# CSC 252: Computer Organization Spring 2021: Lecture 26

## Instructor: Yuhao Zhu

Department of Computer Science
University of Rochester

# Announcements

| SUN 25 | MON 26 | TUE 27 | WED 28 | THU 29 | FRI 30 | SAT May 1 |
|---|---|---|---|---|---|---|
| 2 | 3 | 4 | 5 | 6 **Today Last Lecture** | 7 | 8 **A5 Due** |
| 9 | 10 | 11 | 12 **Final** | 13 | 14 | 15 |

# Announcements

- Final exam: May 12, 19:15 PM -- 22:15 PM; online.

- Past exam & Problem set: https://www.cs.rochester.edu/courses/252/spring2021/handouts.html

- Exam will be electronic using Gradescope, 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.

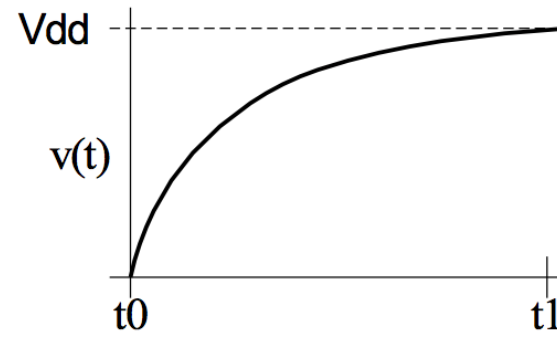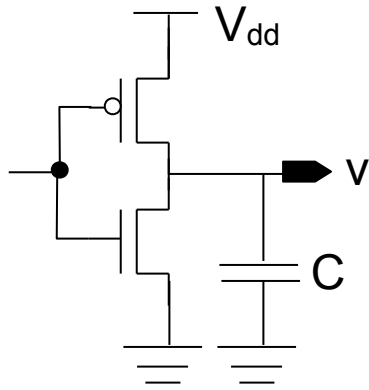  - Write down the answers on a scratch paper, take pictures, and send us the pictures

# Announcements

- Open book test: any sort of paper-based product, e.g., book, **notes**, magazine, old tests.

- Exams are designed to test your ability to apply what you have learned and not your memory (though a good memory could help).

- **Nothing electronic (including laptop, cell phone, calculator, etc) other than the computer you use to take the exam**.

- **Nothing biological**, including your roommate, husband, wife, your hamster, another professor, etc.

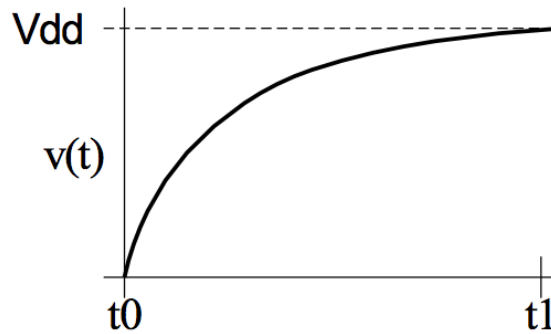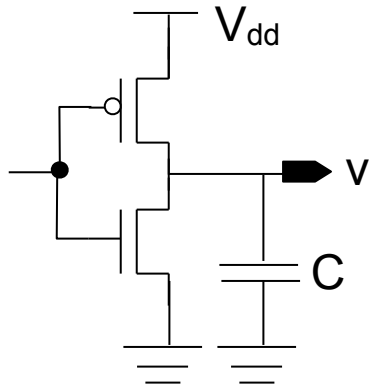- "**I don't know**" gets15% partial credit. Must erase everything else.
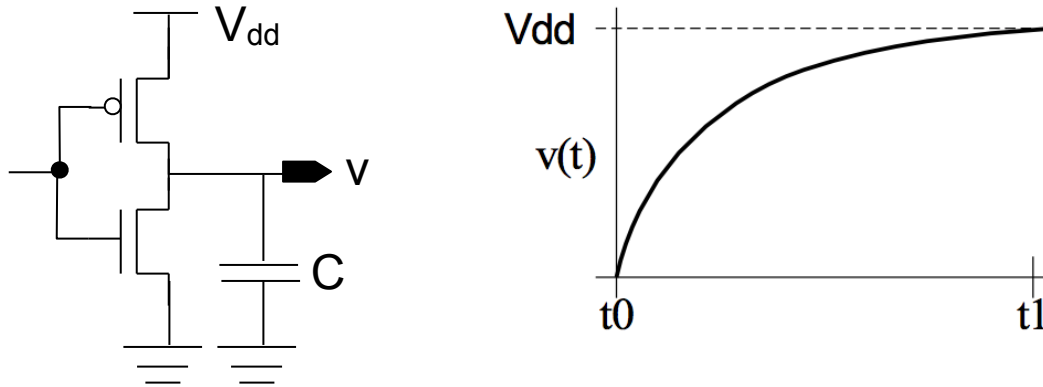
# Dynamic Power

# Dynamic Power

# Dynamic Power



$$E_{sw} = \int_{t_0}^{t_1} P(t)dt = \int_{t_0}^{t_1} (V_{dd} - v) \cdot i(t)dt = \int_{t_0}^{t_1} (V_{dd} - v) \cdot c \, (dv/dt) \, dt =$$

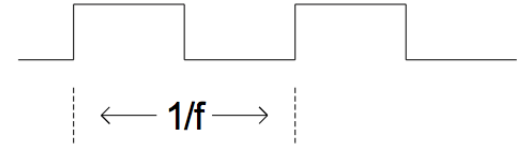$$= cV_{dd} \int_{t_0}^{t_1} dv - c \int_{t_0}^{t_1} v \cdot dv = cV_{dd}^{\,2} - 1/2 \, cV_{dd}^{\,2} = \boxed{1/2 \, cV_{dd}^{\,2}}$$

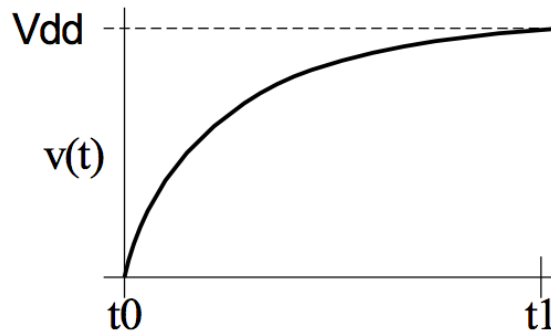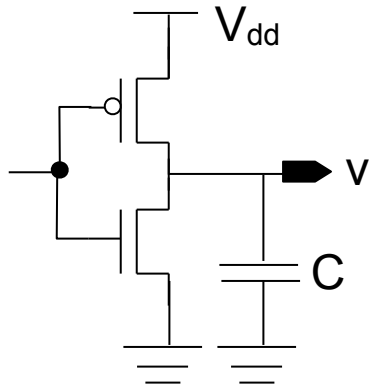# Dynamic Power



Energy dissipated for every transition (0->1 or 1->0)

$$E_{sw} = \int_{t_0}^{t_1} P(t)dt = \int_{t_0}^{t_1} (V_{dd} - v) \cdot i(t)dt = \int_{t_0}^{t_1} (V_{dd} - v) \cdot c \left( dv/dt \right) dt =$$

$$= cV_{dd} \int_{t_0}^{t_1} dv - c \int_{t_0}^{t_1} v \cdot dv = cV_{dd}^2 - 1/2 cV_{dd}^2 = \boxed{1/2 cV_{dd}^2}$$

# Dynamic Power



**Average dynamic power of a transistor:**
$$P = \alpha \cdot (E / T) = \alpha \cdot E\, f = \tfrac{1}{2}\, \alpha\, C\, V_{dd}^2\, f$$

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

# Dynamic Power

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

# Dynamic Power

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

• Increasing f requires V to be increased proportionally

# Dynamic Power

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

- Increasing f requires V to be increased proportionally
  - Intuitively: a higher frequency means a shorter cycle time, which means the critical path of your processor needs to be shorter, which requires faster transistors, which you get by increasing the voltage

# Dynamic Power

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

- Increasing $f$ requires $V$ to be increased proportionally
  - Intuitively: a higher frequency means a shorter cycle time, which means the critical path of your processor needs to be shorter, which requires faster transistors, which you get by increasing the voltage
  - "Overclocking" just increases the clock speed without increasing voltage => machine might crash (cycle time shorter than the critical path delay)

# Dynamic Power

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

- Increasing f requires V to be increased proportionally
  - Intuitively: a higher frequency means a shorter cycle time, which means the critical path of your processor needs to be shorter, which requires faster transistors, which you get by increasing the voltage
  - "Overclocking" just increases the clock speed without increasing voltage => machine might crash (cycle time shorter than the critical path delay)
  - Corollary: reducing voltage requires reducing frequency

# Dynamic Power

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

- Increasing $f$ requires $V$ to be increased proportionally
  - Intuitively: a higher frequency means a shorter cycle time, which means the critical path of your processor needs to be shorter, which requires faster transistors, which you get by increasing the voltage
  - "Overclocking" just increases the clock speed without increasing voltage => machine might crash (cycle time shorter than the critical path delay)
  - Corollary: reducing voltage requires reducing frequency
  - 15% reduction in voltage requires about 15% slow down in frequency

# Dynamic Power

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

- Increasing f requires V to be increased proportionally
  - Intuitively: a higher frequency means a shorter cycle time, which means the critical path of your processor needs to be shorter, which requires faster transistors, which you get by increasing the voltage
  - "Overclocking" just increases the clock speed without increasing voltage => machine might crash (cycle time shorter than the critical path delay)
  - Corollary: reducing voltage requires reducing frequency
  - 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.

# Dynamic Power Favors Parallelisms

$$P = k\ C\ f^3$$

- Dynamic power favors parallel processing over higher clock rate

# Dynamic Power Favors Parallelisms

$$P = k\, C\, f^3$$

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

# Dynamic Power Favors Parallelisms

$$P = k \, C \, f^3$$

- 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**!

# Dynamic Power Favors Parallelisms

$$P = k\ C\ f^3$$

- 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

# Dynamic Power Favors Parallelisms

$$P = k\,C\,f^3$$

- 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
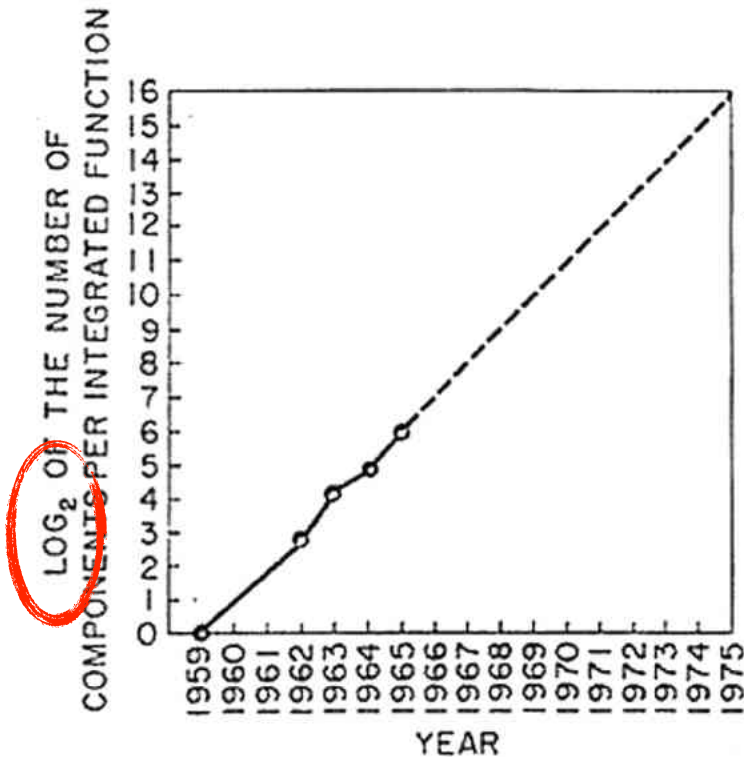
# Dynamic Power Favors Parallelisms

$$P = k\,C\,f^3$$

- 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
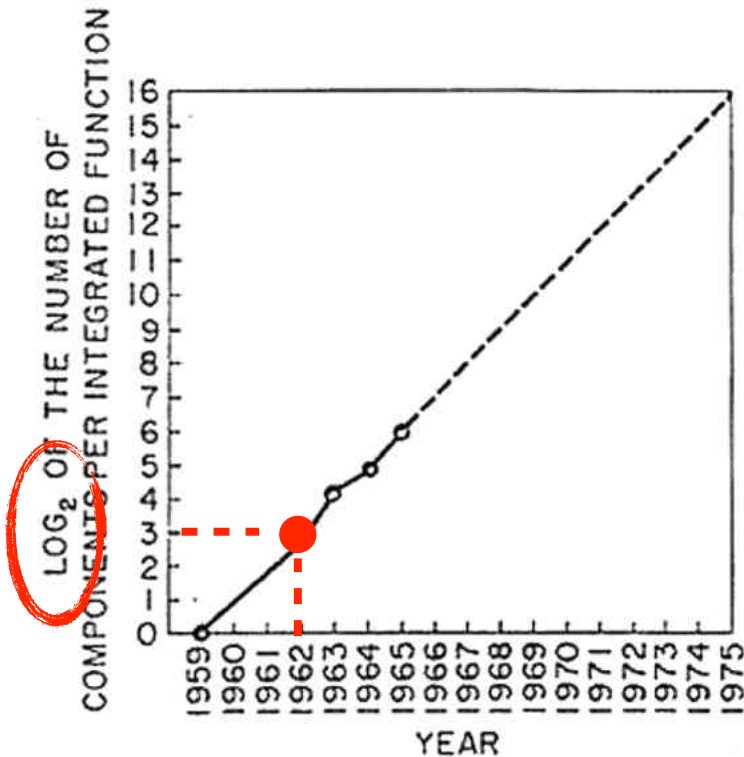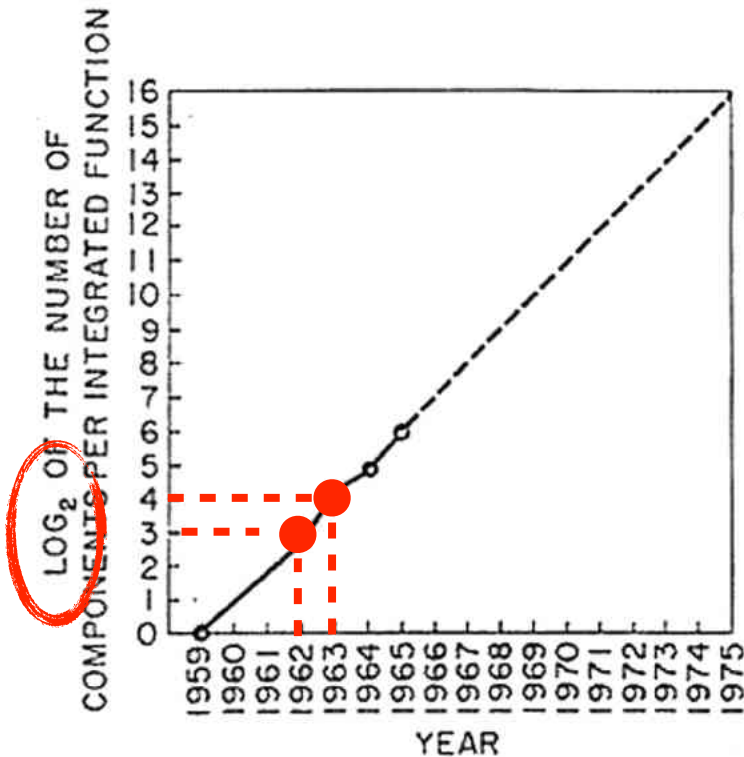  - Dynamic power becomes $4 \times (1/4)^3 = 1/16$

# 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
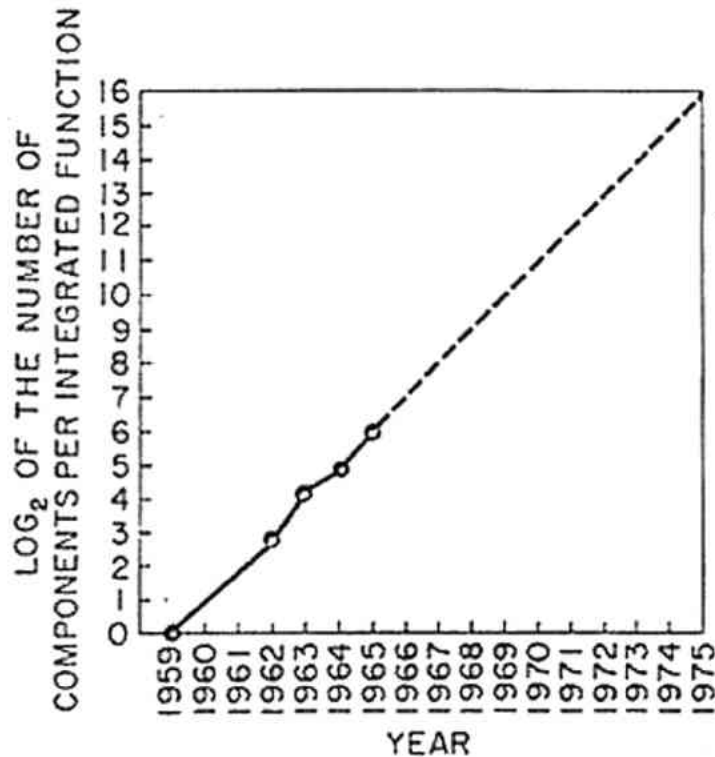
# Moore's Law

- 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
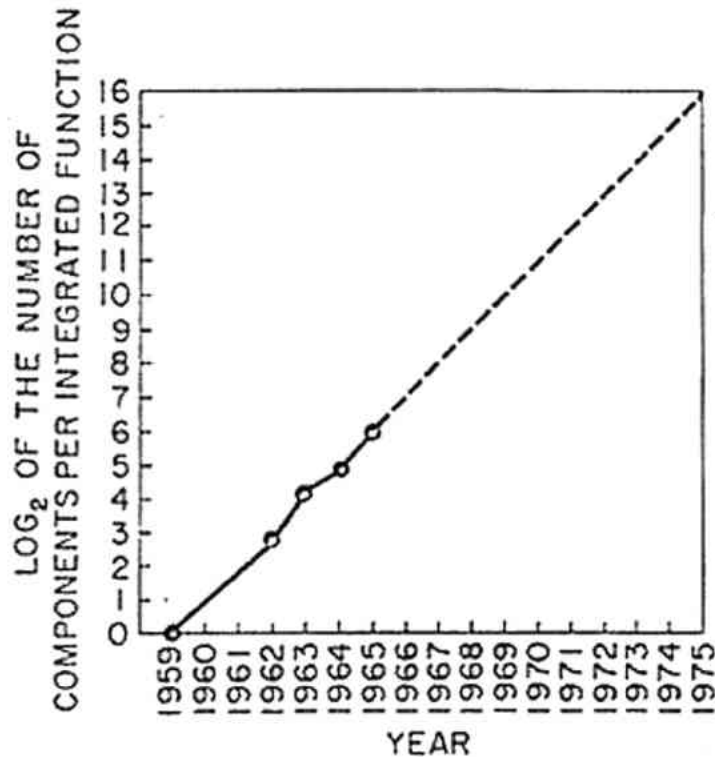
# Moore's Law

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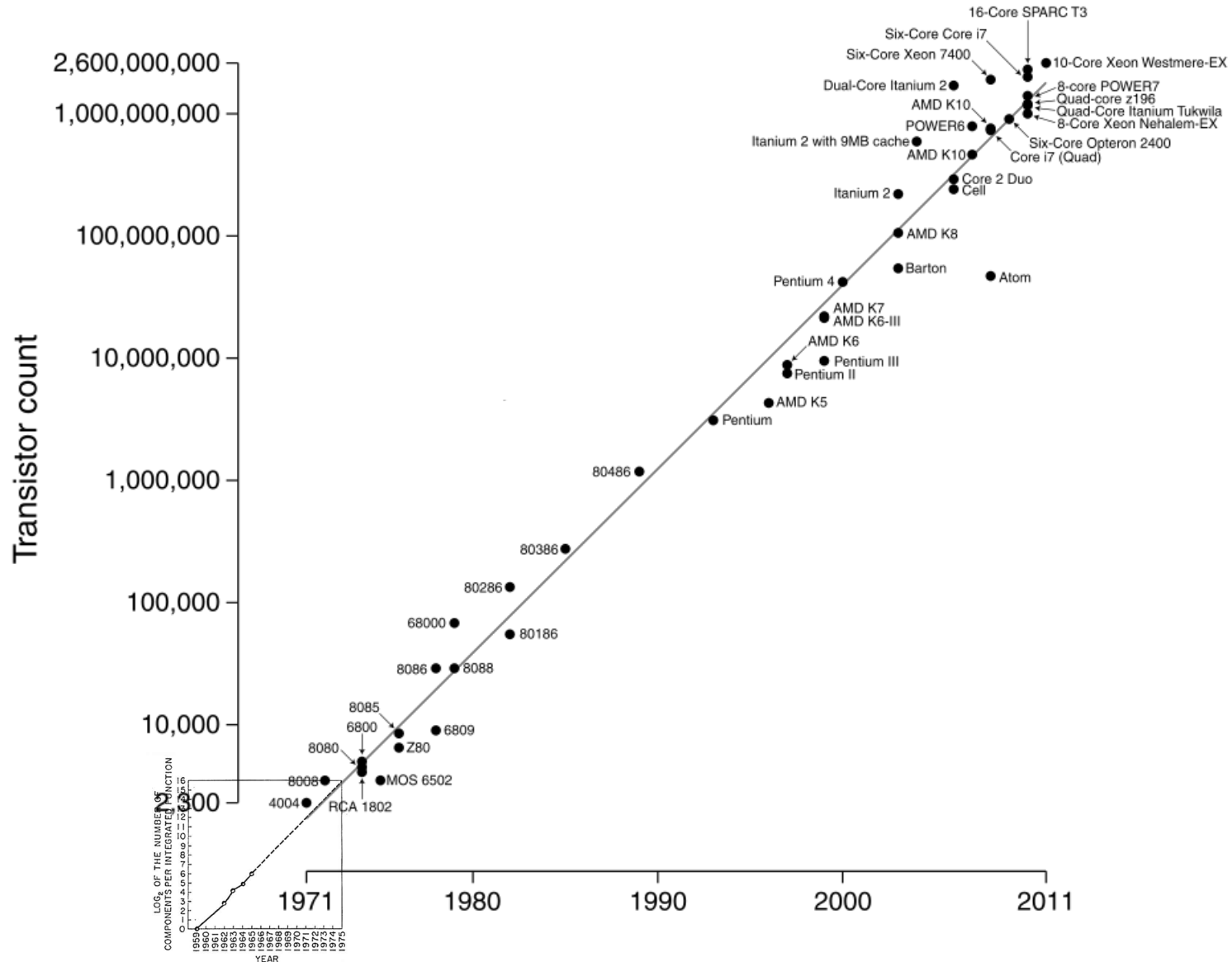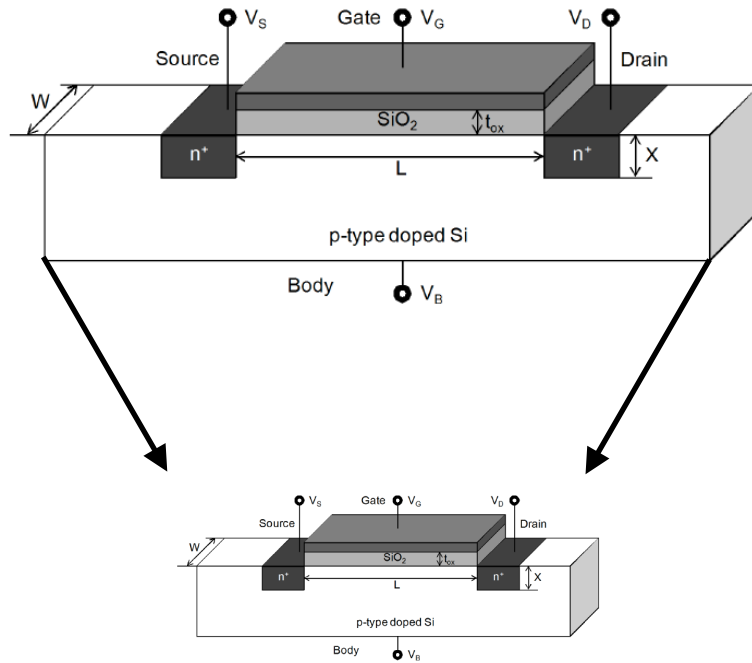
# Moore's Law

- 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…)
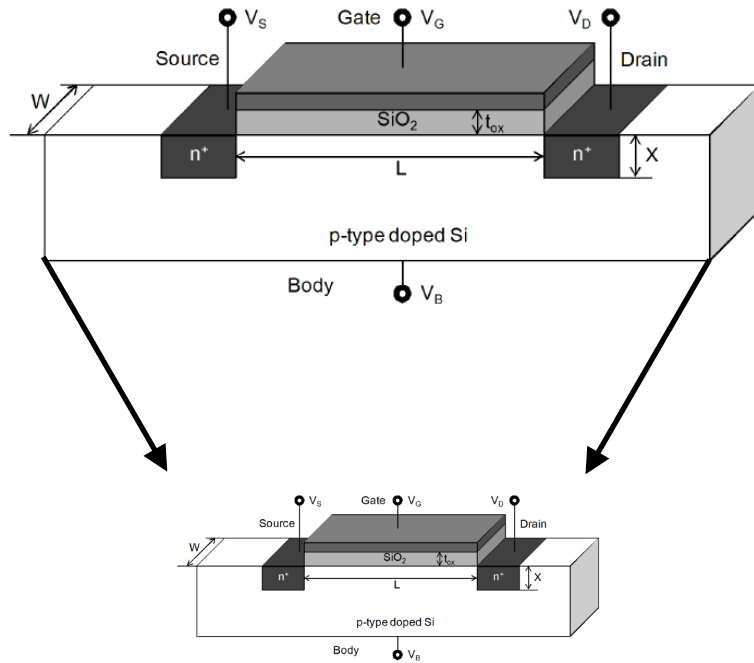
# Moore's Law

# Dennard Scaling



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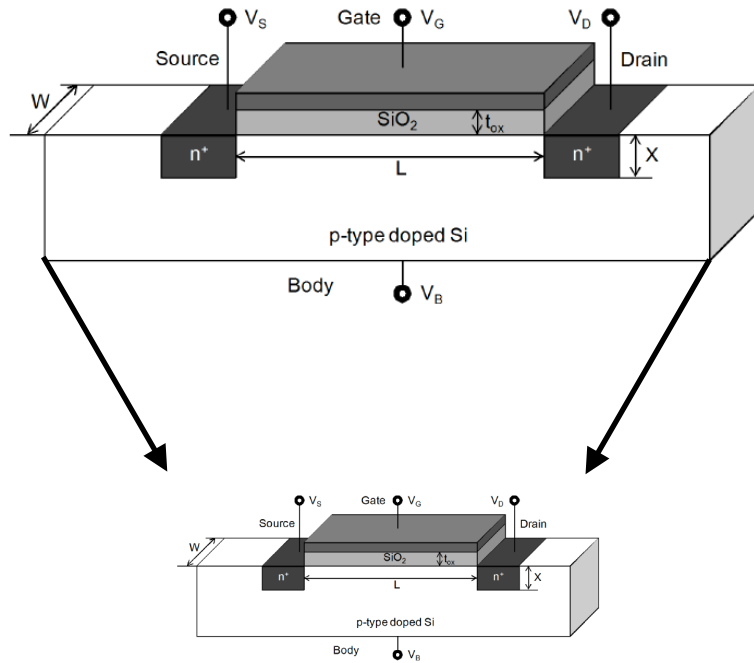
$\alpha = 0.7 \Rightarrow 2$X more transistors!

Design of Ion-Implanted MOSFET's with Very Small Physical Dimensions. Proc. of IEEE. Dennart et al. 1974.

# Dennard Scaling



| Parameter | Value | Scaled Value |
|---|---|---|
| Dopant concentrations | Na, Nd | Na/α, Nd/α |
| Dimensions | L, W, Tox | αL, αW, αTox |
| Field | E | E |
| Voltage | V | αV |
| Capacitance | C | αC |
| Current | I | αI |

*Scale factor* $\alpha < 1$

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

Design of Ion-Implanted MOSFET's with Very Small Physical Dimensions. Proc. of IEEE. Dennart et al. 1974.

# Dennard Scaling



| Parameter | Value | Scaled Value |
|---|---|---|
| Dopant concentrations | Na, Nd | Na/α, Nd/α |
| Dimensions | L, W, Tox | αL, αW, αTox |
| Field | E | E |
| Voltage | V | αV |
| Capacitance | C | αC |
| Current | I | αI |

| Transistors/Area | d | d/α² |
|---|---|---|

*Scale factor $\alpha < 1$*

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

Design of Ion-Implanted MOSFET's with Very Small Physical Dimensions. Proc. of IEEE. Dennart et al. 1974.

# Dennard Scaling



| Parameter | Value | Scaled Value |
|---|---|---|
| Dopant concentrations | Na, Nd | Na/α, Nd/α |
| Dimensions | L, W, Tox | αL, αW, αTox |
| Field | E | E |
| Voltage | V | αV |
| Capacitance | C | αC |
| Current | I | αI |

| Transistors/Area | d | d/α² |
|---|---|---|

| Propagation time (~CV/I) | t | αt |
|---|---|---|
| Frequency (1/t) | f | f/α |

*Scale factor* $\alpha < 1$

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

Design of Ion-Implanted MOSFET's with Very Small Physical Dimensions. Proc. of IEEE. Dennart et al. 1974.

# Dennard Scaling



| Parameter | Value | Scaled Value |
|---|---|---|
| Dopant concentrations | Na, Nd | Na/α, Nd/α |
| Dimensions | L, W, Tox | αL, αW, αTox |
| Field | E | E |
| Voltage | V | αV |
| Capacitance | C | αC |
| Current | I | αI |

| | | |
|---|---|---|
| Transistors/Area | d | d/α$^2$ |

| | | |
|---|---|---|
| Propagation time (~CV/I) | t | αt |
| Frequency (1/t) | f | f/α |

| | | |
|---|---|---|
| Power (CV$^2$f) | P | α$^2$P |
| **Power/area (Power density)** | **Pd** | **Pd** |

*Scale factor α<1*

α = 0.7 => 2X more transistors!

Design of Ion-Implanted MOSFET's with Very Small Physical Dimensions. Proc. of IEEE. Dennart et al. 1974.

# Implications of Dennard Scaling and Moore's Law

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

# Implications of Dennard Scaling and Moore's Law

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

- More transistors means better microarchitecture, which leads to better performance even under the same frequency.

# Implications of Dennard Scaling and Moore's Law

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

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

- Each new processor generation can have more transistors (Moore's Law), and will run at higher frequency but won't consume more power under 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;
>
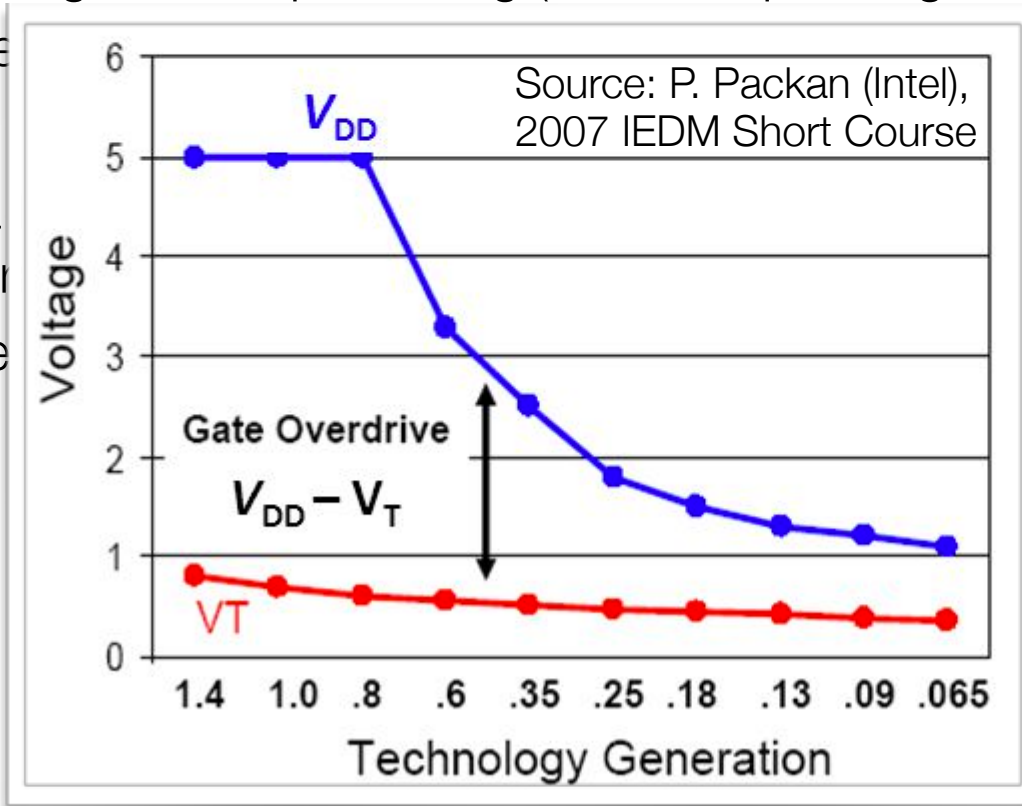> Dennard scaling made them useful.
>
> Bob Colwell, DAC 2013, June 4, 2013
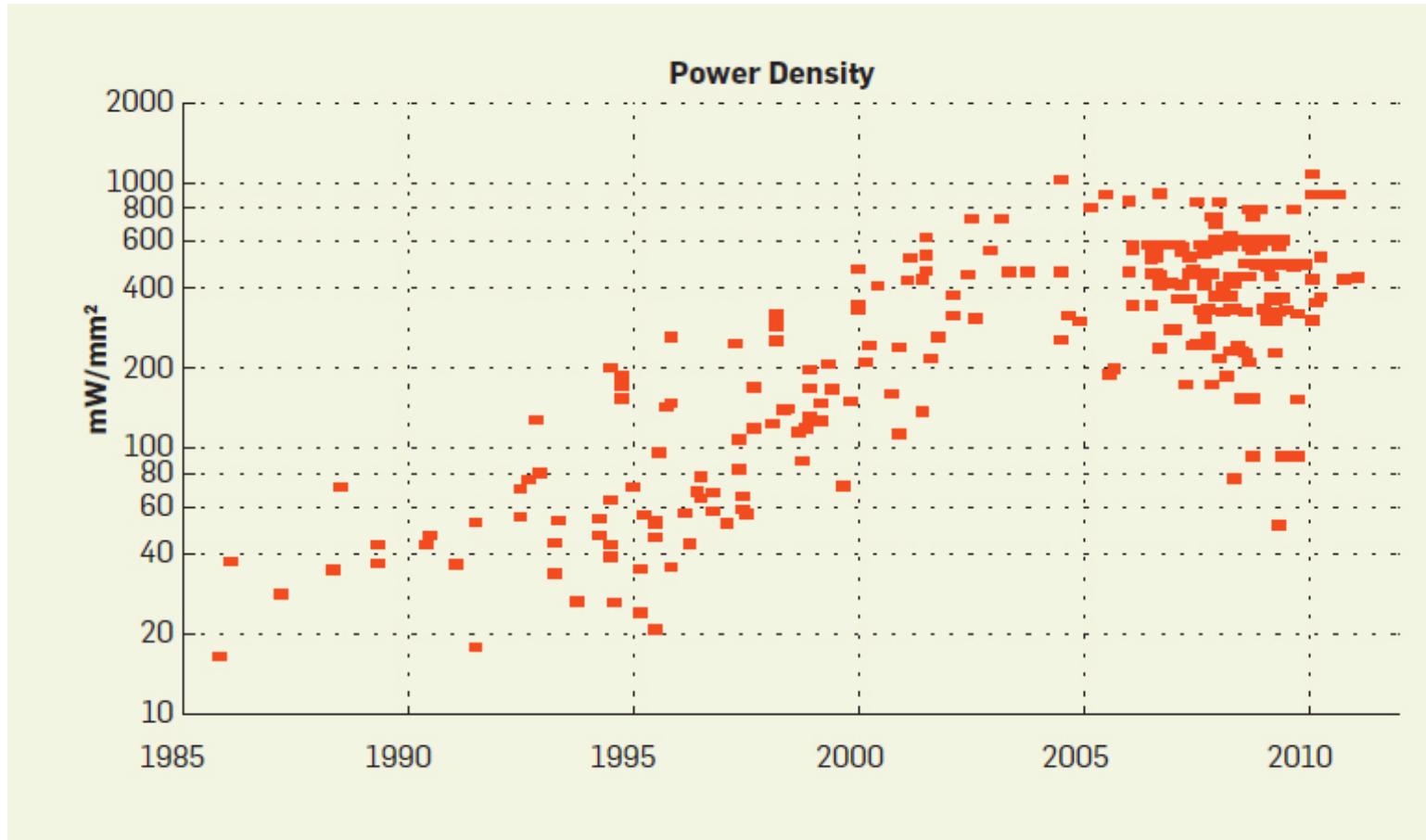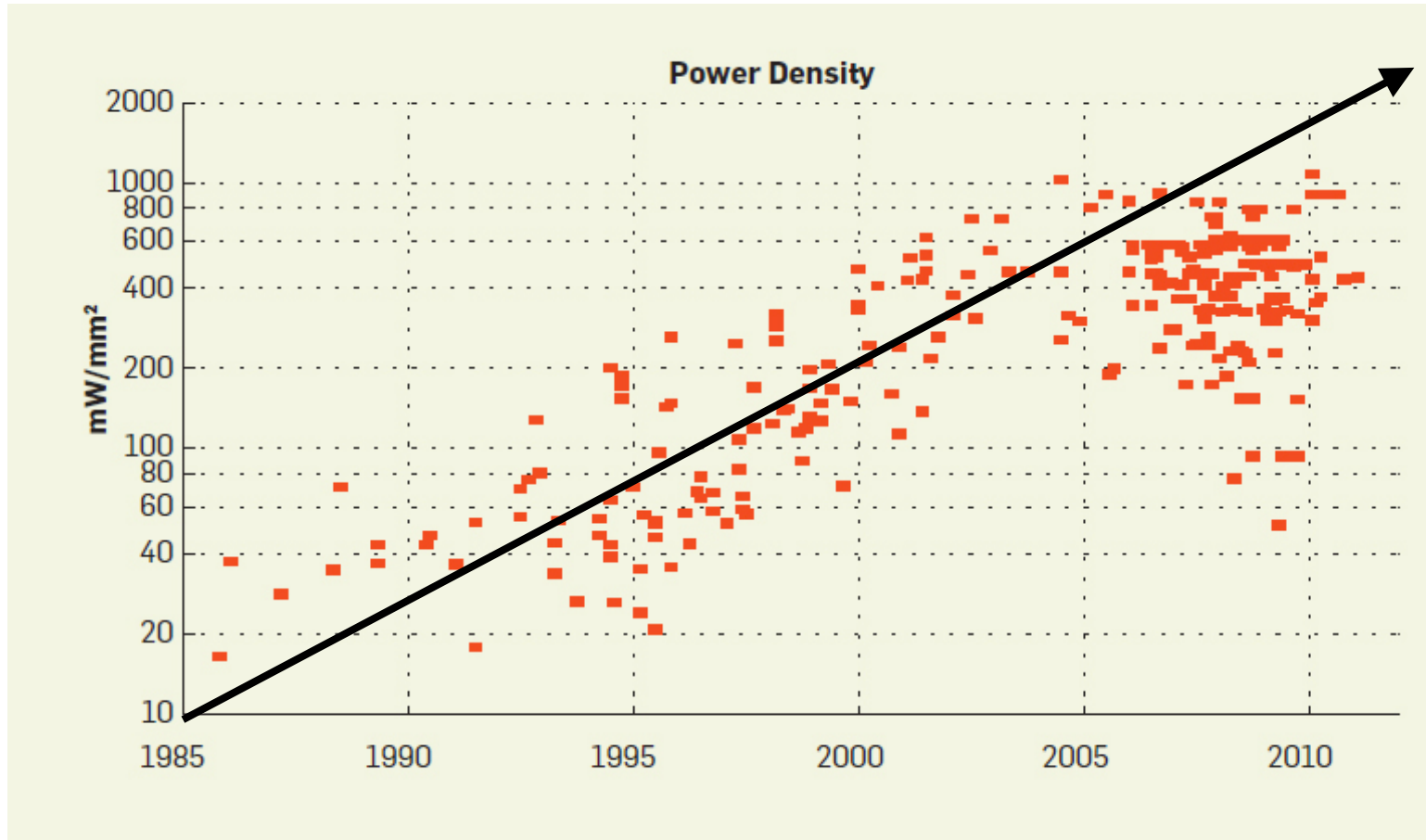
# 2005: End of Dennard Scaling

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

# 2005: End of Dennard Scaling

- What Happened?
  - Supply voltage $V_{dd}$ stops scaling (Can't drop voltage below ~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 ~1 V)
  - Remembe
- Why?
  - There is a
    switch a tr
  - $V_{th}$ stoppe
    becomes



Source: P. Packan (Intel), 2007 IEDM Short Course

# 2005: End of Dennard Scaling

- What Happened?
  - Supply voltage $V_{dd}$ stops scaling (Can't drop voltage below ~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
- 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



Power Density

# 2005: End of Dennard Scaling

# 2005: End of Dennard Scaling

# 2005: End of Dennard Scaling

# 2005: End of Dennard Scaling

# Dark Silicon
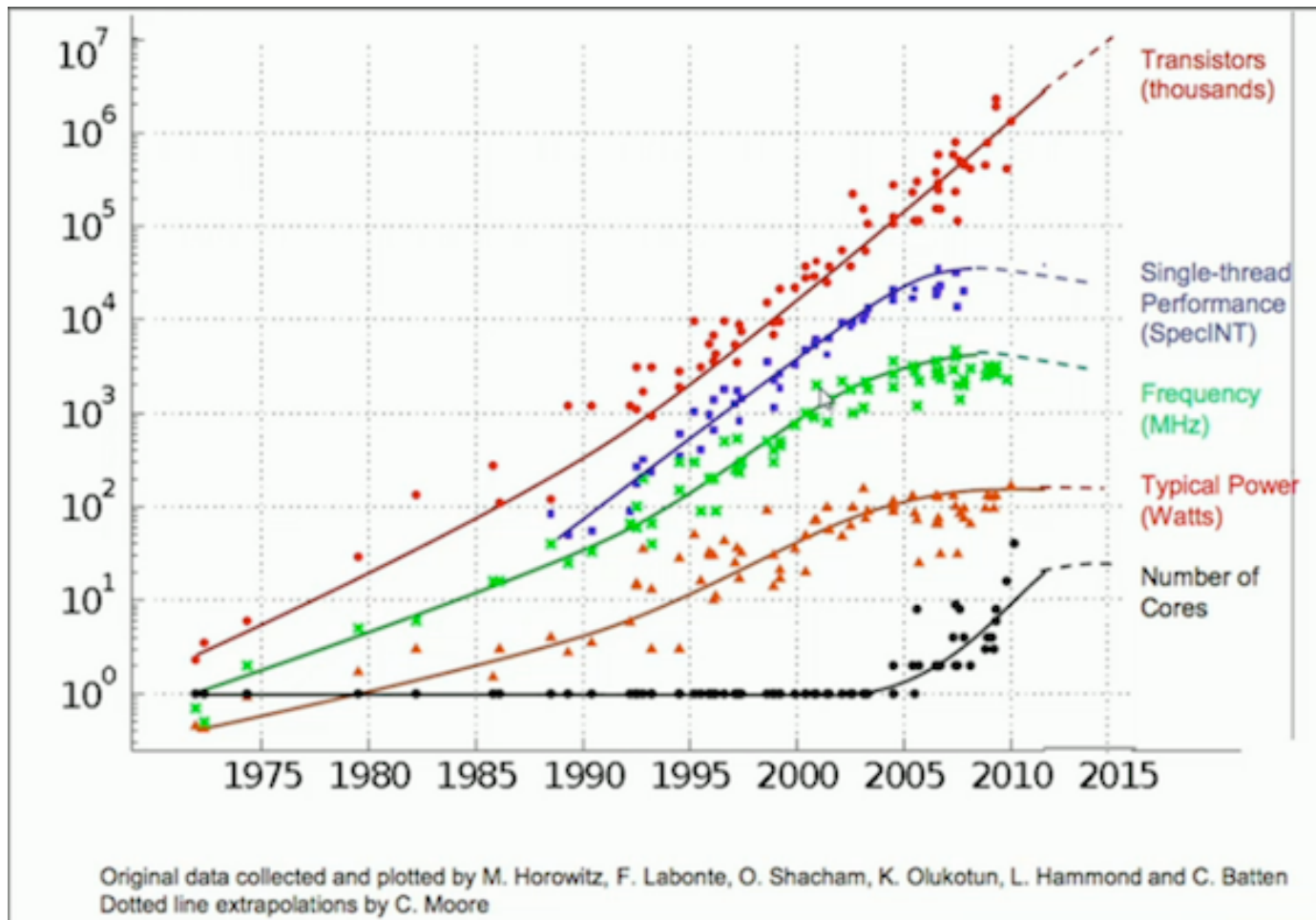
*n.* [därk, sĭl′ĭ-kən, -kŏn′]

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

# 2005: End of Dennard Scaling

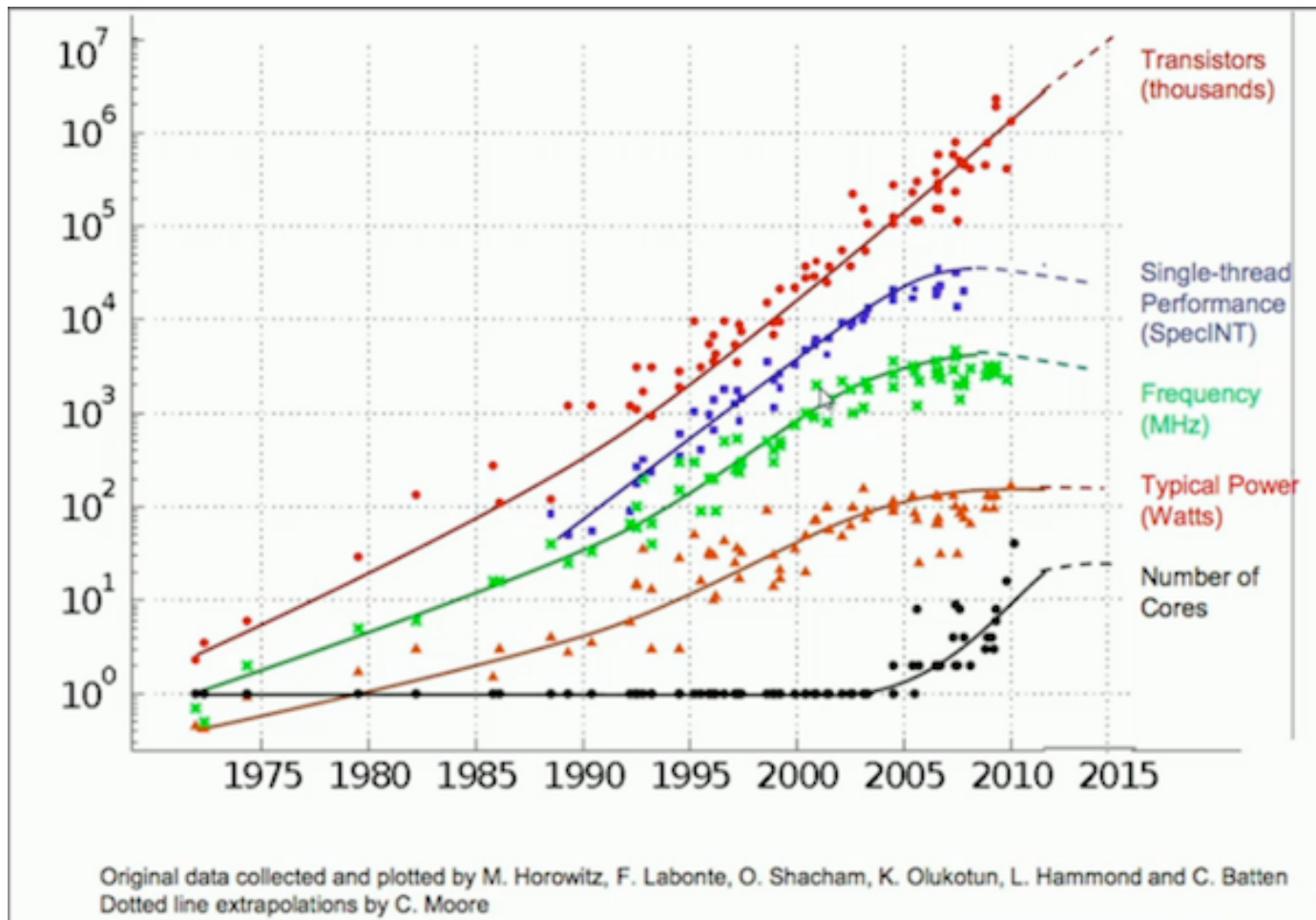- Initial response has been to lower frequency and increase cores / chip

# 2005: End of Dennard Scaling

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



Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten
Dotted line extrapolations by C. Moore

# 2005: End of Dennard Scaling

- Initial response has been to lower frequency and increase cores / chip
- There is a limit to core scaling. Why?



Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten
Dotted line extrapolations by C. Moore

# 2007: A Revolutionary New Computer

# No Moore's Law for batteries

**Fred Schlachter[1]**

*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-
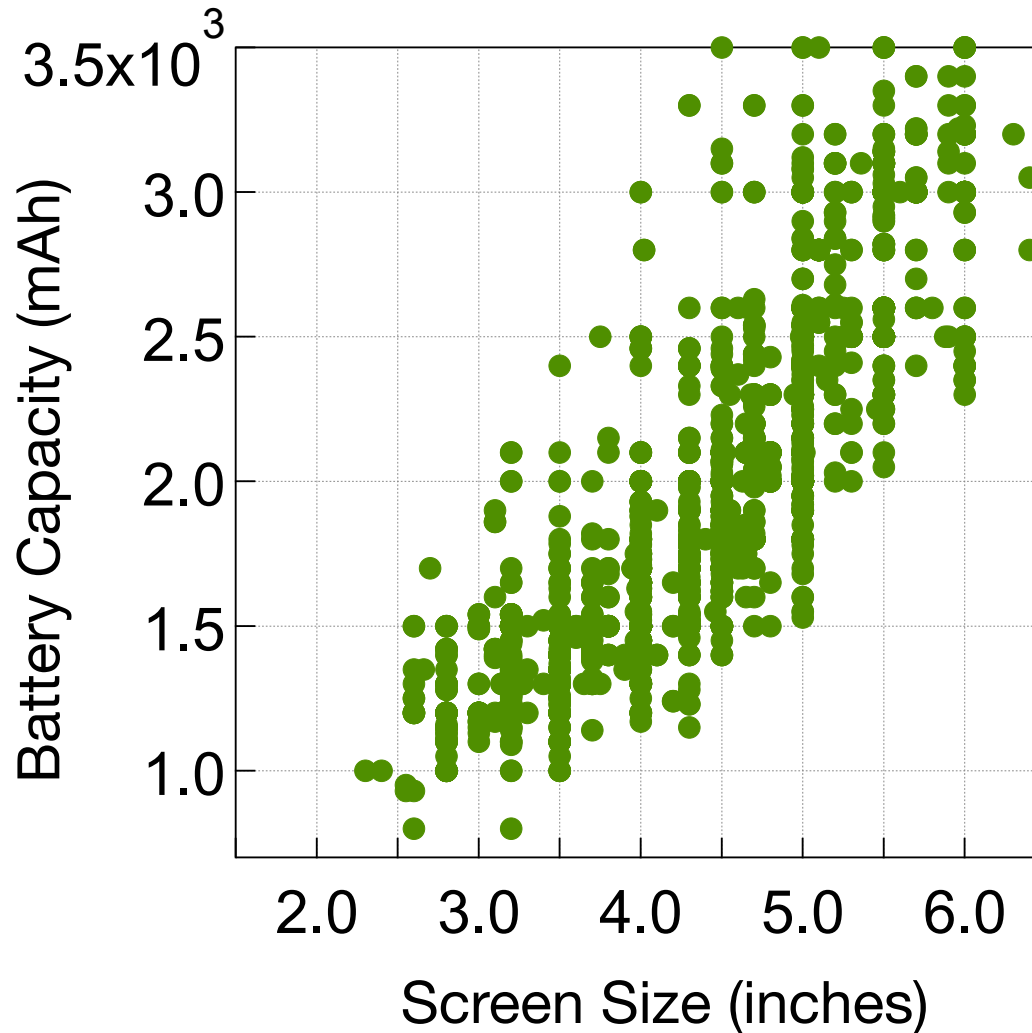
# "Improving" Energy Capacity

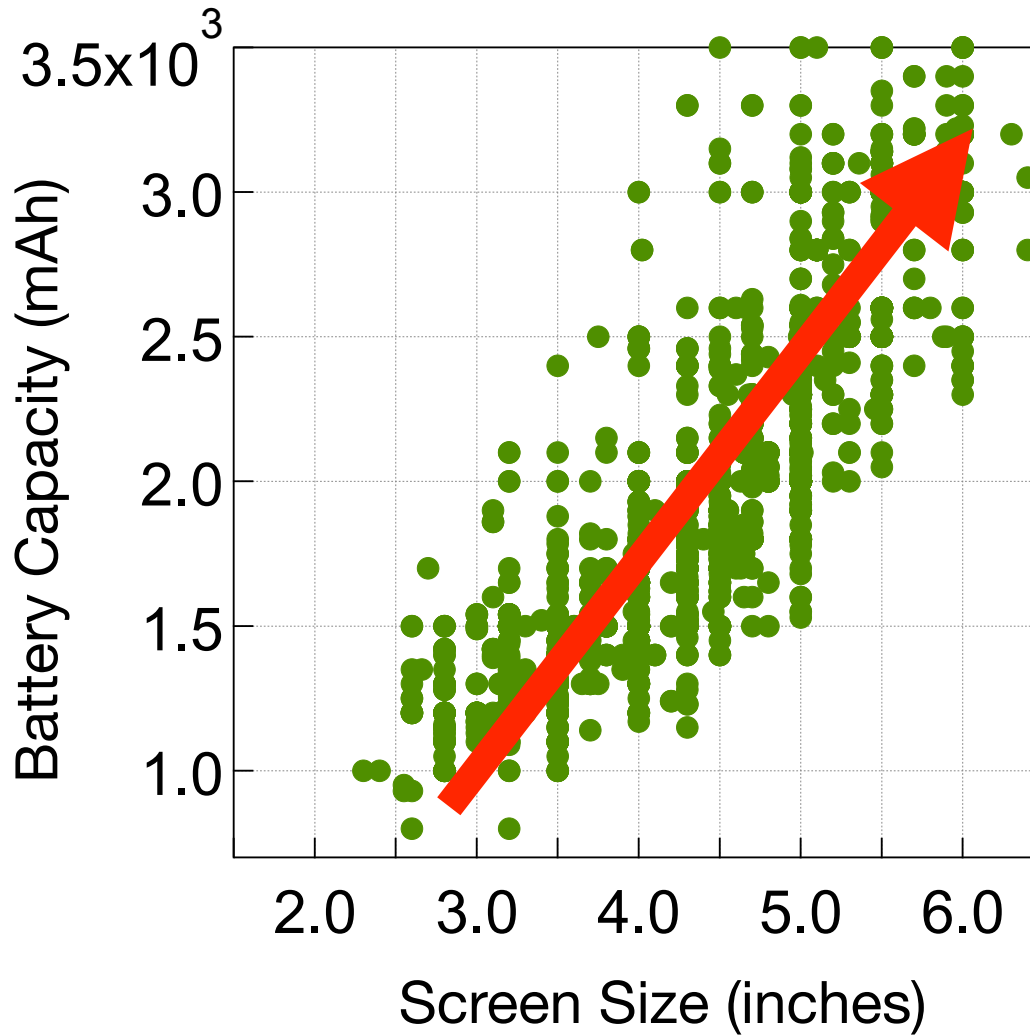600 smartphone from 2006 to 2014 on http://www.gsmarena.com/makers.php3

# "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



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

19

# "Improving" Energy Capacity

# "Improving" Energy Capacity



SMARTPHONE    PHABLET    TABLET

# "Improving" Energy Capacity



SMARTPHONE     PHABLET     TABLET

# "Improving" Energy Capacity



SMARTPHONE    PHABLET    TABLET

# "Improving" Energy Capacity

# Sources of Energy-Inefficiencies

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

Understanding sources of inefficiency in general-purpose chips, Hameed et al., ISCA 2010

# 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

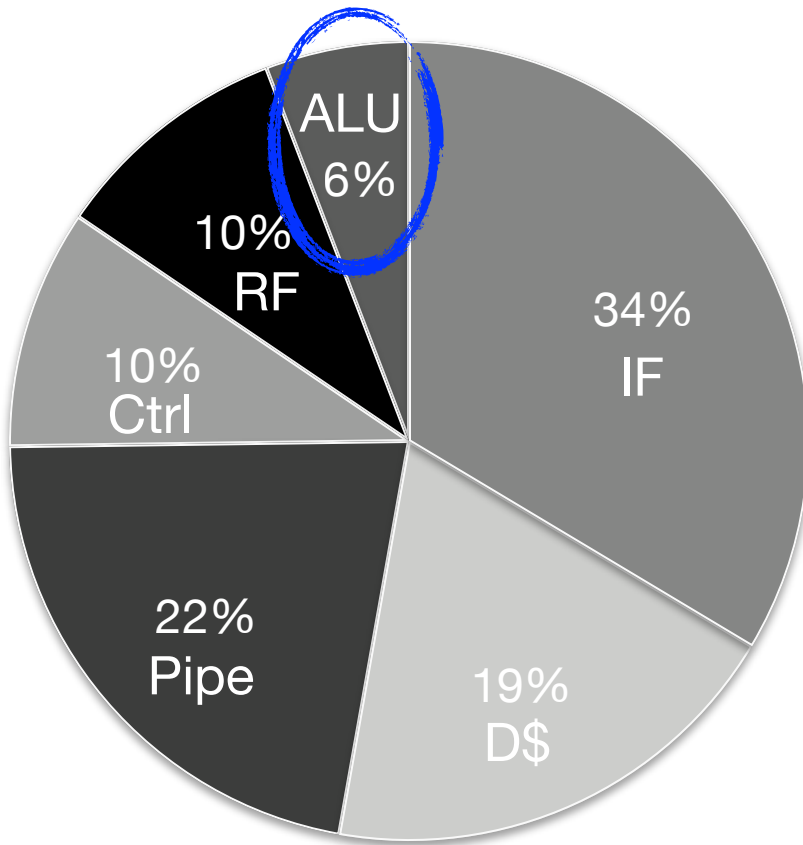**Ctrl**: Other control logics

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

**D$**: Data cache

**RF**: Register file

**ALU**: Functional units

Understanding sources of inefficiency in general-purpose chips, Hameed et al., ISCA 2010

# 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
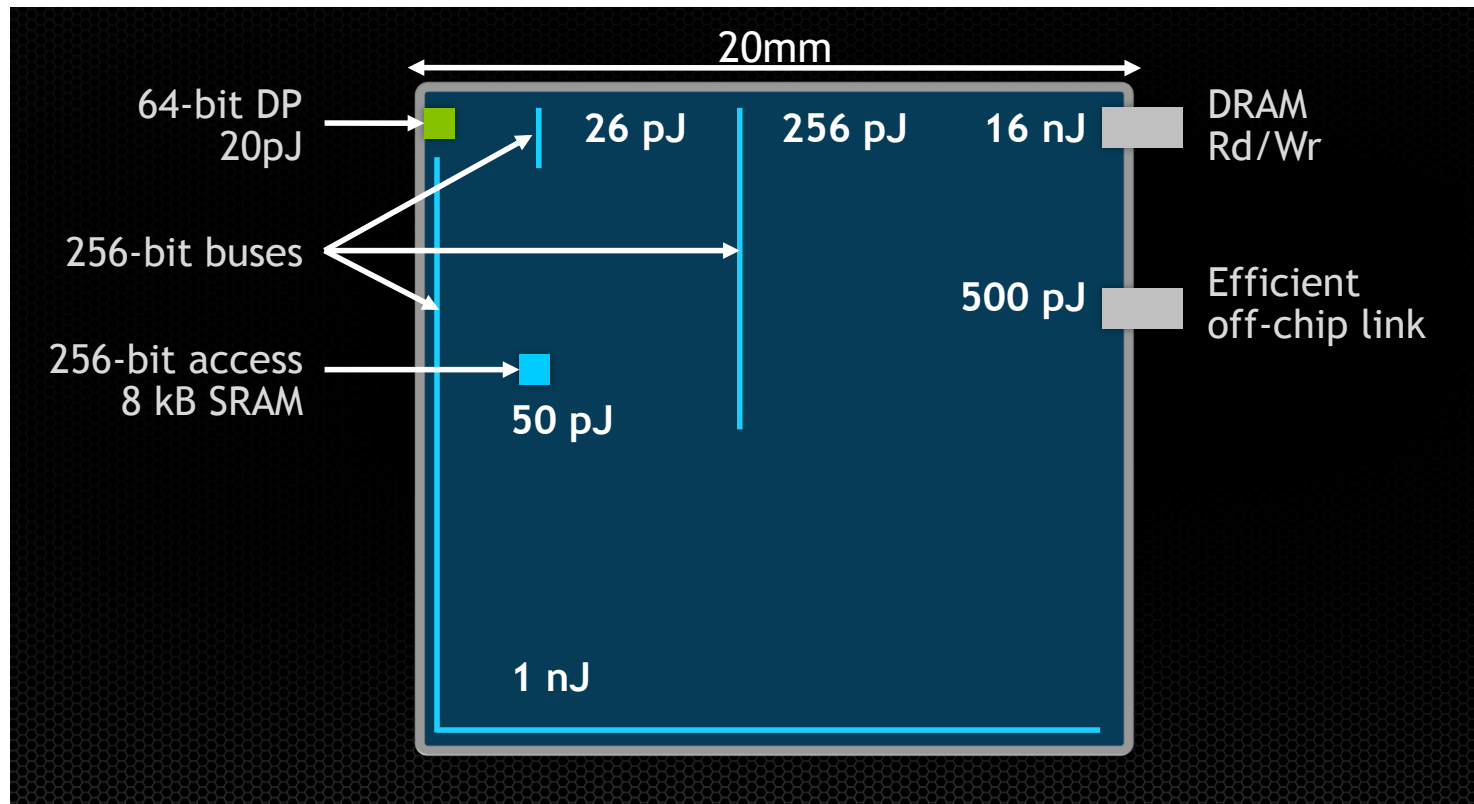
**D$**: Data cache

**RF**: Register file
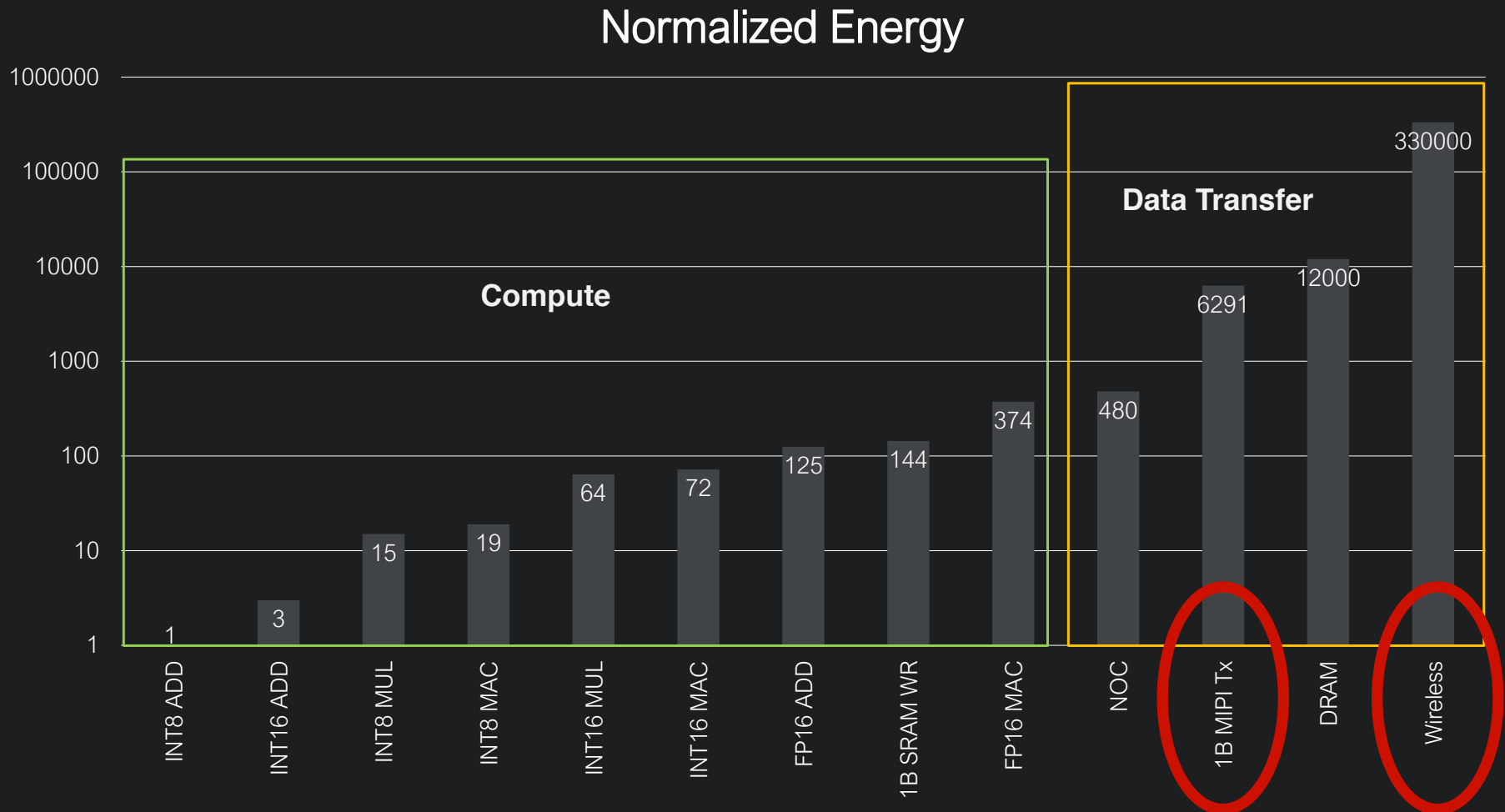
**ALU**: Functional units

Understanding sources of inefficiency in general-purpose chips, Hameed et al., ISCA 2010

# 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

Understanding sources of inefficiency in general-purpose chips, Hameed et al., ISCA 2010

# 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**

Understanding sources of inefficiency in general-purpose chips, Hameed et al., ISCA 2010

# Computation vs. Data Movement

Data movement energy >> computation energy



Challenges for future computing systems, Bill Dally, 2015

# Computation vs. Data Movement

Data movement energy >> computation energy



Normalized Energy

# SIMD

- Single Instruction (operating on) Multiple Data
- Amortizing the cost of instruction delivery/control across many execution units (even cores).
- Almost all modern ISAs provide such instructions:
  - x86: MMX/SSE/AVX
  - Arm: Neon



Pie chart labels: ALU 6%, 10% RF, 10% Ctrl, 34% IF, 19% D$, 22% Pipe

**Scalar Process**

b0 b1 b2 b3    c0 c1 c2 c3

X X X X

a0 a1 a2 a3

**4 instructions 4 elements**

**Vector Process (N=8)**

b0 b1 b2 b3 b4 b5 b6 b7    c0 c1 c2 c3 c4 c5 c6 c7

vmulps

a0 a1 a2 a3 a4 a5 a6 a7

**1 instruction 8 elements (AVX)**

# Graphics Processing Units/GPUs (SIMT)

- Designed for graphics rendering, which is massively parallel.



**Graphics rendering pipeline based on rasterization**



Same program for all vertices

Same program for all pixels

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007 ECE 498AL, University of Illinois, Urbana-Champaign

Computer Architecture

# Execute shader

```
<diffuseShader>:
sample r0, v4, t0, s0
mul  r3, v0, cb0[0]
madd r3, v1, cb0[1], r3
madd r3, v2, cb0[2], r3
clmp r3, r3, l(0.0), l(1.0)
mul  o0, r0, r3
mul  o1, r1, r3
mul  o2, r2, r3
mov  o3, l(1.0)
```

26

Kayvon Fatahalian

Kayvon Fatahalian, 2008

15

# Two cores   (two fragments in parallel)

fragment 1

fragment 2



```
<diffuseShader>:
sample r0, v4, t0, s0
mul  r3, v0, cb0[0]
madd r3, v1, cb0[1], r3
madd r3, v2, cb0[2], r3
clmp r3, r3, l(0.0), l(1.0)
mul  o0, r0, r3
mul  o1, r1, r3
mul  o2, r2, r3
mov  o3, l(1.0)
```

```
<diffuseShader>:
sample r0, v4, t0, s0
mul  r3, v0, cb0[0]
madd r3, v1, cb0[1], r3
madd r3, v2, cb0[2], r3
clmp r3, r3, l(0.0), l(1.0)
mul  o0, r0, r3
mul  o1, r1, r3
mul  o2, r2, r3
mov  o3, l(1.0)
```

**Fetch/Decode**

**ALU** (Execute)

**Execution Context**

**Fetch/Decode**

**ALU** (Execute)

**Execution Context**

27

# Four cores (four fragments in parallel)

Kayvon Fatahalian

EE382N: Principles of Computer Architecture

Kayvon Fatahalian, 2008

# Sixteen cores (sixteen fragments in parallel)



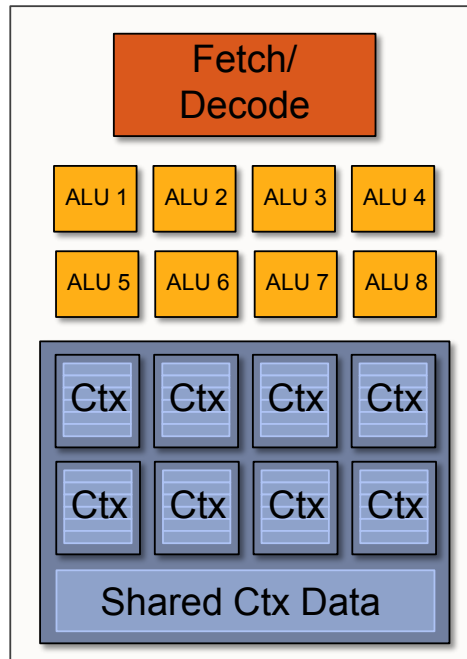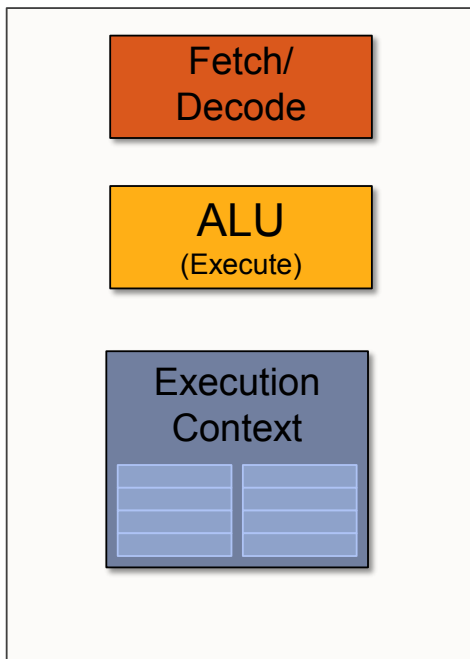16 cores = 16 simultaneous instruction streams

# Instruction stream coherence



But… many fragments should be able to share an instruction stream!

```
<diffuseShader>:
sample r0, v4, t0, s0
mul  r3, v0, cb0[0]
madd r3, v1, cb0[1], r3
madd r3, v2, cb0[2], r3
clmp r3, r3, l(0.0), l(1.0)
mul  o0, r0, r3
mul  o1, r1, r3
mul  o2, r2, r3
mov  o3, l(1.0)
```
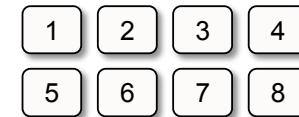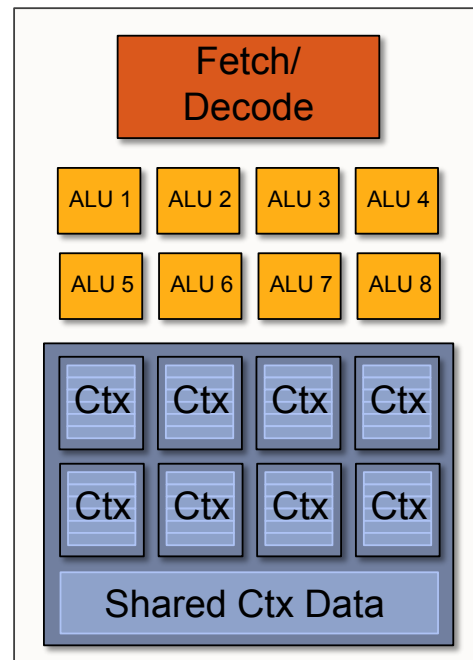
30

Fetch/Decode

ALU
(Execute)

Execution Context

Fetch/Decode

| ALU 1 | ALU 2 | ALU 3 | ALU 4 |
| ALU 5 | ALU 6 | ALU 7 | ALU 8 |

| Ctx | Ctx | Ctx | Ctx |
| Ctx | Ctx | Ctx | Ctx |

Shared Ctx Data

Amortize cost/complexity of managing an instruction stream across many ALUs
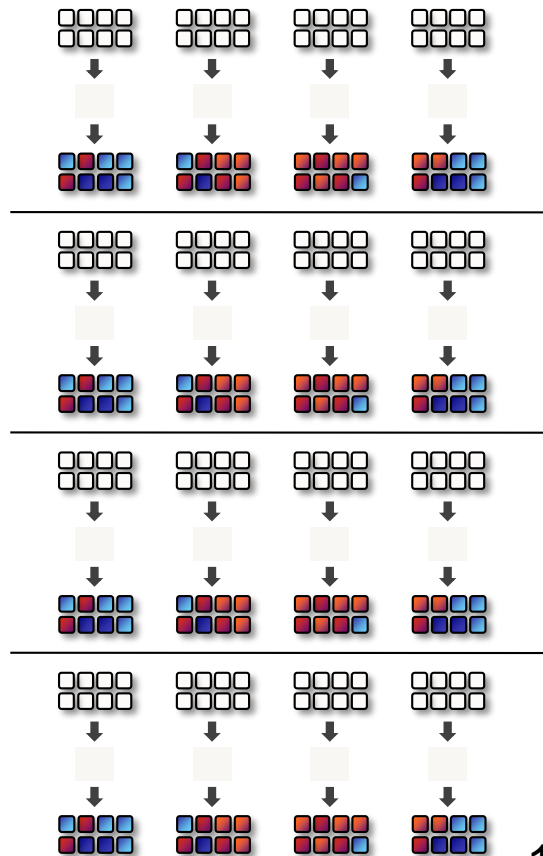
# SIMD processing

SIMD/vector instructions, each operates on a vector of 8 elements here.



| 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 |

```
VEC8_DiffuseShader>:
VEC8_sample vec_r0, vec_v4, t0, vec_s0
VEC8_mul   vec_r3, vec_v0, cb0[0]
VEC8_madd  vec_r3, vec_v1, cb0[1], vec_r3
VEC8_madd  vec_r3, vec_v2, cb0[2], vec_r3
VEC8_clmp  vec_r3, vec_r3, l(0.0), l(1.0)
VEC8_mul   vec_o0, vec_r0, vec_r3
VEC8_mul   vec_o1, vec_r1, vec_r3
VEC8_mul   vec_o2, vec_r2, vec_r3
VEC8_mov   vec_o3, l(1.0)
```
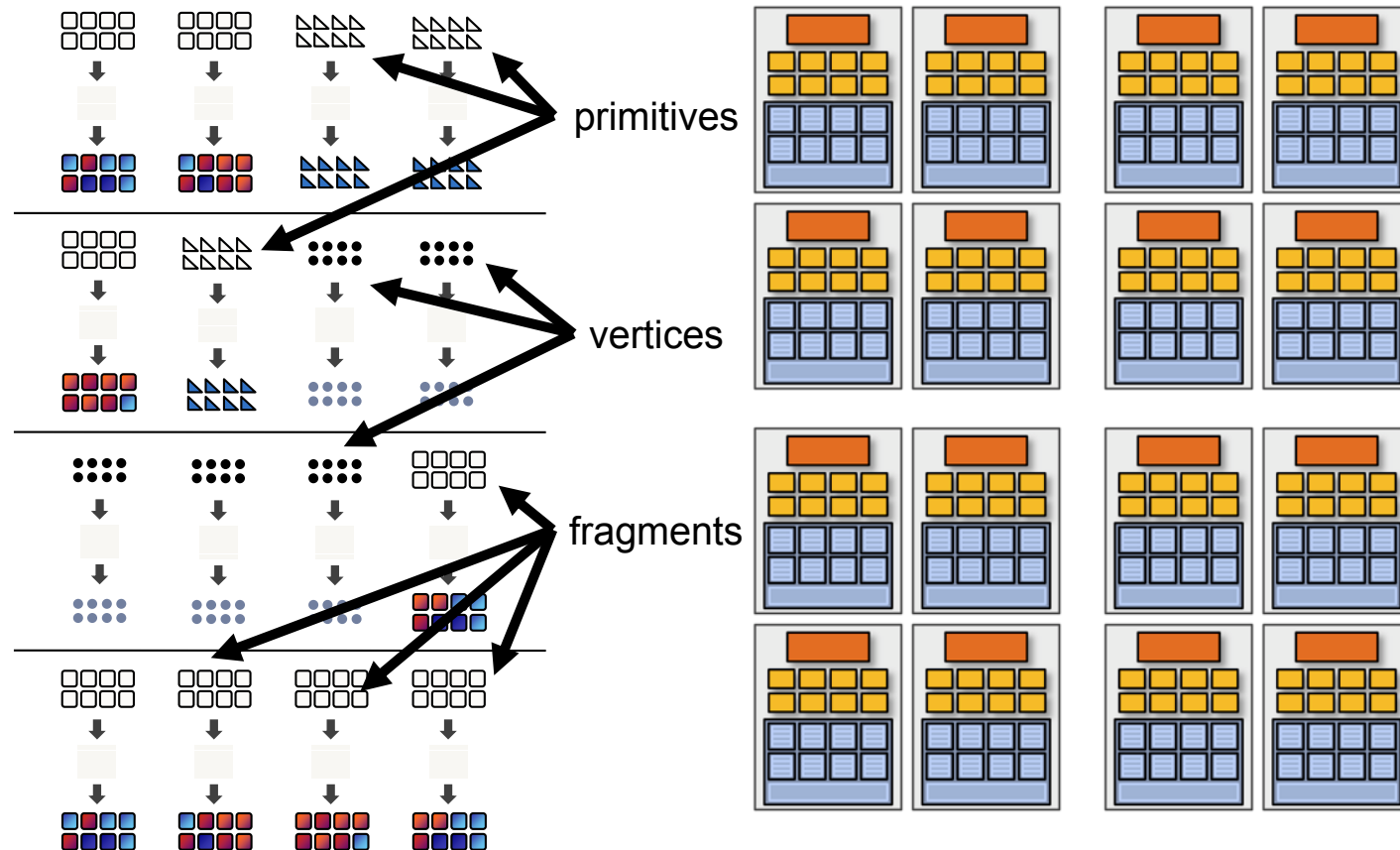
Fetch/
Decode

| ALU 1 | ALU 2 | ALU 3 | ALU 4 |
| ALU 5 | ALU 6 | ALU 7 | ALU 8 |

| Ctx | Ctx | Ctx | Ctx |
| Ctx | Ctx | Ctx | Ctx |

Shared Ctx Data

# 16 cores, each with 8 ALUs. Each core here runs the same program (fragment shader)



16 cores = 128 ALUs
= 16 simultaneous instruction streams

Kayvon Fatahalian

EE382N: Principles of Computer Architecture

Kayvon Fatahalian, 2008

# 16 cores, each with 8 ALUs. Cores here run different programs
(some are processing vertices, some are processing fragments)



primitives

vertices
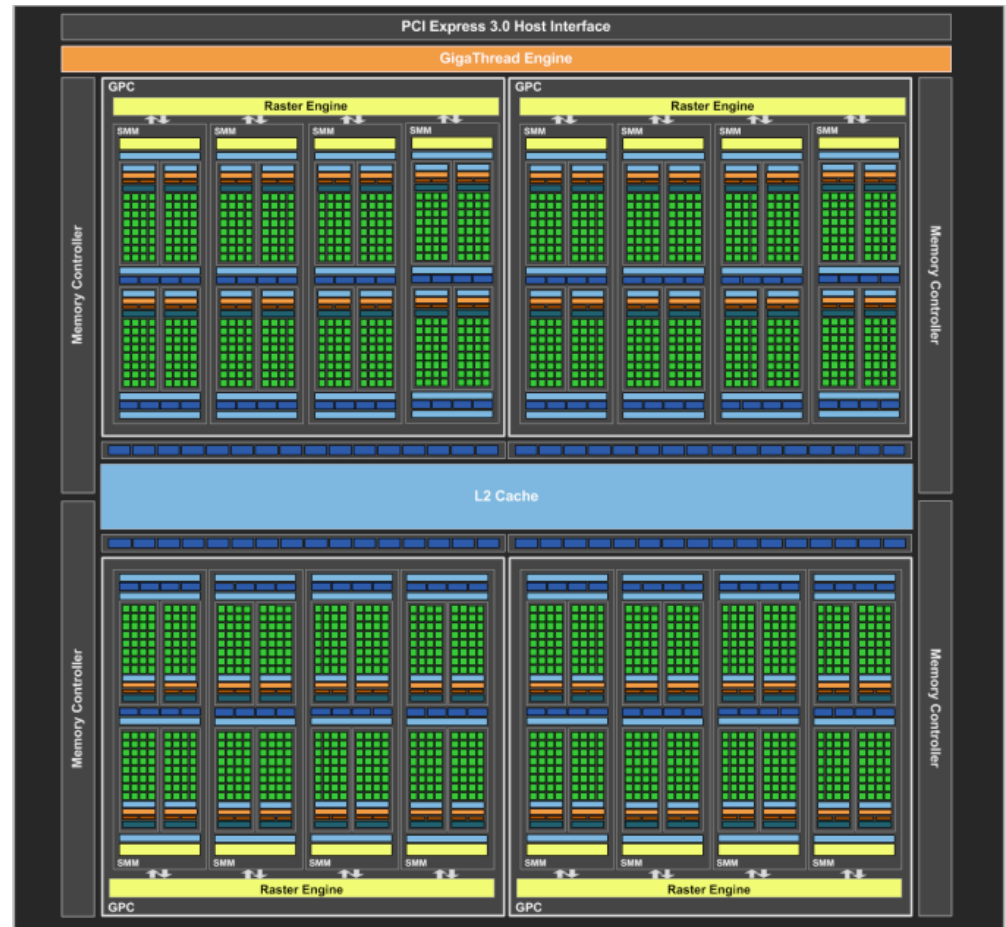
fragments

34

# Each Core Does Fine-Grained Multi-threading

Wrap: a group of threads (8 here)



Time

No need for branch prediction and out-of-order execution. Simple core design.

Warp 6, instruction 14

Warp 3, instruction 60

Warp 11, instruction 4

Warp 6, instruction 15

Warp 11, instruction 5

Warp 3, instruction 61

Kayvon Fatahalian, 2008

33

35

# Nvidia Maxwell GPU (2014)

- Today: General Purpose GPU (GPGPU), used for any massive parallel applications:
  - Physics simulation
  - Deep learning
  - Computer vision

# Nvidia Maxwell GPU (2014)

- Today: General Purpose GPU (GPGPU), used for any massive parallel applications:

**NVIDIA Corporation**

**$578.34** ↑ 35,164.63%  +576.70 MAX

After Hours: $577.00 (↓ -0.23%) -$1.34

Closed: May 5, 7:59:33 PM UTC-4 · USD · NASDAQ · Disclaimer

| 1D | 5D | 1M | 6M | YTD | 1Y | 5Y | MAX |

# Entering the Era of Specialization

# Entering the Era of Specialization

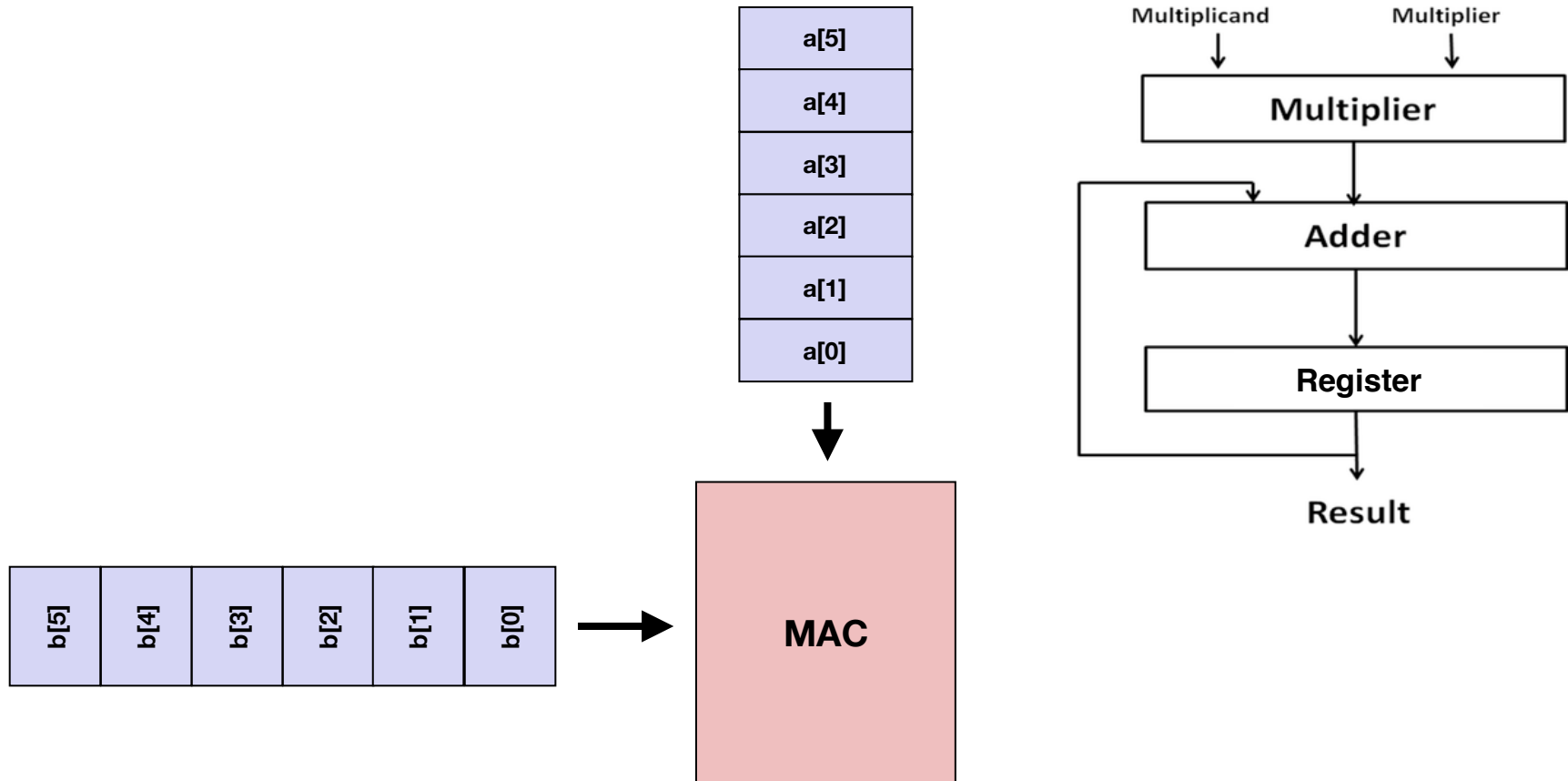- GPUs are very efficient for massively parallel program

# Entering the Era of Specialization

- GPUs are very efficient for massively parallel program
- But are still fairly general, so there are still many inefficiencies
  - Still need to fetch and decode instructions
  - Still have (very large) caches, so data delivery isn't efficient

# Entering the Era of Specialization

- GPUs are very efficient for massively parallel program
- But are still fairly general, so there are still many inefficiencies
  - Still need to fetch and decode instructions
  - Still have (very large) caches, so data delivery isn't efficient
- Idea: 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.
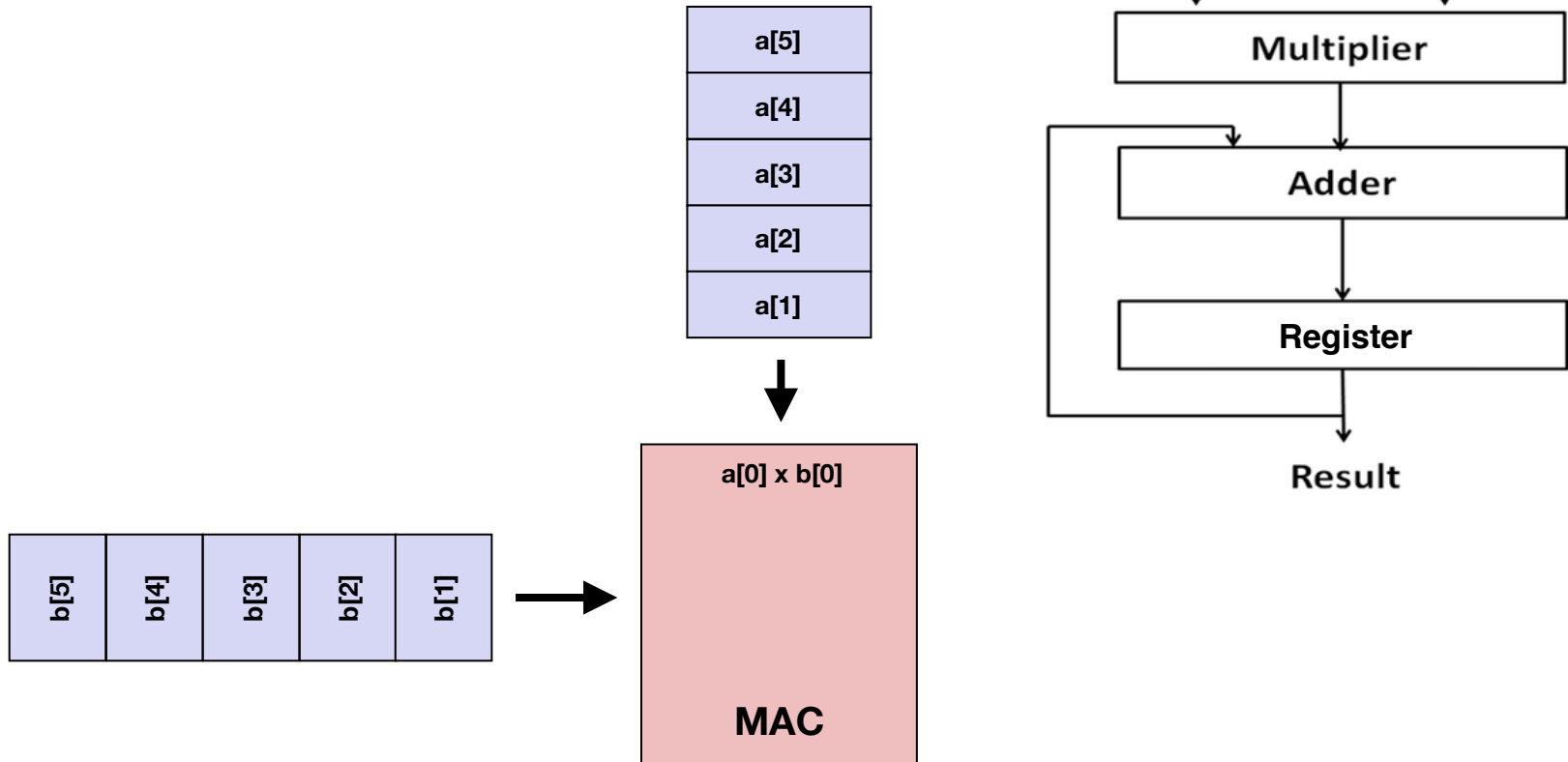
# Entering the Era of Specialization

- GPUs are very efficient for massively parallel program
- But are still fairly general, so there are still many inefficiencies
    - Still need to fetch and decode instructions
    - Still have (very large) caches, so data delivery isn't efficient
- Idea: 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.
- A.k.a., domain-specific accelerators

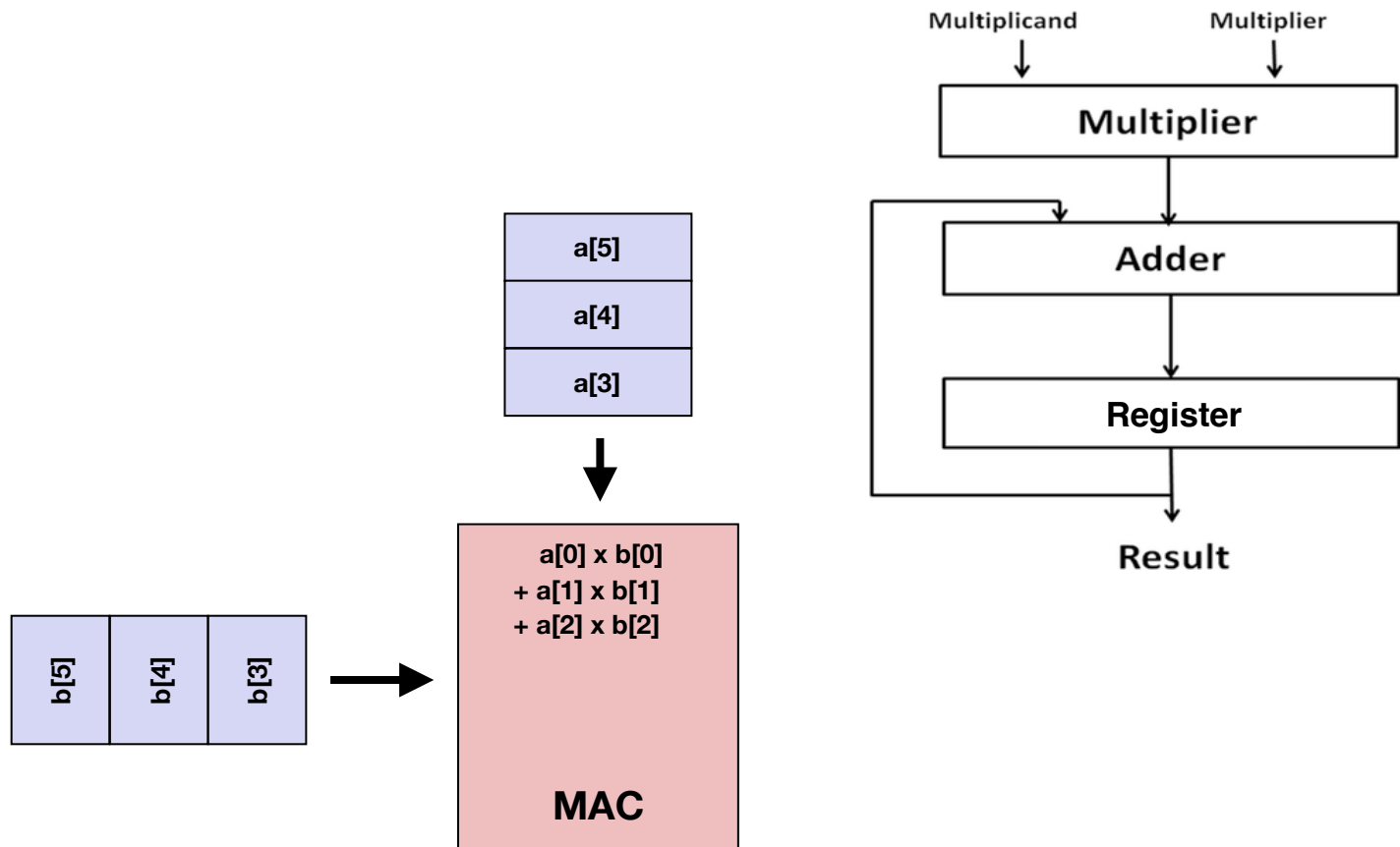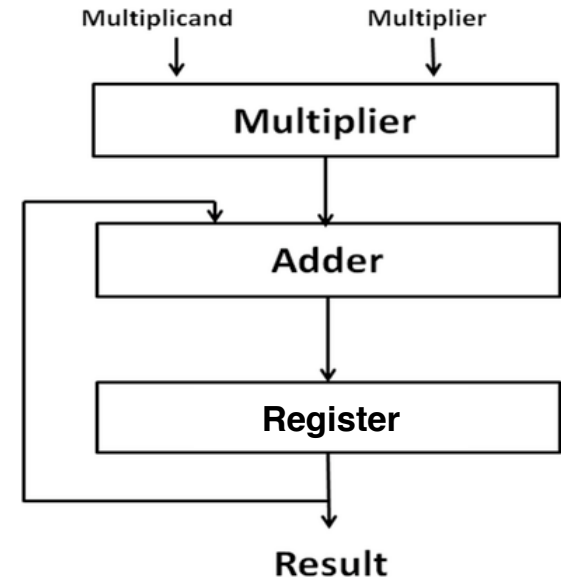# Example: Vector Dot Product

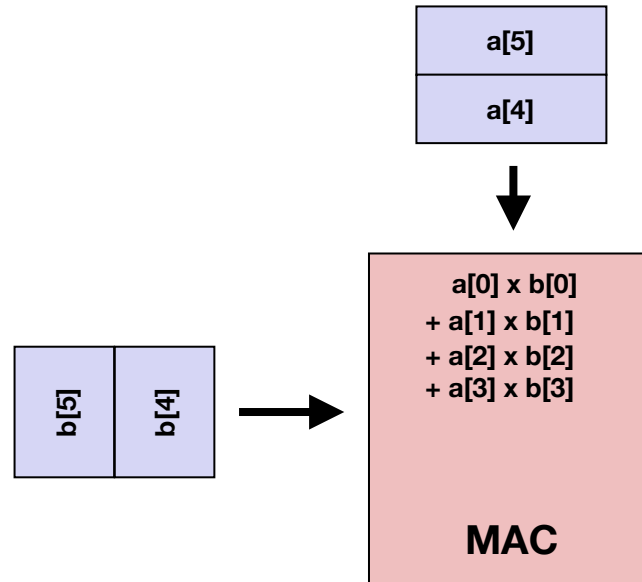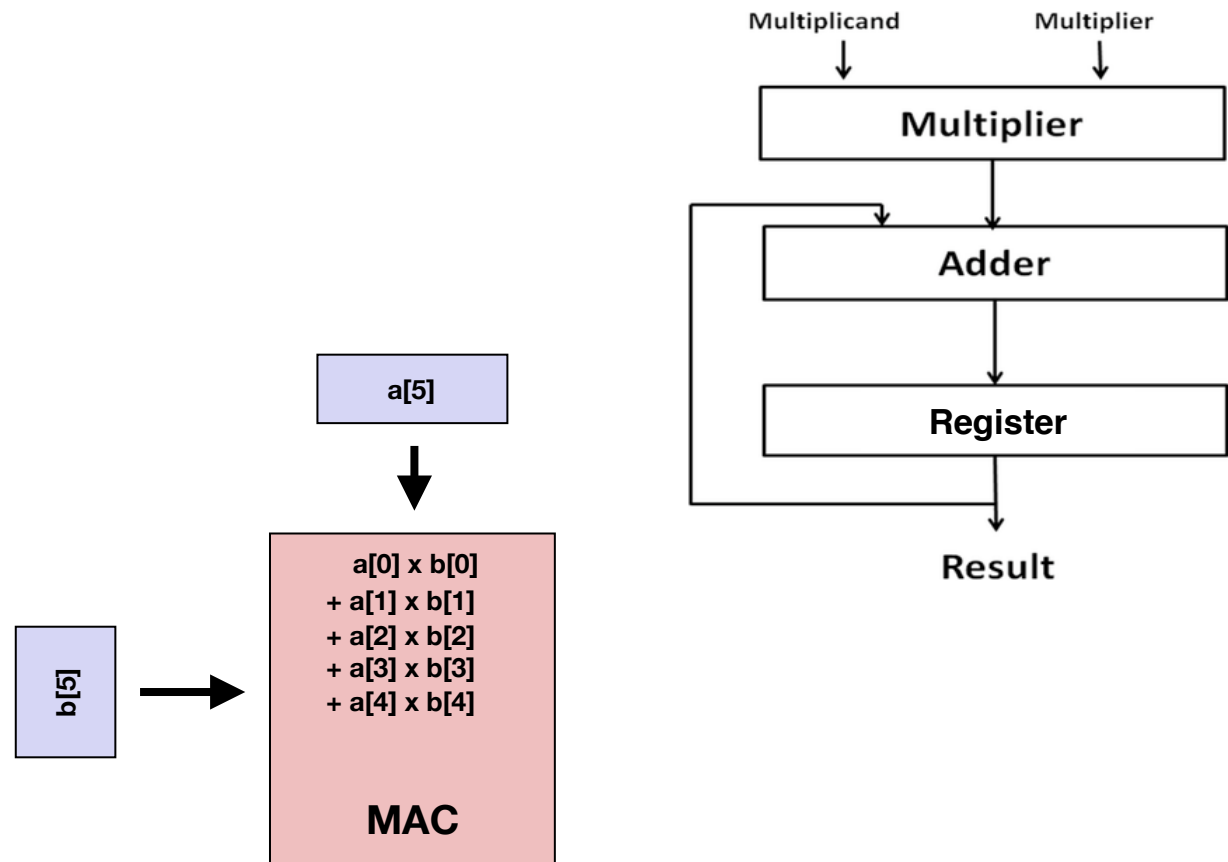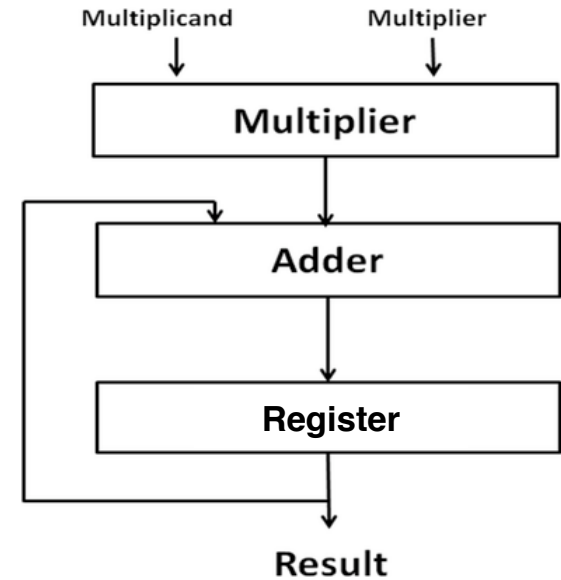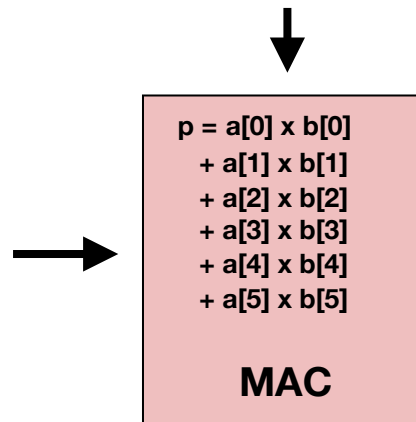# Example: Vector Dot Product

# Example: Vector Dot Product
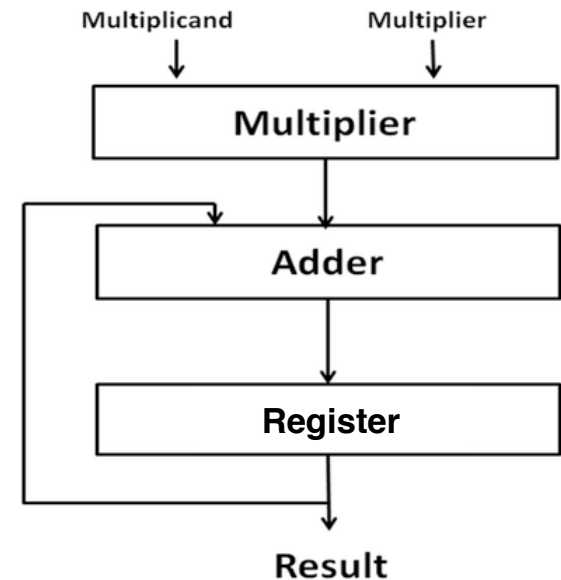
# Example: Vector Dot Product



a[5]

a[4]

a[3]

a[2]

a[0] x b[0]
+ a[1] x b[1]

b[5]  b[4]  b[3]  b[2]

**MAC**

Multiplicand          Multiplier

**Multiplier**

**Adder**

**Register**

Result

# Example: Vector Dot Product

a[5]

a[4]

a[3]

b[5]   b[4]   b[3]

a[0] x b[0]
+ a[1] x b[1]
+ a[2] x b[2]

**MAC**

Multiplicand          Multiplier

**Multiplier**

**Adder**

**Register**

Result

# Example: Vector Dot Product

Multiplicand      Multiplier

**Multiplier**

**Adder**

**Register**

Result

| a[5] |
| a[4] |

| b[5] | b[4] |

a[0] x b[0]
+ a[1] x b[1]
+ a[2] x b[2]
+ a[3] x b[3]

**MAC**

# Example: Vector Dot Product

Multiplicand          Multiplier

**Multiplier**

**Adder**

**Register**

Result

**a[5]**

**b[5]**

a[0] x b[0]
+ a[1] x b[1]
+ a[2] x b[2]
+ a[3] x b[3]
+ a[4] x b[4]

**MAC**

# Example: Vector Dot Product



Multiplicand          Multiplier

**Multiplier**

**Adder**

**Register**

Result

p = a[0] x b[0]
  + a[1] x b[1]
  + a[2] x b[2]
  + a[3] x b[3]
  + a[4] x b[4]
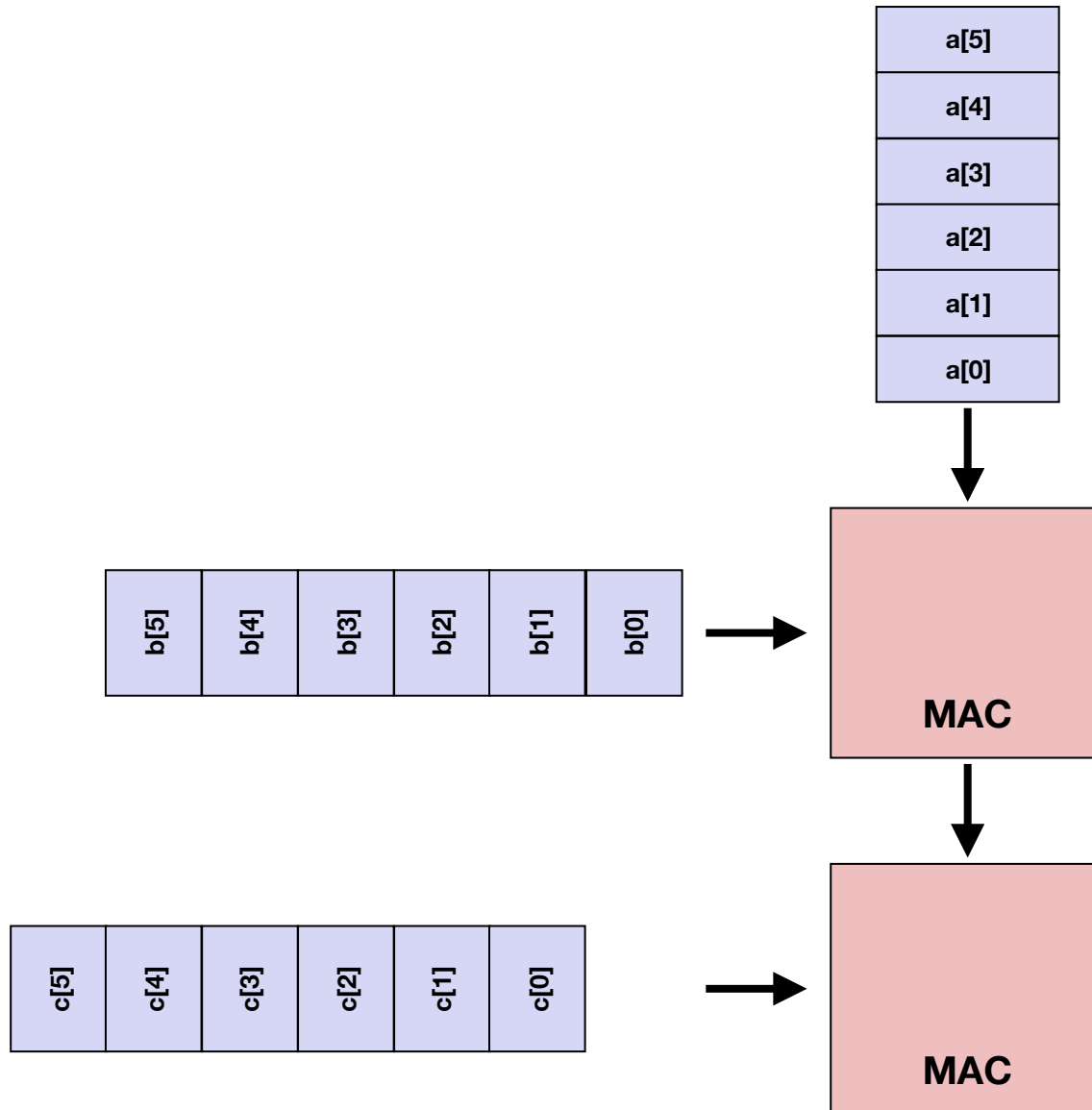  + a[5] x b[5]

**MAC**

# Example: Vector Dot Product

- Does nothing but vector dot product
  - No instruction fetch and decode (there is no instruction)
  - The register is close to the ALU and gets reused over and over: good data delivery efficiency
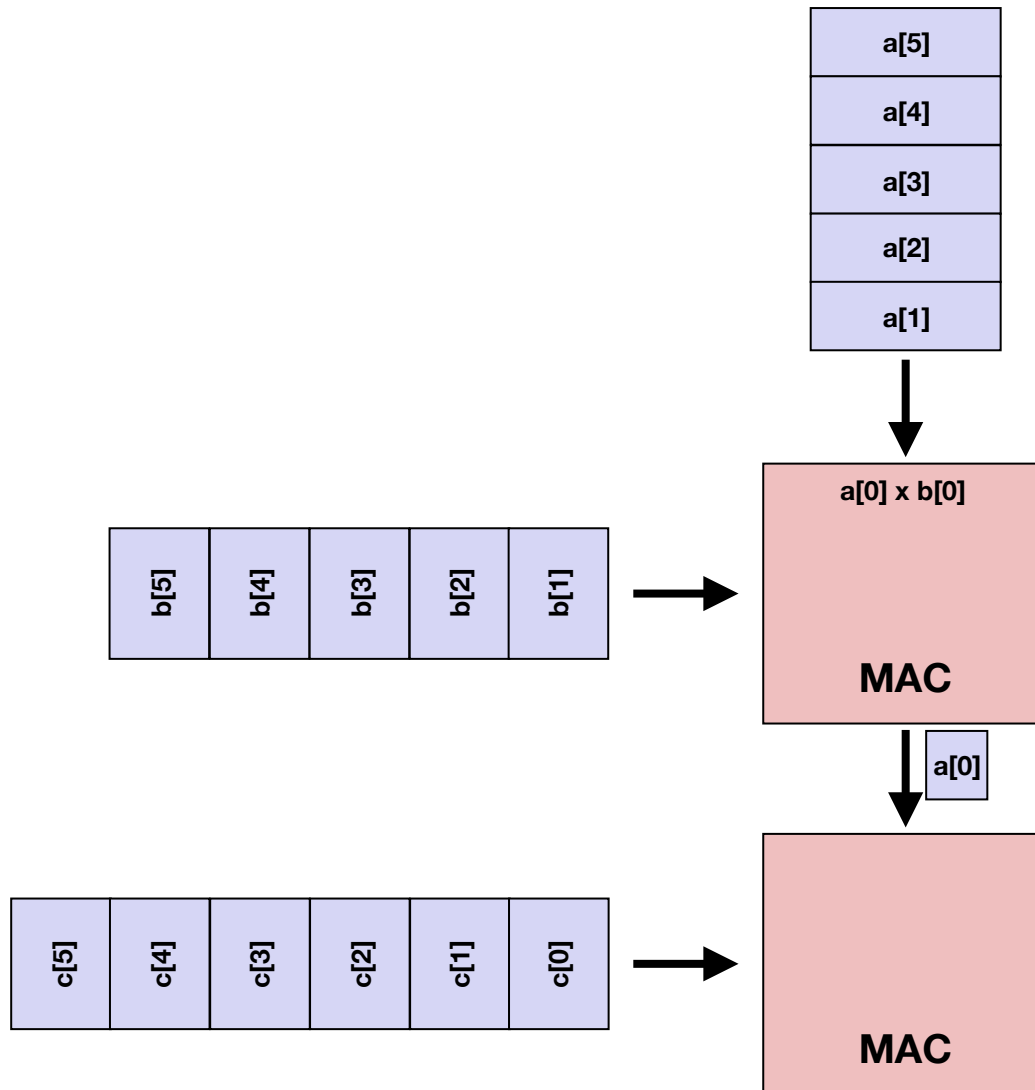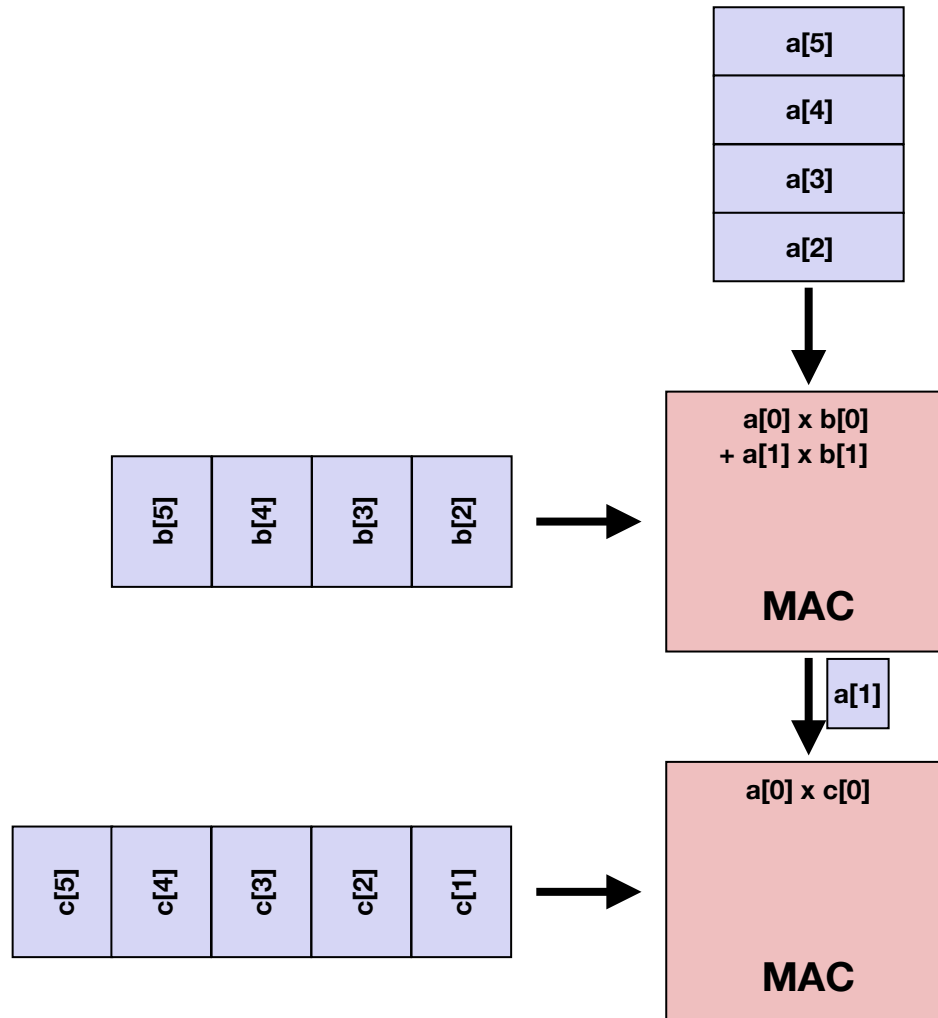  - Very simple control



```
p = a[0] x b[0]
  + a[1] x b[1]
  + a[2] x b[2]
  + a[3] x b[3]
  + a[4] x b[4]
  + a[5] x b[5]

    MAC
```

# Matrix Vector Multiplication

# Matrix Vector Multiplication

| a[5] |
|---|
| a[4] |
| a[3] |
| a[2] |
| a[1] |
| a[0] |

| b[5] | b[4] | b[3] | b[2] | b[1] | b[0] |
|---|---|---|---|---|---|

**MAC**

| c[5] | c[4] | c[3] | c[2] | c[1] | c[0] |
|---|---|---|---|---|---|

**MAC**

# Matrix Vector Multiplication

a[5]

a[4]

a[3]

a[2]

a[1]

a[0] x b[0]

b[5] b[4] b[3] b[2] b[1]

**MAC**

a[0]

c[5] c[4] c[3] c[2] c[1] c[0]

**MAC**

# Matrix Vector Multiplication

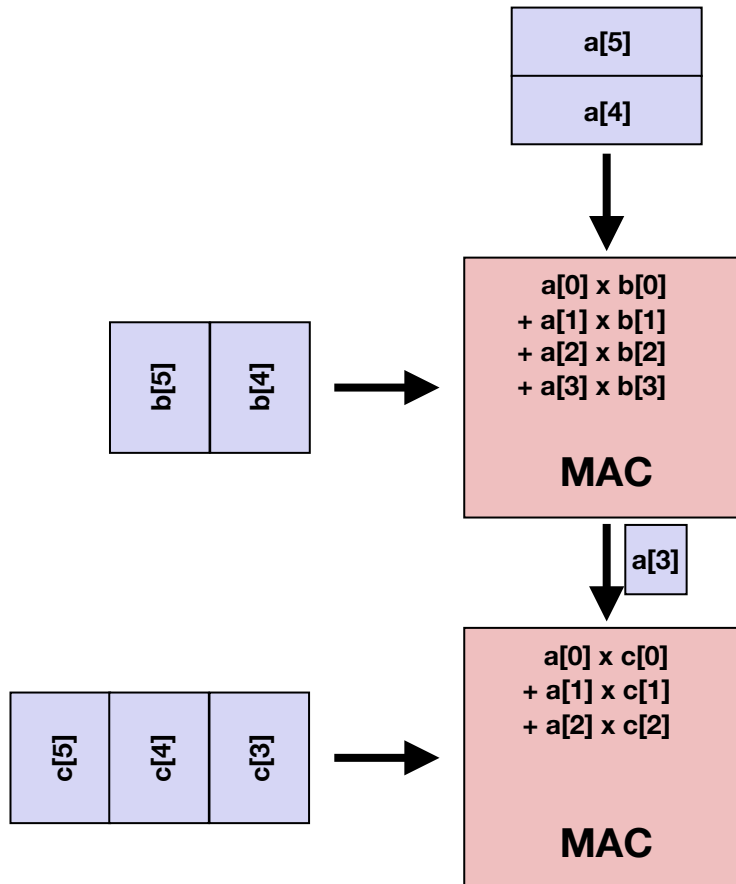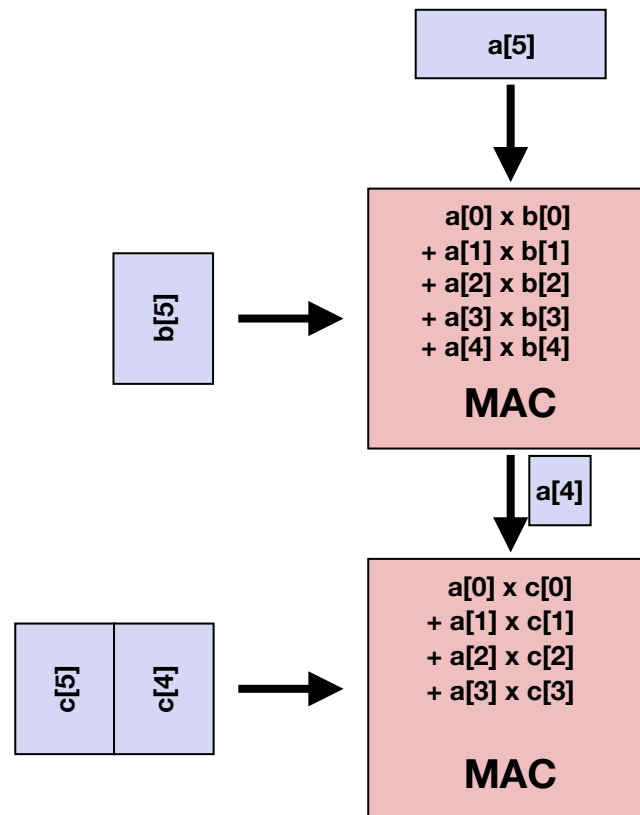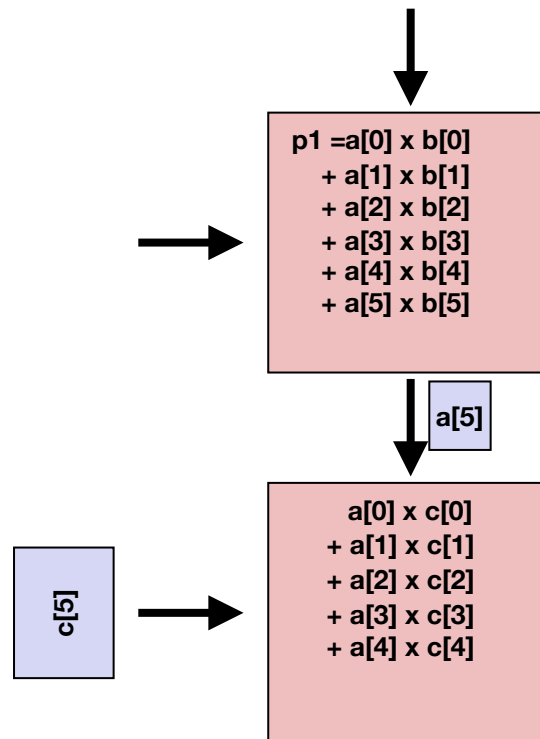# Matrix Vector Multiplication
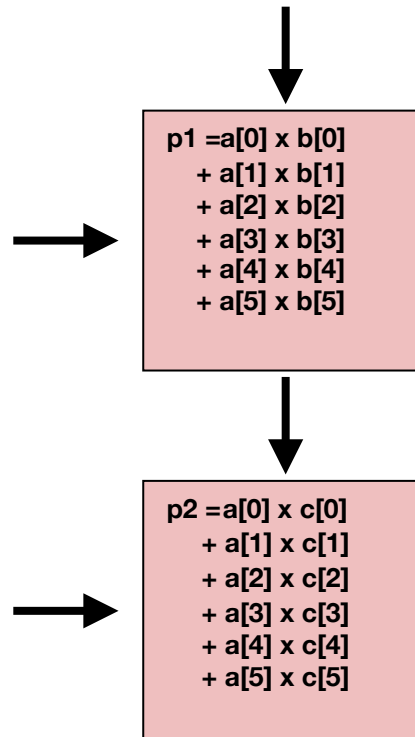
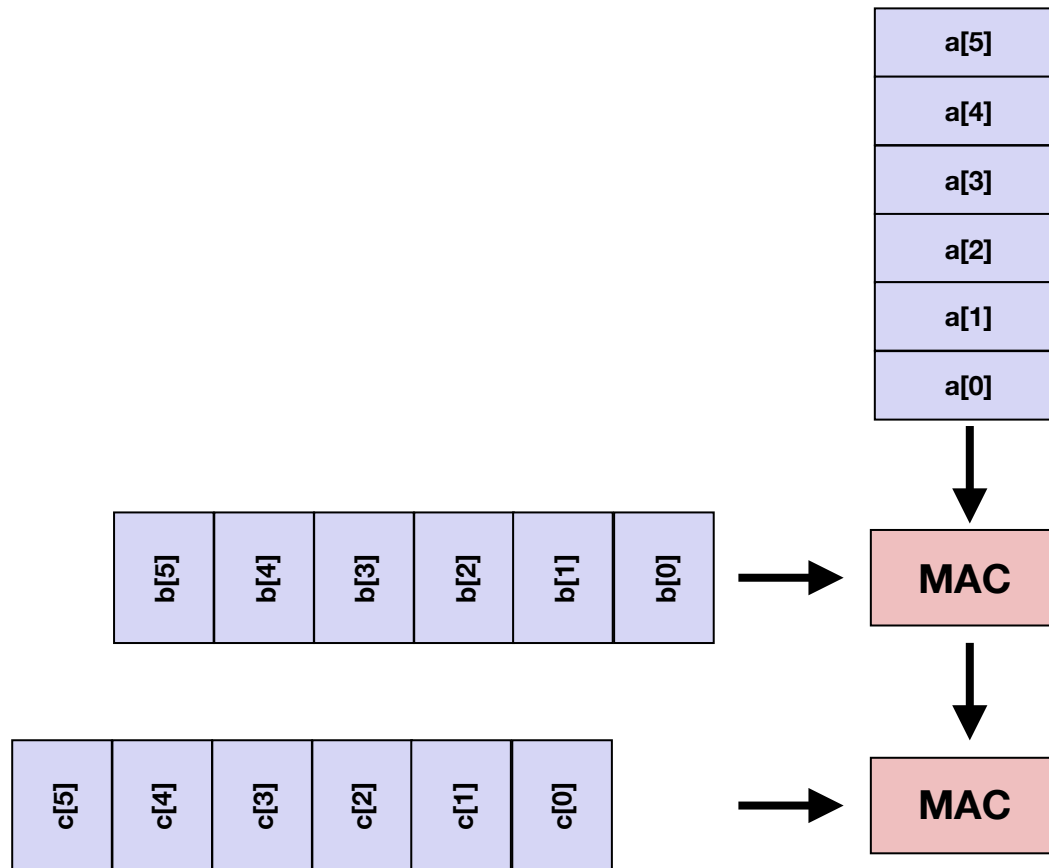# Matrix Vector Multiplication
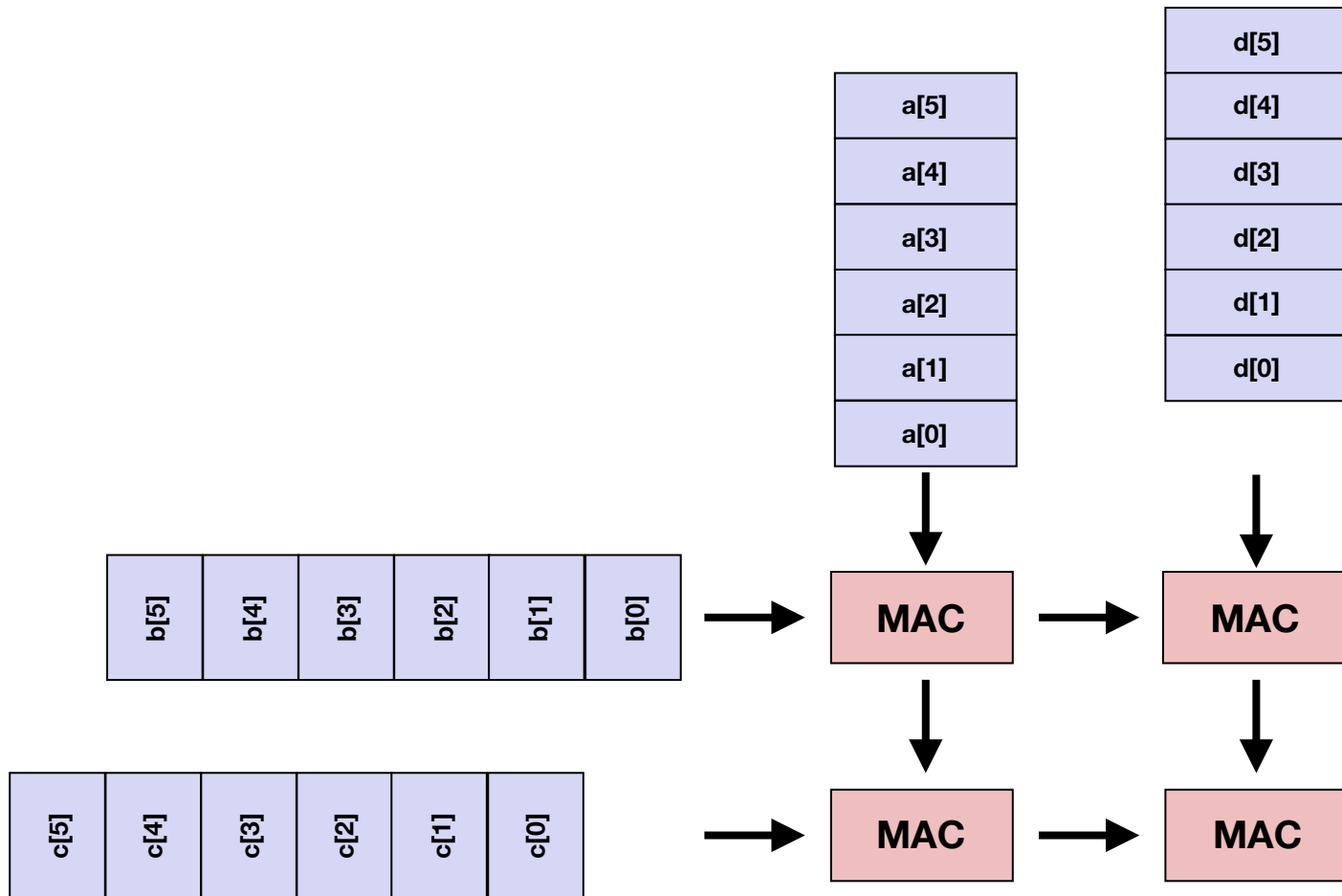
# Matrix Vector Multiplication

# Matrix Vector Multiplication

p1 =a[0] x b[0]
+ a[1] x b[1]
+ a[2] x b[2]
+ a[3] x b[3]
+ a[4] x b[4]
+ a[5] x b[5]

a[5]

a[0] x c[0]
+ a[1] x c[1]
+ a[2] x c[2]
+ a[3] x c[3]
+ a[4] x c[4]

c[5]

# Matrix Vector Multiplication

**p1 =a[0] x b[0]**
**+ a[1] x b[1]**
**+ a[2] x b[2]**
**+ a[3] x b[3]**
**+ a[4] x b[4]**
**+ a[5] x b[5]**

**p2 =a[0] x c[0]**
**+ a[1] x c[1]**
**+ a[2] x c[2]**
**+ a[3] x c[3]**
**+ a[4] x c[4]**
**+ a[5] x c[5]**

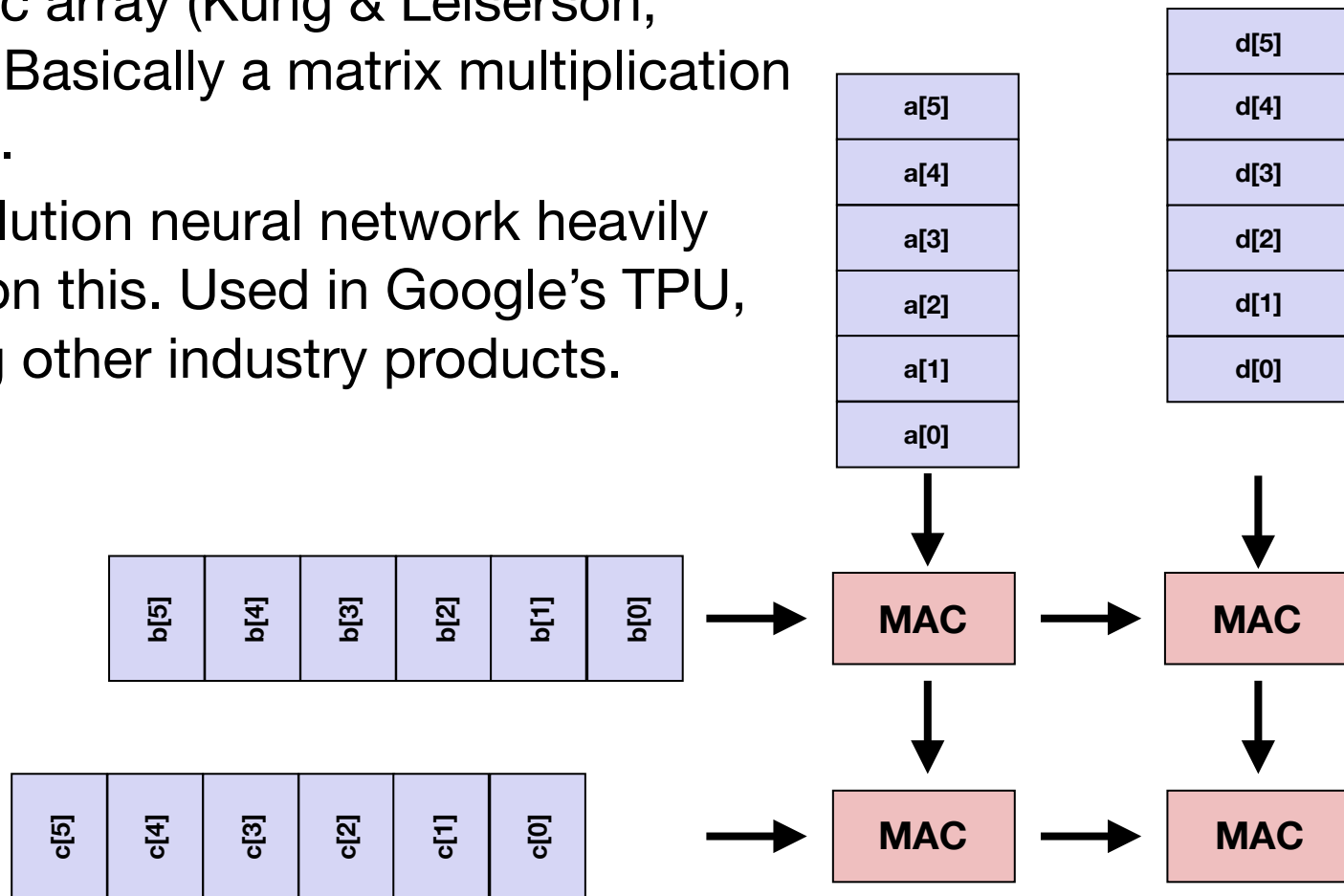# Matrix Matrix Multiplication
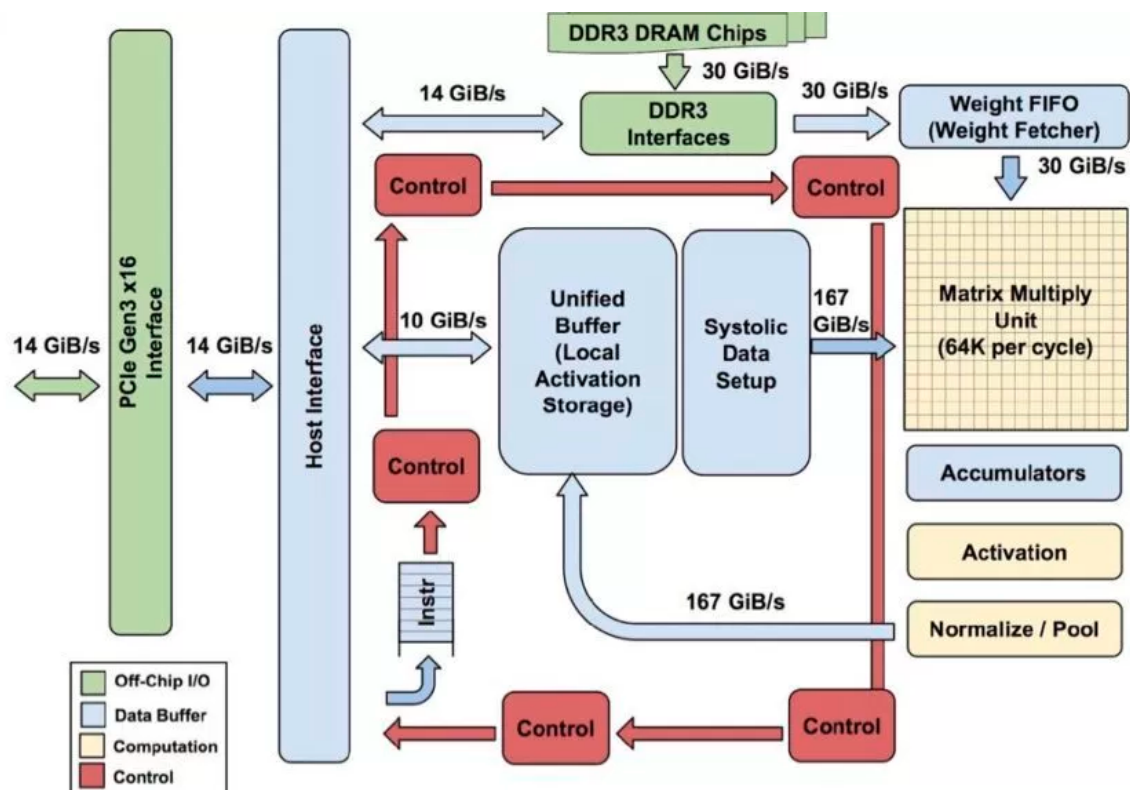
# Matrix Matrix Multiplication

# Matrix Matrix Multiplication

- Systolic array (Kung & Leiserson, 1978). Basically a matrix multiplication engine.

- Convolution neural network heavily relies on this. Used in Google's TPU, among other industry products.

# Google Tensor Processing Unit

- Convolution in deep learning can be transformed to matrix multiplication.
- TPU: specialized processor (i.e., systolic array architecture) for tensor processing (matrix multiply)
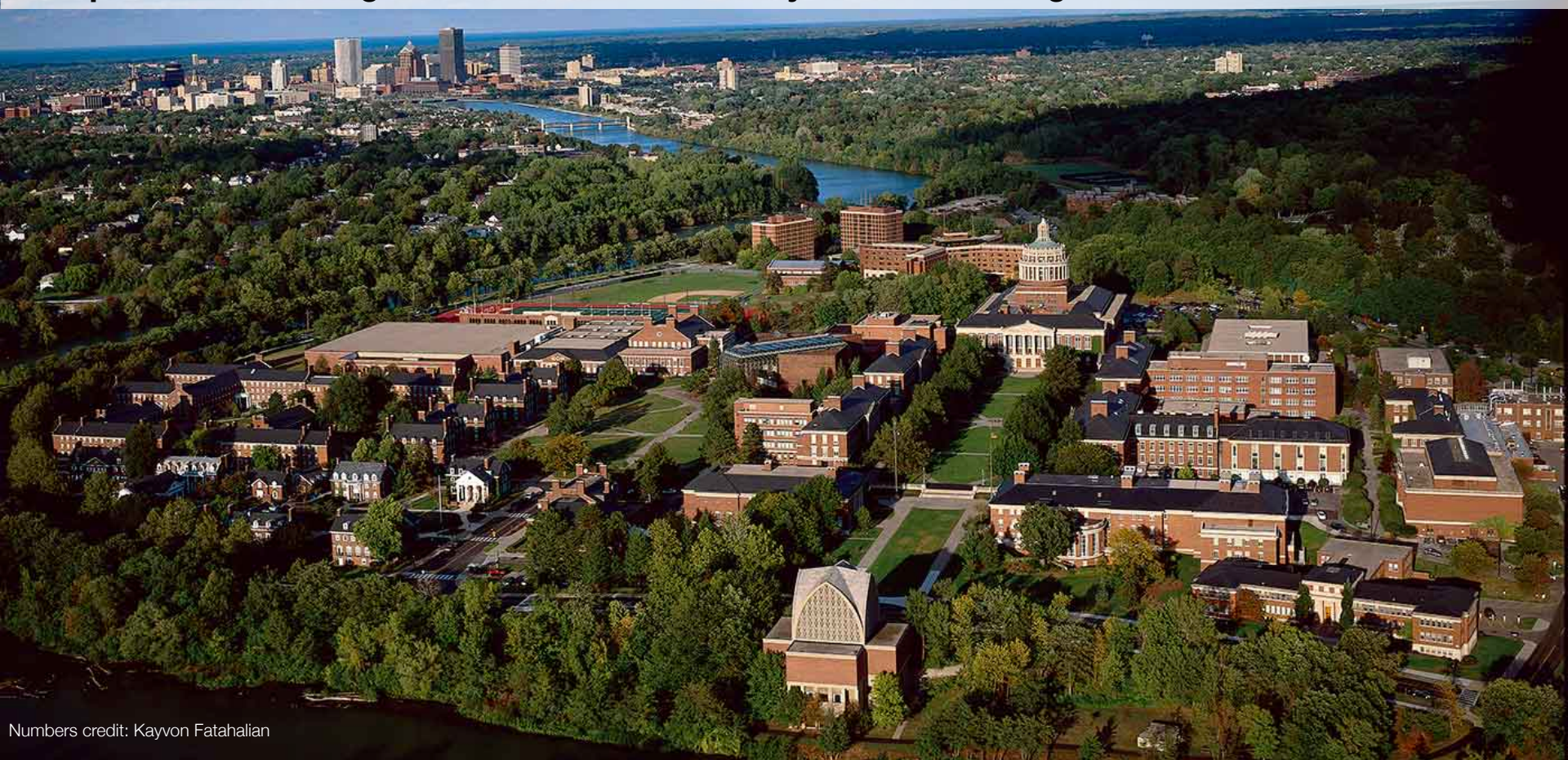  - 30x~80x more power-efficient than GPU

# Another Domain: Video Compression

# Another Domain: Video Compression

30-second video @ 1080p resolution (1920 x 1080 pixels per frame) @ 30 frames per second (FPS)
3 colors per pixel + 1 byte per color → 6.2 MB/frame → 6.2 MB x 30 s x 30 FPS = 5.2 GB total size
Actual H.264 video file size: 65.4 MB (**80-to-1 compression ratio**).
**Compression/encoding done in real-time without you even realizing it!**



Numbers credit: Kayvon Fatahalian

# Another Domain: Computational Photography

- Use computational algorithms to mimic a DSLR.
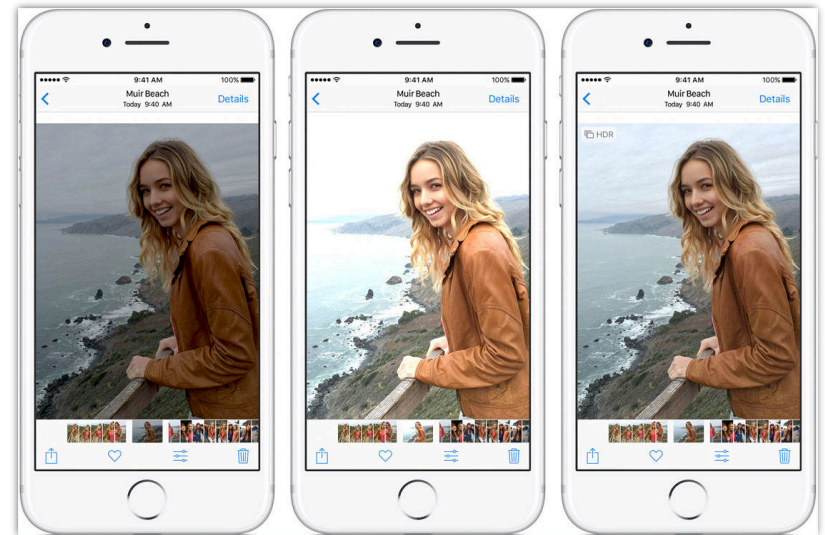- Must be done in real-time. Executed on a dedicated Image Signal Processor (ISP).



**Conventional cameras**



**Today's "cameras"**

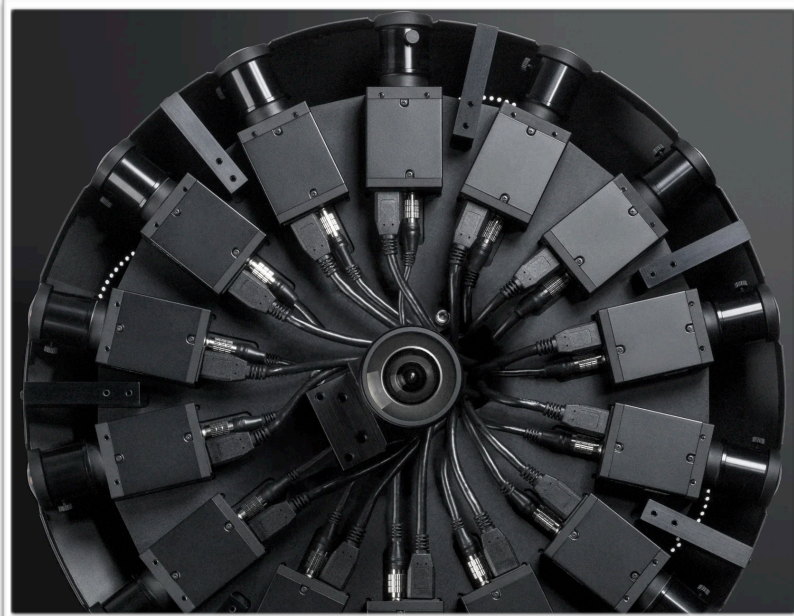# Another Domain: Computational Photography

**Portrait mode: simulate a large aperture**



**HDR mode: simulate a high dynamic range sensor**

# 360° (VR/Panoramic) Videos and Photos



https://www.cursosfotografiabarcelona.com/como-hacer-una-foto-esferica-360-grados-reflex/

# VR Video Capturing



**Facebook Surround 360**

**Google Jump VR**

# Autonomous Machines

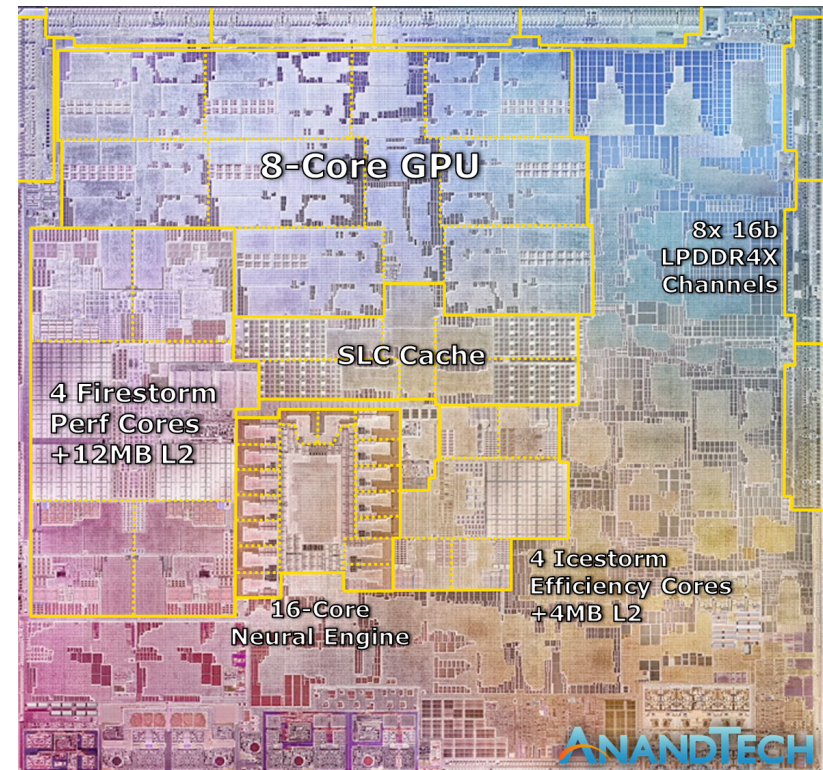# Autonomous Machines

61

# Photorealistic Rendering

# Photorealistic Rendering

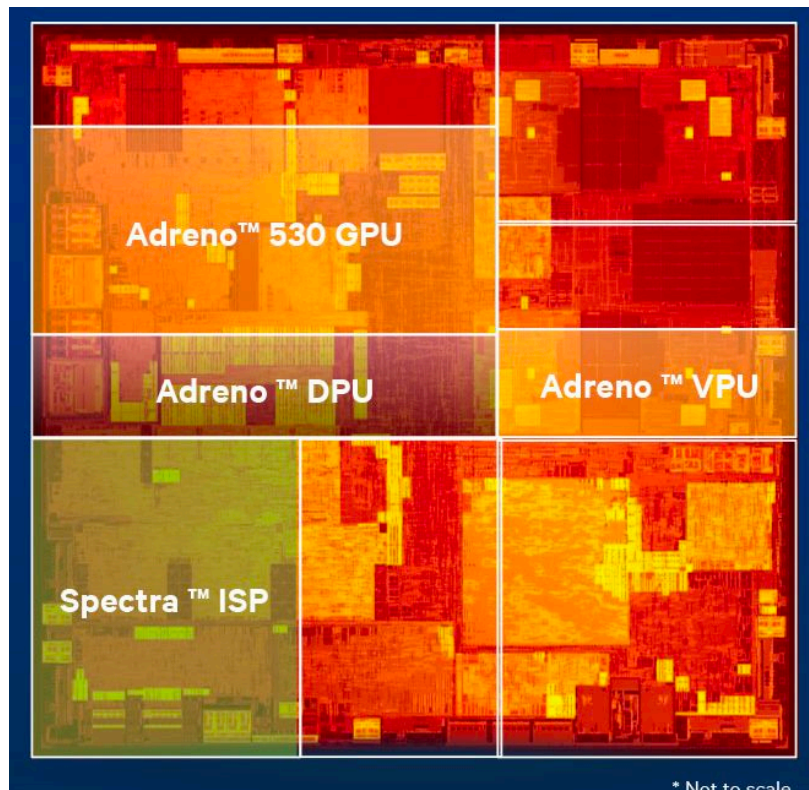# Photorealistic Rendering

# Today's Processor Chips are Full of Accelerators
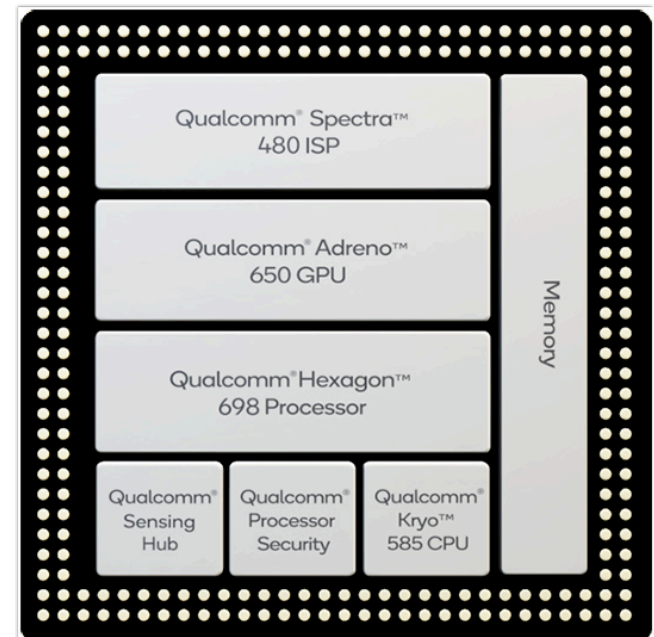
# Today's Processor Chips are Full of Accelerators

Qualcomm
Snapdragon
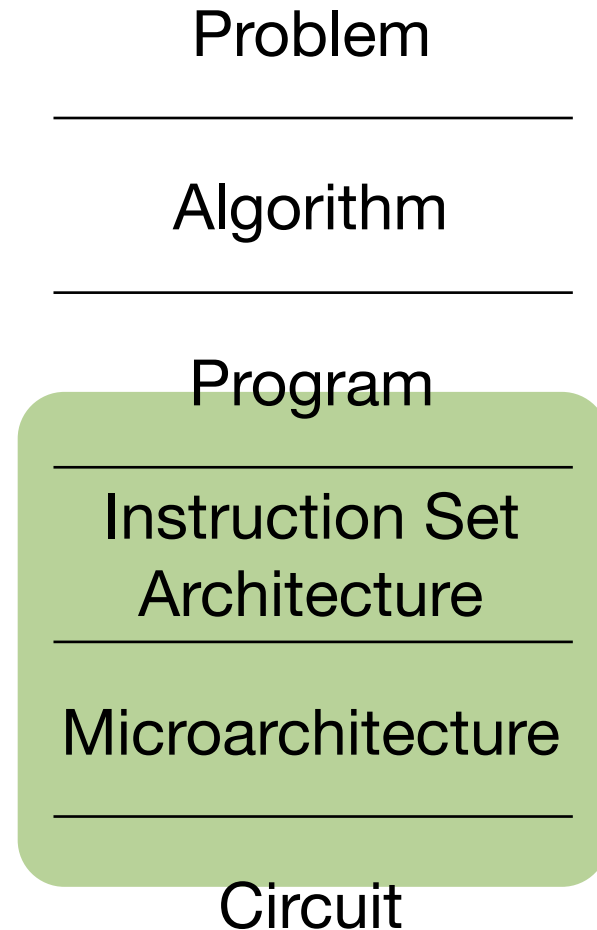820 SoC



Qualcomm
Snapdragon
835 SoC

# Traditional Scope of Computer Systems

- Take a program and try to figure out how to best execute on the hardware

Problem

Algorithm

Program

Instruction Set Architecture

Microarchitecture

Circuit

# Real Scope of Computer Systems

- Understand the problem to be solved, design algorithms, understand algorithms characteristics to design the best computer systems.

- It's no longer enough to work with a given program without understanding it.

| Problem |
| --- |
| Algorithm |
| Program |
| Instruction Set Architecture |
| Microarchitecture |
| Circuit |

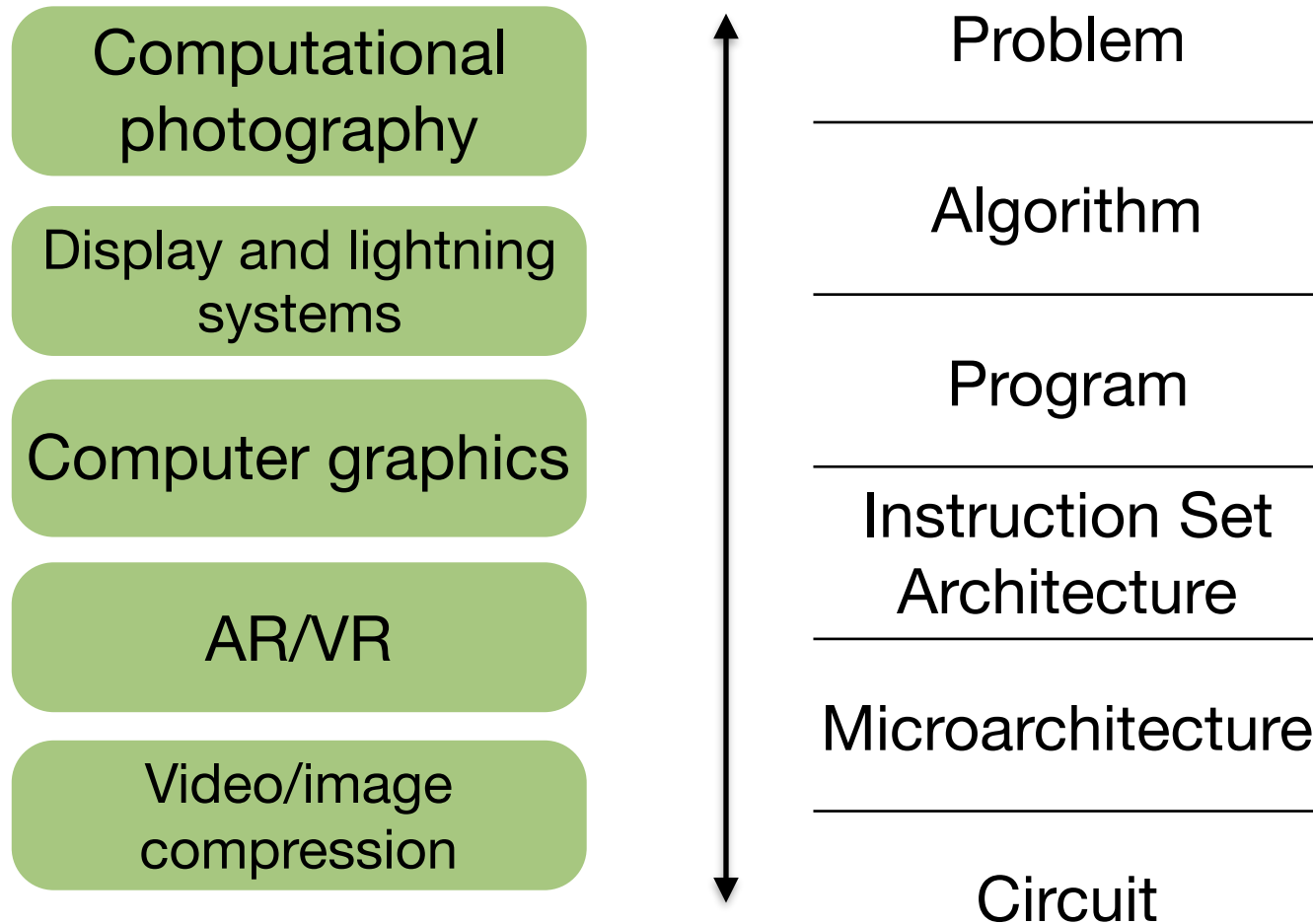# CSC 292/572: Mobile Visual Computing

Computational photography

Display and lightning systems

Computer graphics

AR/VR

Video/image compression

Problem

Algorithm

Program

Instruction Set Architecture

Microarchitecture

Circuit

# The Most Important Take Away of 252

- "There is no magic."
- Every thing can be derived from first principles. Trust your logical reasoning.
- Apply to virtually everything in science and engineering.

# The Second Most Important Take Away of 252

- "Things don't have to be this way."

- As long as you don't violate physics, you can design a computer however you want.

- But every design decision you make usually involves certain trade-offs. Be clear what your design goal is.

# The Third Most Important Take Away of 252

- Virtual all computer system design practices follow a small set of basic principles.

- It is these basic principles that are important, not the practices.

**Make common case faster**
{
Locality

Parallelism

Speculation

Specialization
}

**Combine the best of both worlds**
{
Heterogeneity

Hierarchy

#pragma
}

**Virtualization**