLR-like quality in a smartphone form factor?
- Computational photography: digitally post-processing images
  - E.g., High-dynamic range (HDR)
  - Used to be just in desktop photo editing software (e.g., PhotoShop), but now in high-end smartphones

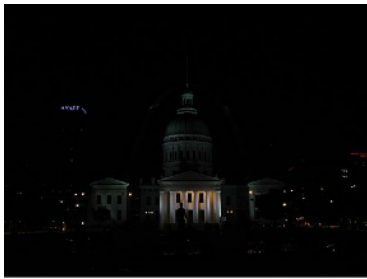


**Short Exposure**



**Long Exposure**





**Short Exposure**



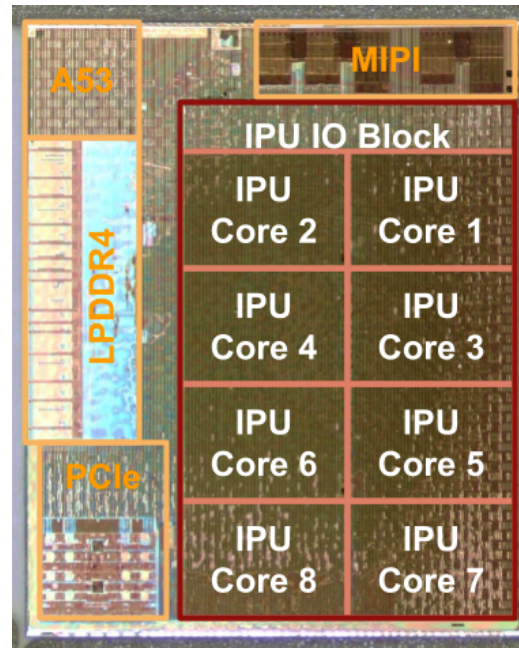
**Long Exposure**



# Example: Computational Photography

- ISP (Image Signal Processor): Specialized processors for imaging and computational photography algorithms.
- Fast, energy-efficient, but only for imaging

**Google  
Pixel 2**





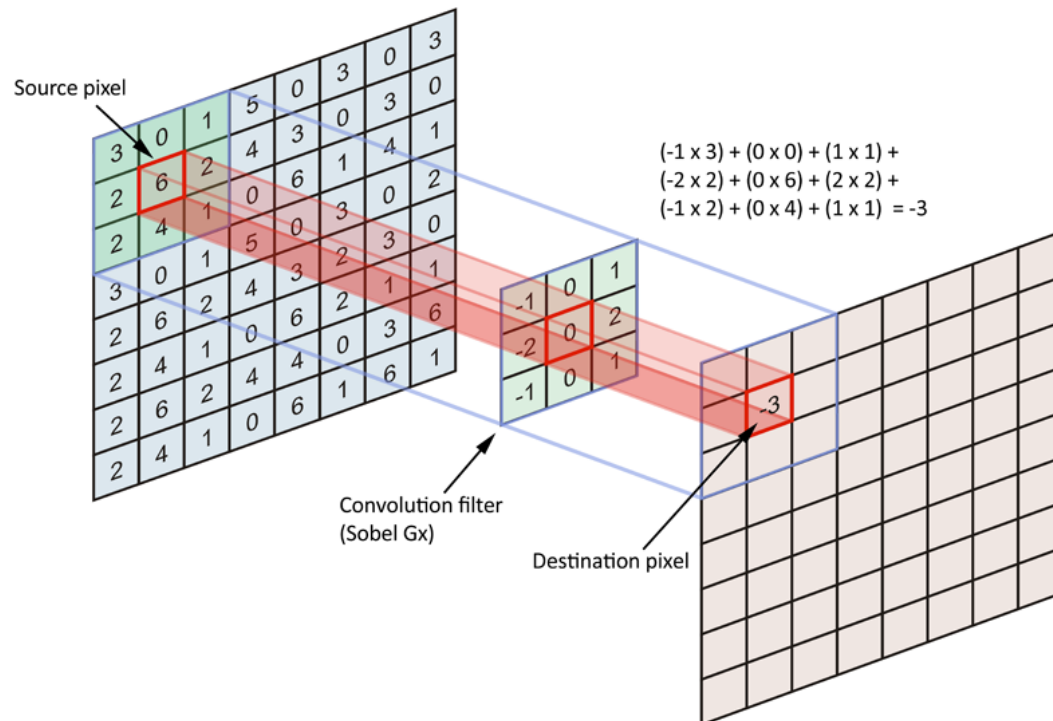
# Example: Machine Learning

- Machine learning (e.g., Neural Networks) takes over the world
- Computer vision is the poster-child example
- Natural language processing, Precision medicine, robotics planning, house rent/load prediction, games, etc.



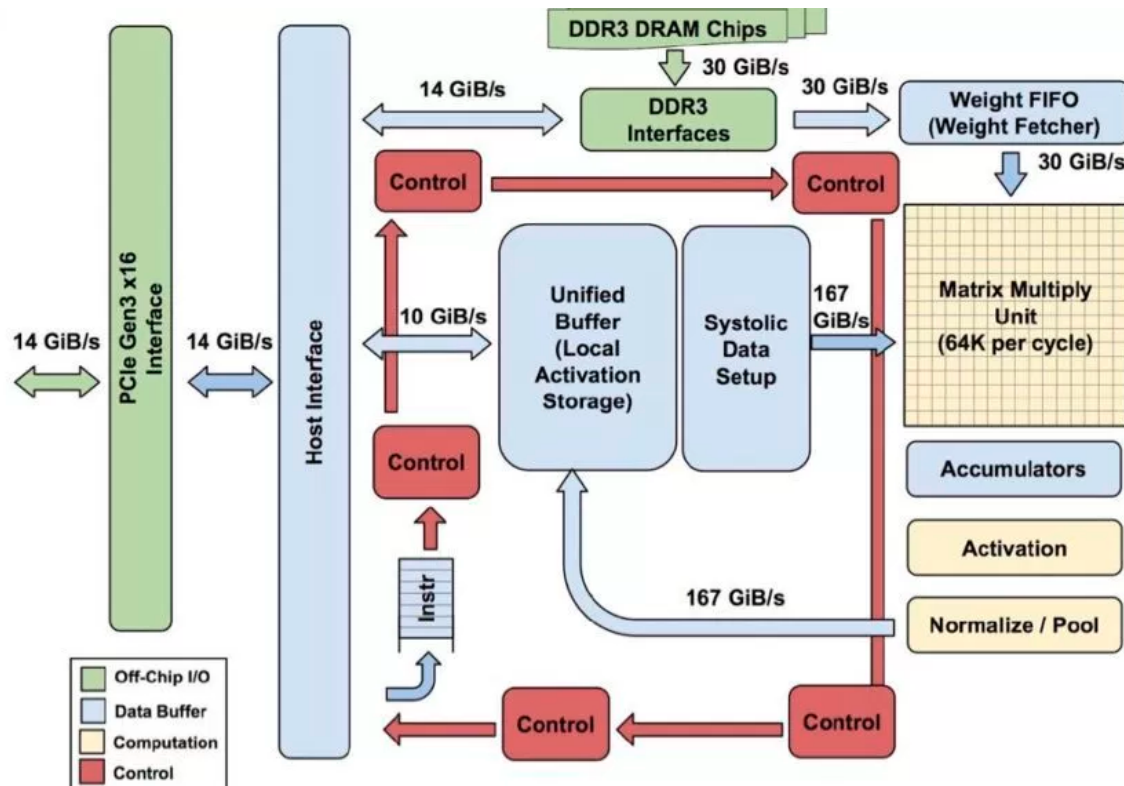
# Example: Machine Learning

- Neural networks make heavy use of the convolution operation
  - Convolution can be transformed to matrix-multiplication
  - Convolutional neural networks (CNN) are by far the most popular NN, especially effectively for computer vision tasks

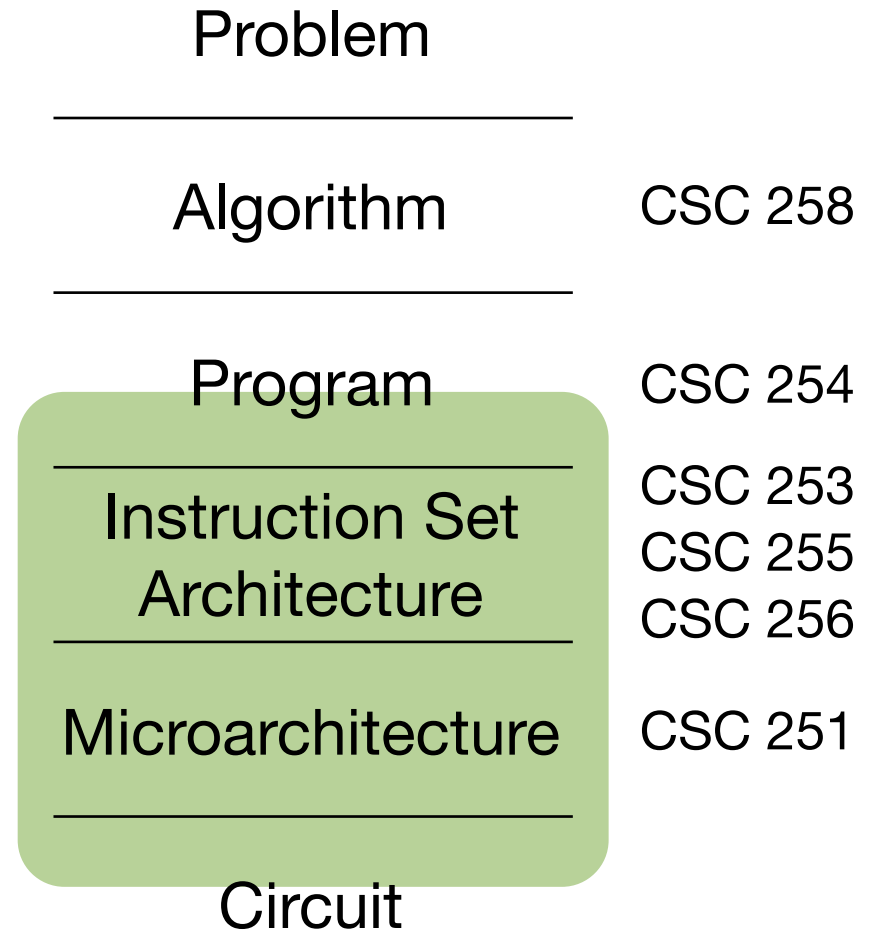


# Example: Machine Learning

- Google Tensor Processing Unit (TPU)
  - Specialized processor (i.e., systolic array architecture) for tensor processing (matrix multiply)
  - (Arguably) 30x~80x more power-efficient than GPU

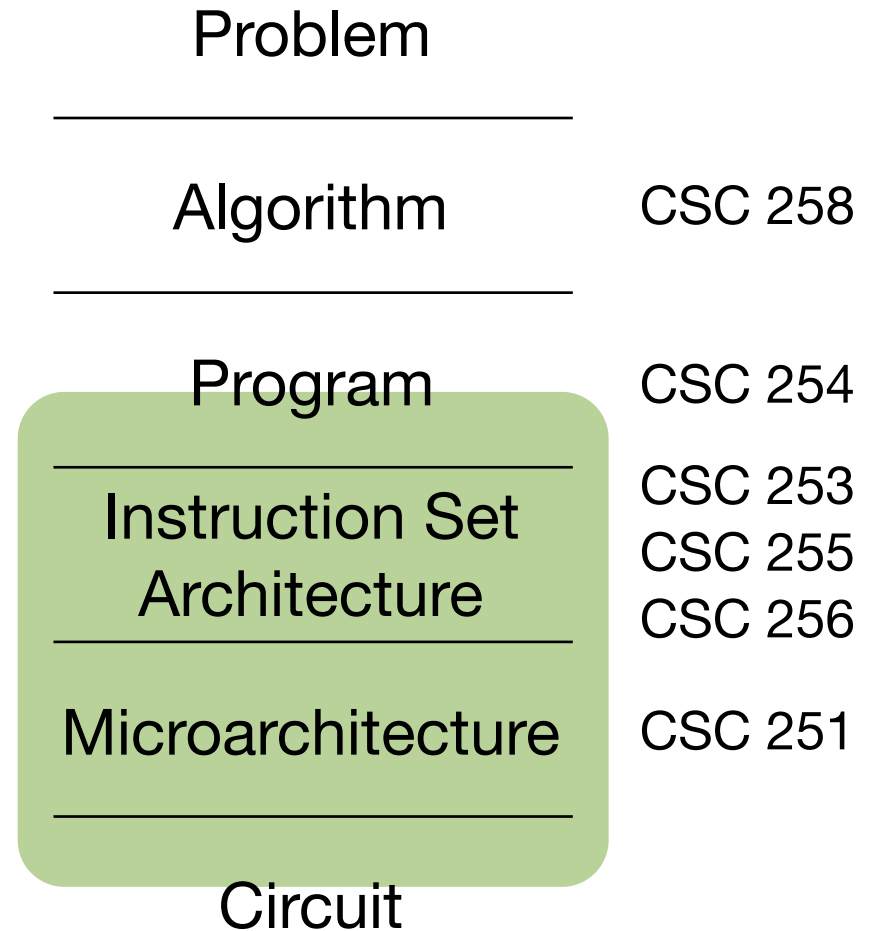


# The Role of a Computer System Designer



# The Role of a Computer System Designer

- Look Up
  - Nature of the problems





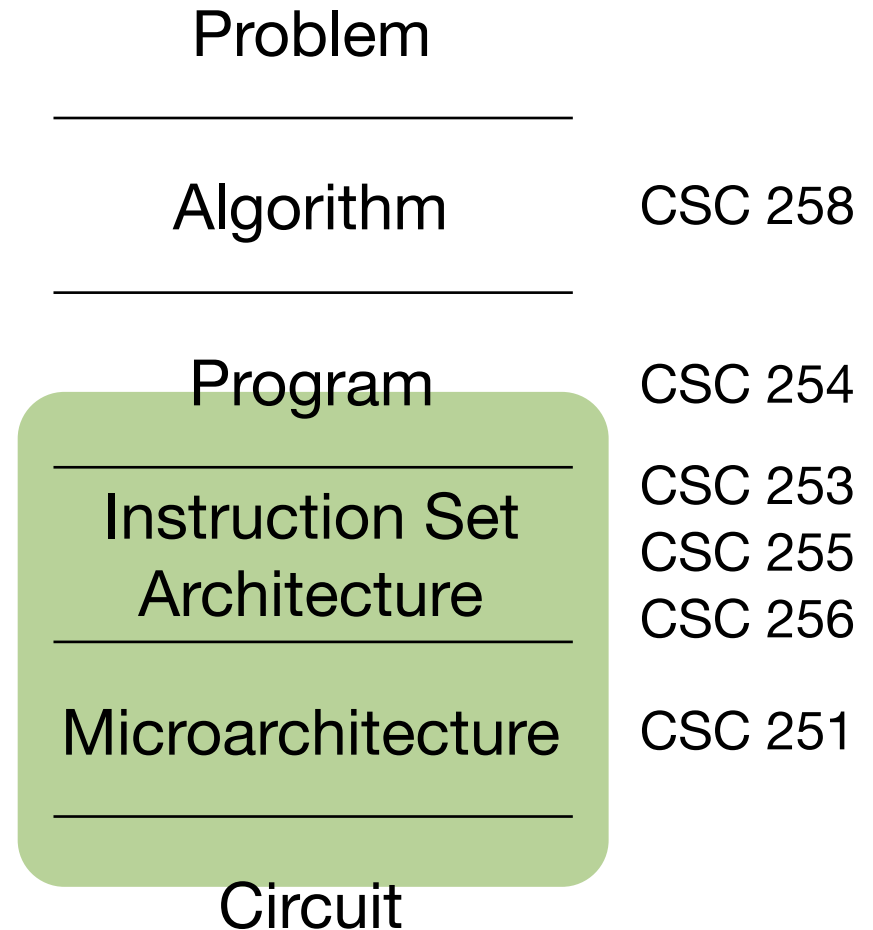
# The Role of a Computer System Designer

- **Look Up**

- Nature of the problems

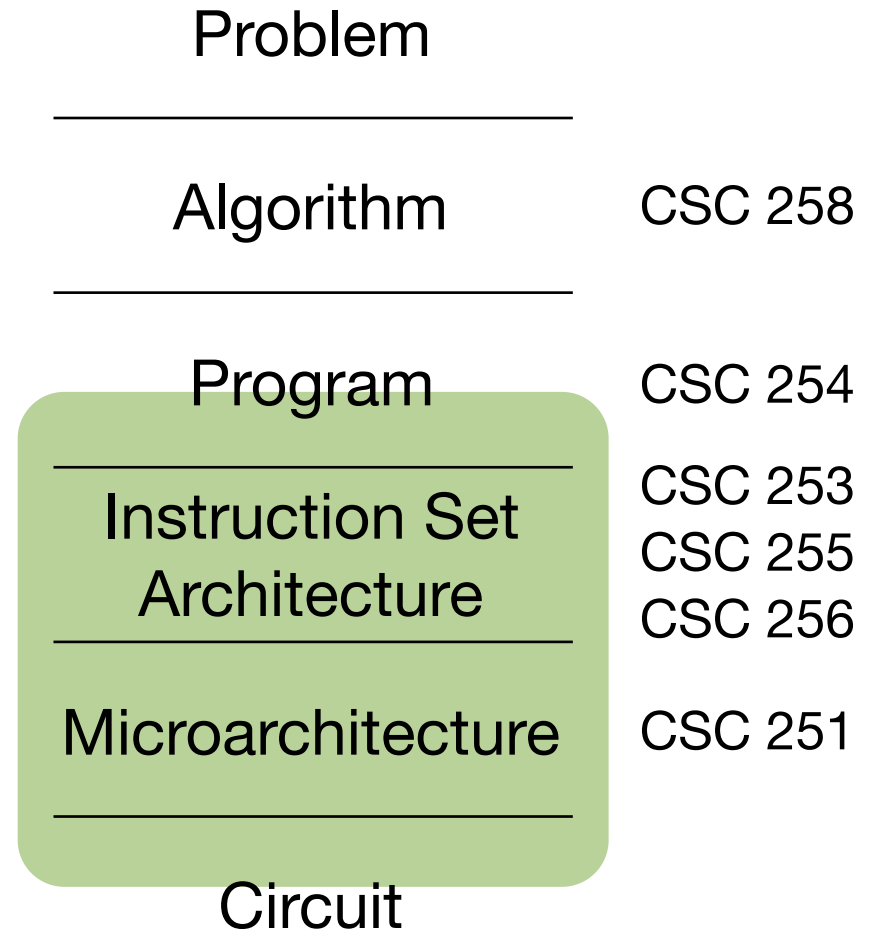
- **Look Down**

- Nature of the circuit technology and physics



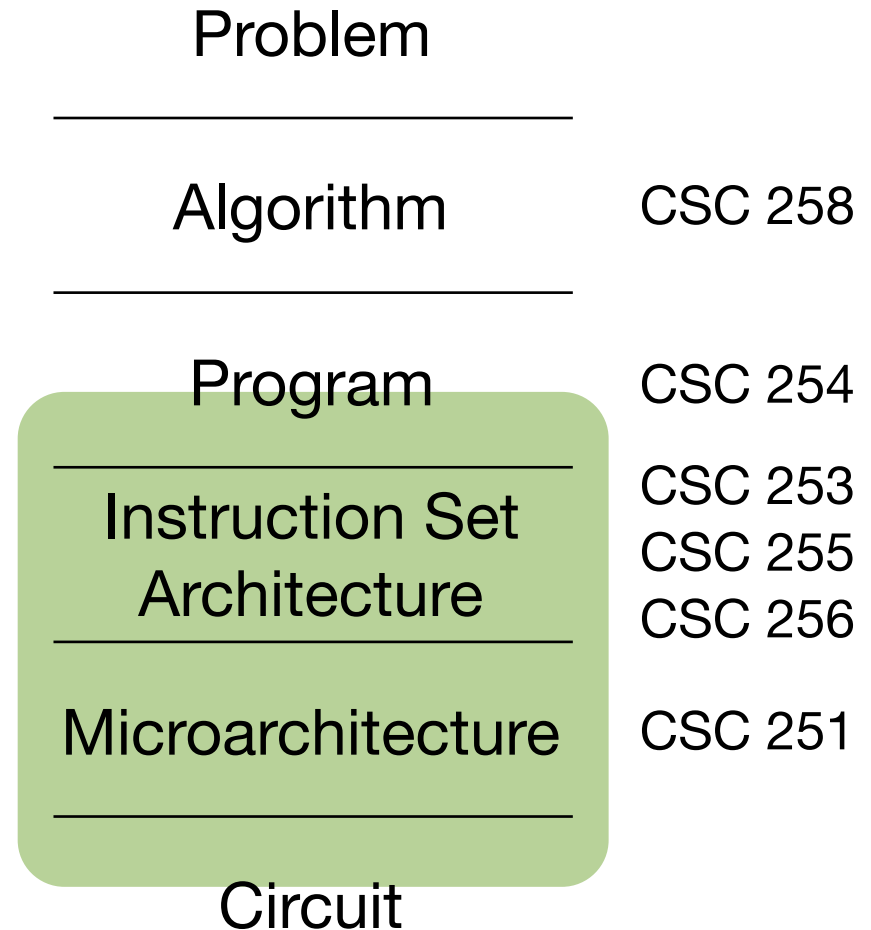
# The Role of a Computer System Designer

- **Look Up**
  - Nature of the problems
- **Look Down**
  - Nature of the circuit technology and physics
- **Look Backward**
  - Evaluating old ideas in light of new technologies



# The Role of a Computer System Designer

- **Look Up**
  - Nature of the problems
- **Look Down**
  - Nature of the circuit technology and physics
- **Look Backward**
  - Evaluating old ideas in light of new technologies
- **Look Forward**
  - Listen to dreamers and predict the future



# CSC 292/572: Mobile Visual Computing

