

# The inexorable rise of hardware specialization

Nathan Beckmann, CMU  
~~CSD faculty meeting, Fall '21~~  
15-418/618 Fall '25

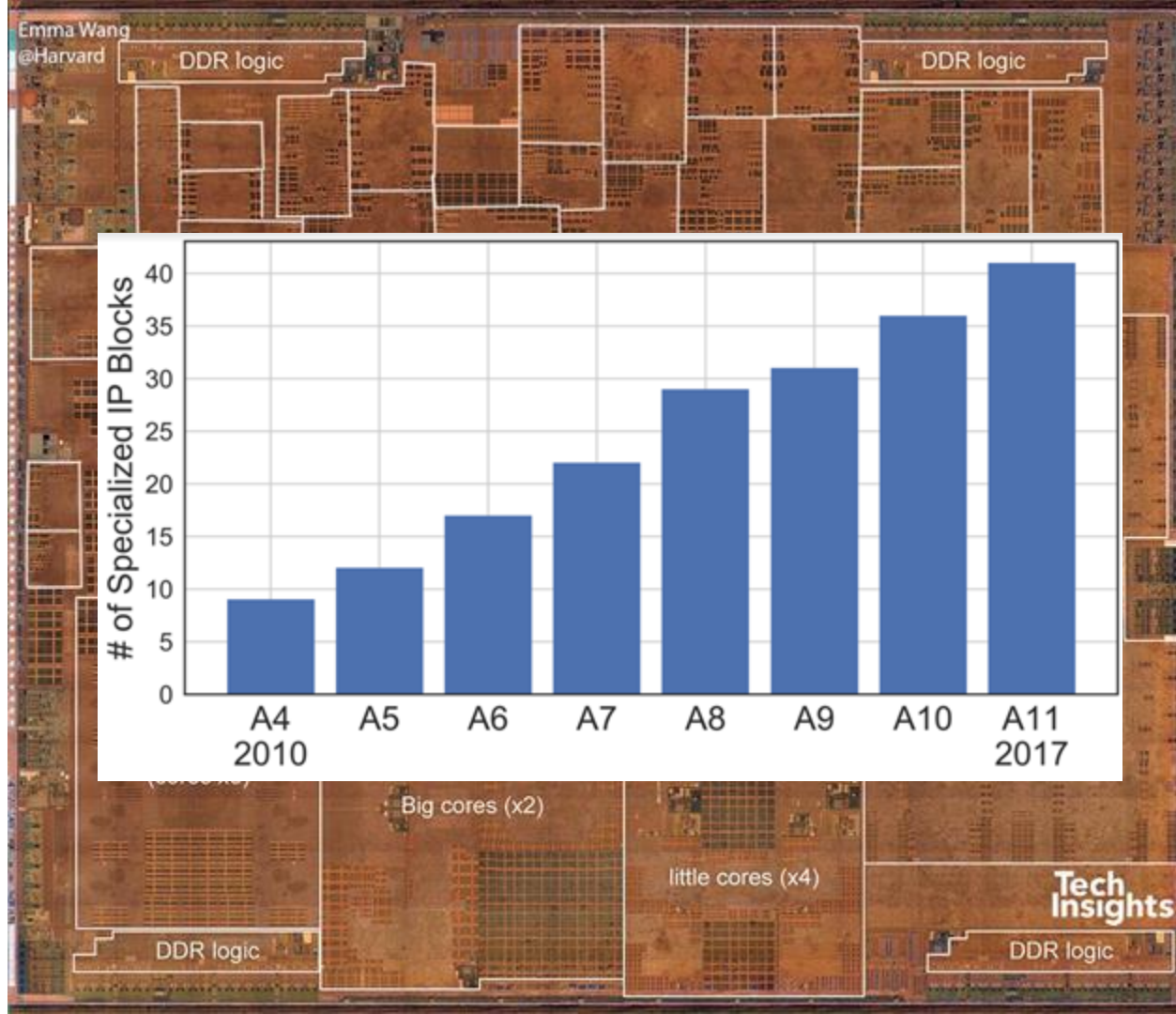
# The world today:

## Apple A12 (2019)

- 42 different “IP blocks” on one processor chip
- High-performance CPUs
- Energy-efficient CPUs
- GPU
- “Neural engine”
- + 38 other specialized widgets

Traditional CPUs are a small fraction of the chip!

Credit: Emma Wang & Sophia Shao,  
<http://vlsiarch.eecs.harvard.edu/research/accelerators/die-photo-analysis/>

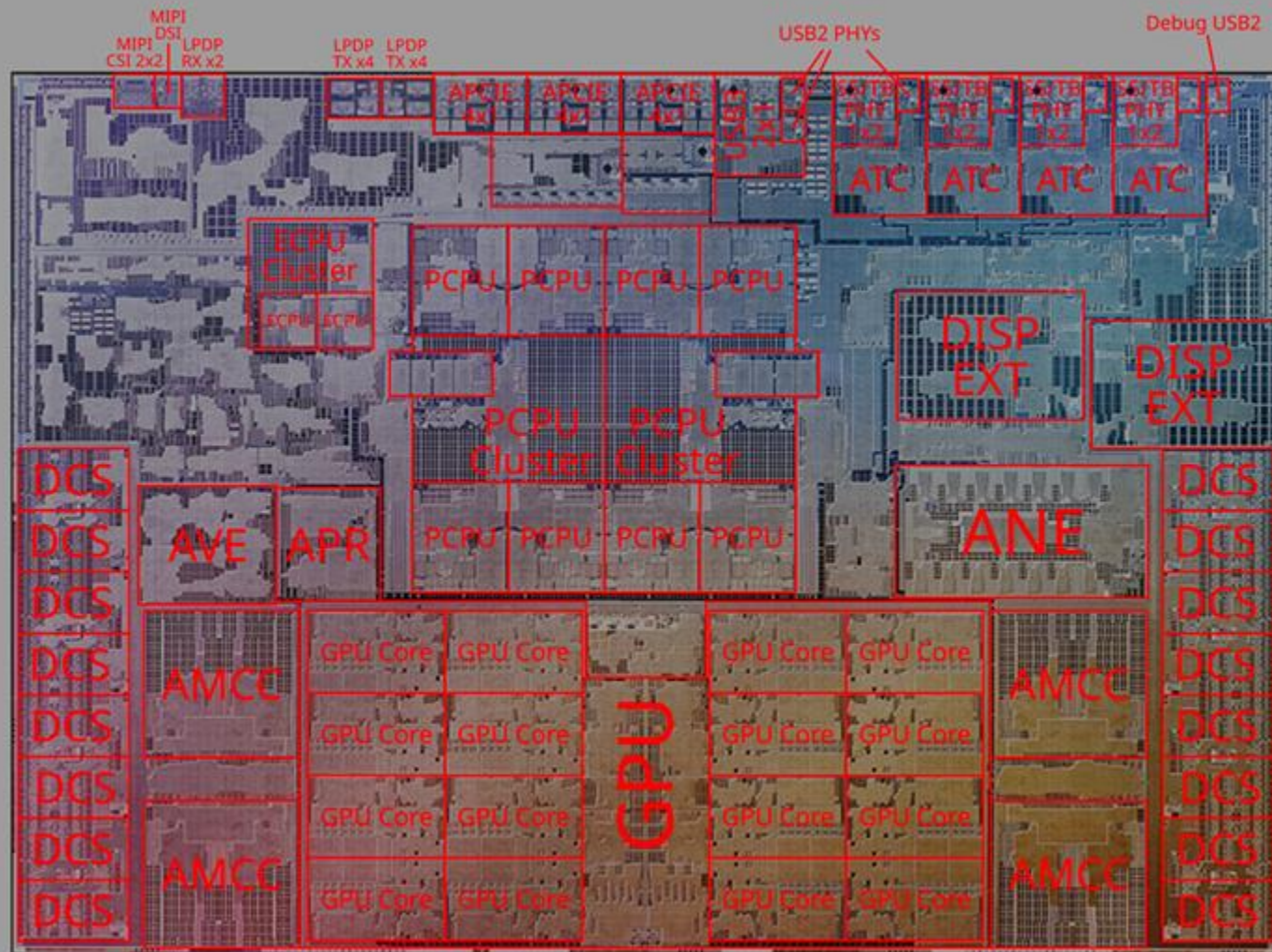


# Laptops, too

Apple M1 chip

- Big CPUs
- Small CPUs
- GPU
- “Neural engine”
- Video encoder
- PCIe controller
- Etc.

Credit: Hector Martin  
<https://twitter.com/marcan42/>





# Also, Intel!

## Intel Alder Lake

- Big CPUs
- Little CPUs
- GPU
- Etc

Credit: Nemez

<https://twitter.com/GPUsAreMagic>



What's going on with hardware?

# Takeaways

1. The world is moving away from traditional, fully programmable cores
2. “Non-von Neumann” computing seems to be the way to continue scaling
3. There isn’t a very compelling story for how one should write programs
4. The best balance between programmability & efficiency is an open question

Why should you care  
about hardware?

# Goal: Performance & efficiency scaling

Hardware performance has grown exponentially for decades

Basic assumption: *This is good, and we want it to continue*

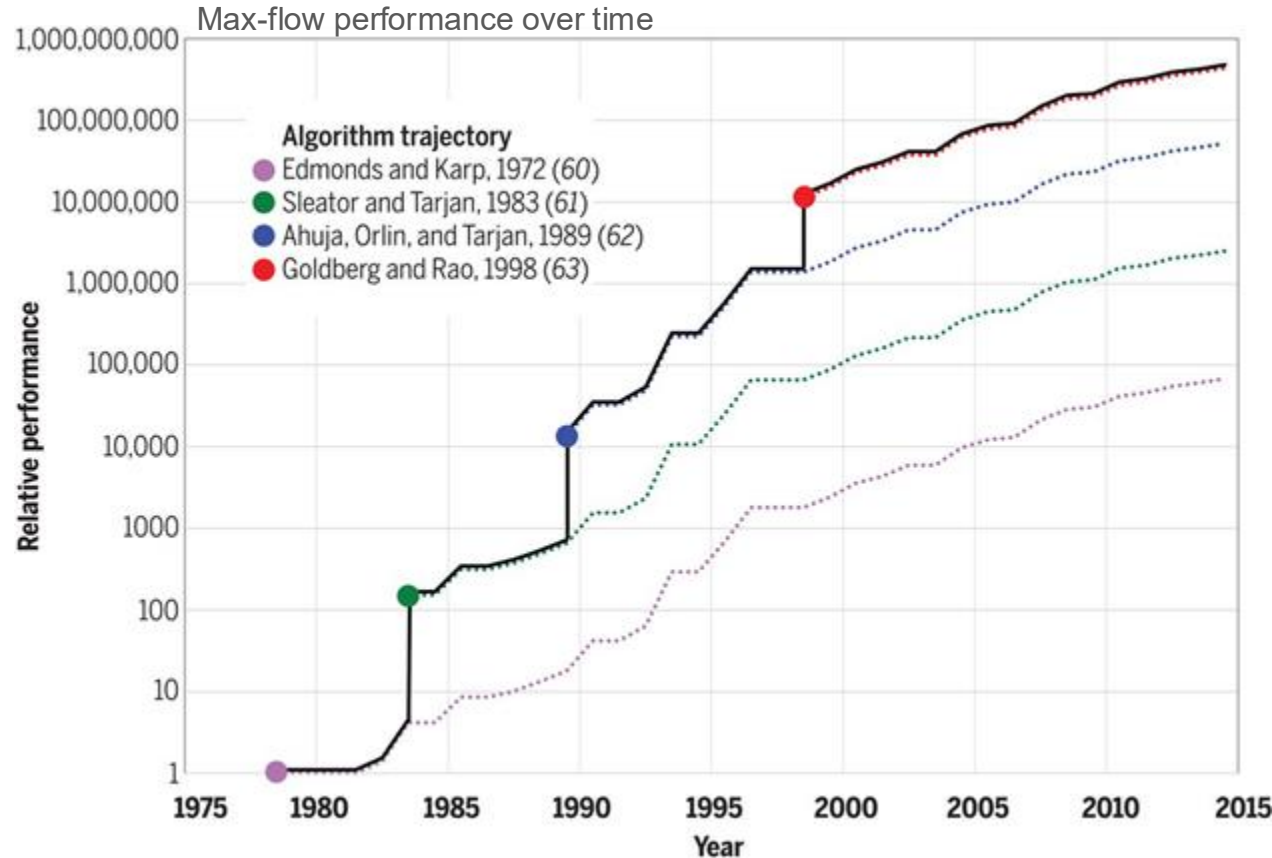


# Algorithms, hardware, or both?

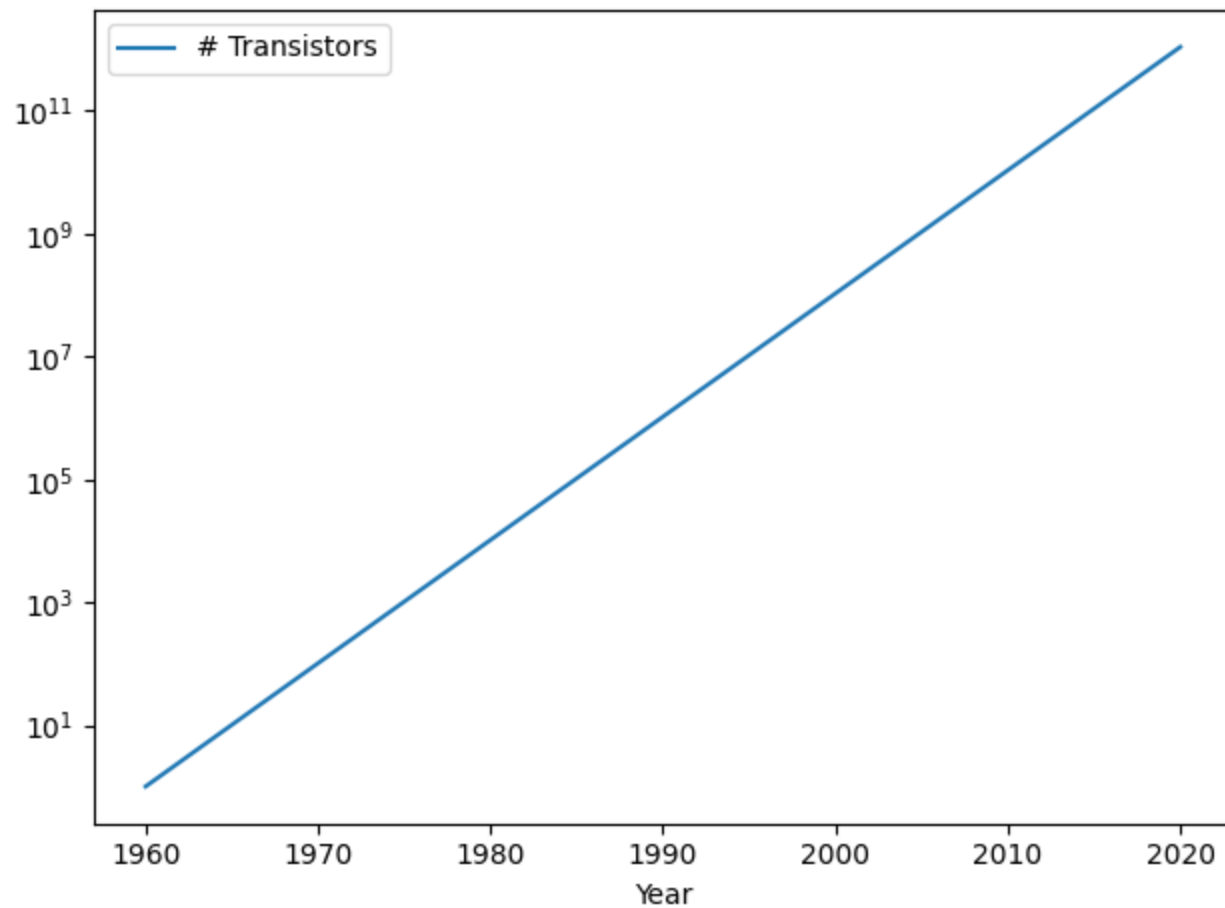
Similar order-of-magnitude improvements from algorithm & hardware improvements

- Diminishing returns from both!

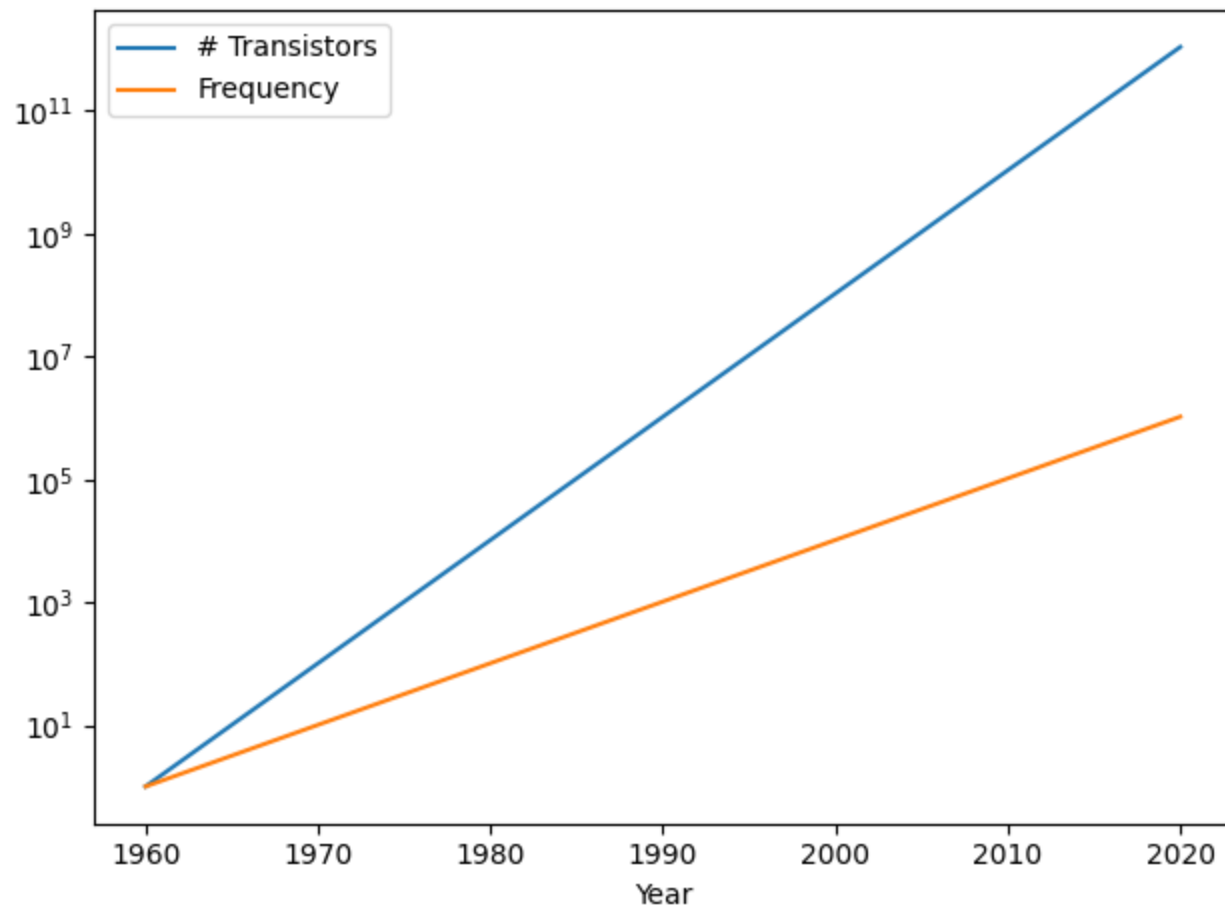
Credit: “There’s plenty of room at the top”, Leiserson+, Science ‘20



OK, so what's going on with hardware?

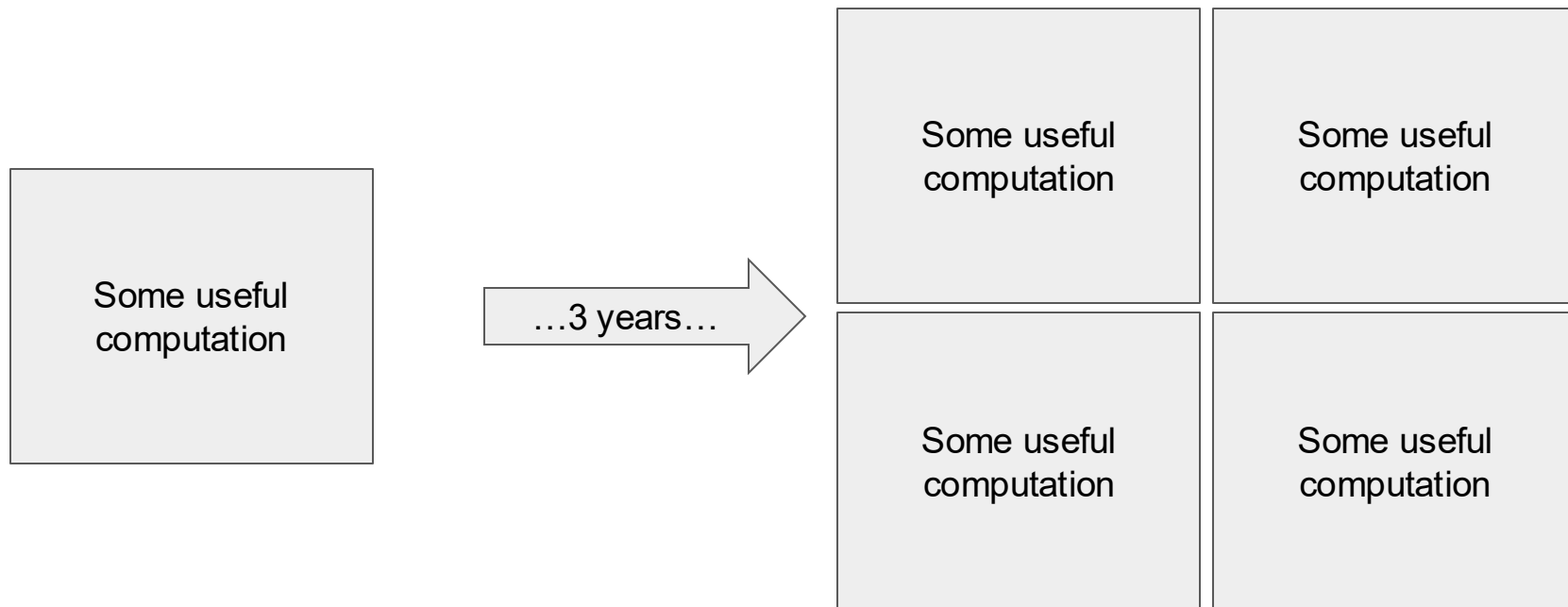


Moore's Law: 2x transistors every 18 months



Smaller transistors are faster, too!

# Optimistic performance model



# How does performance scale?

Optimistic performance = Number of transistors  $\times$  frequency =  $N f$

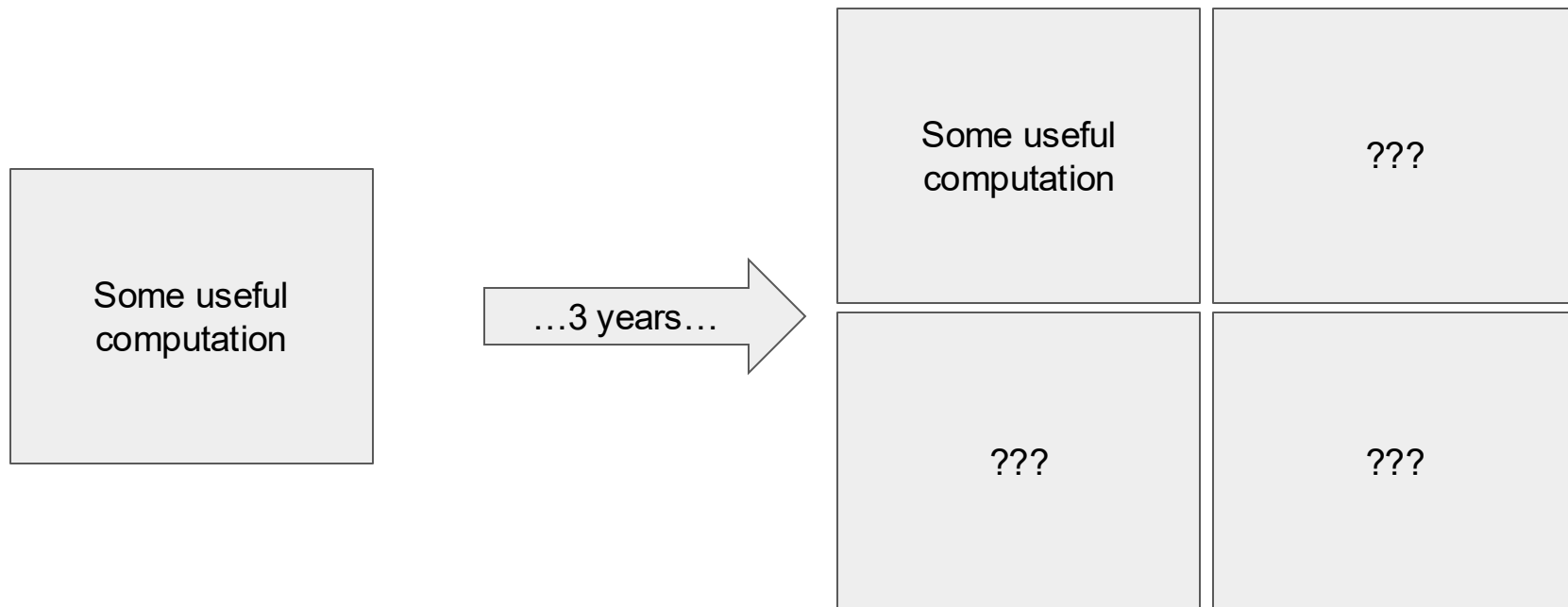
$$N' = 2 N$$

$$f' = \sqrt{2} f$$

$\Rightarrow$  Optimistic scaling factor =  $2\sqrt{2}\times$  per technology generation



# Pessimistic performance model

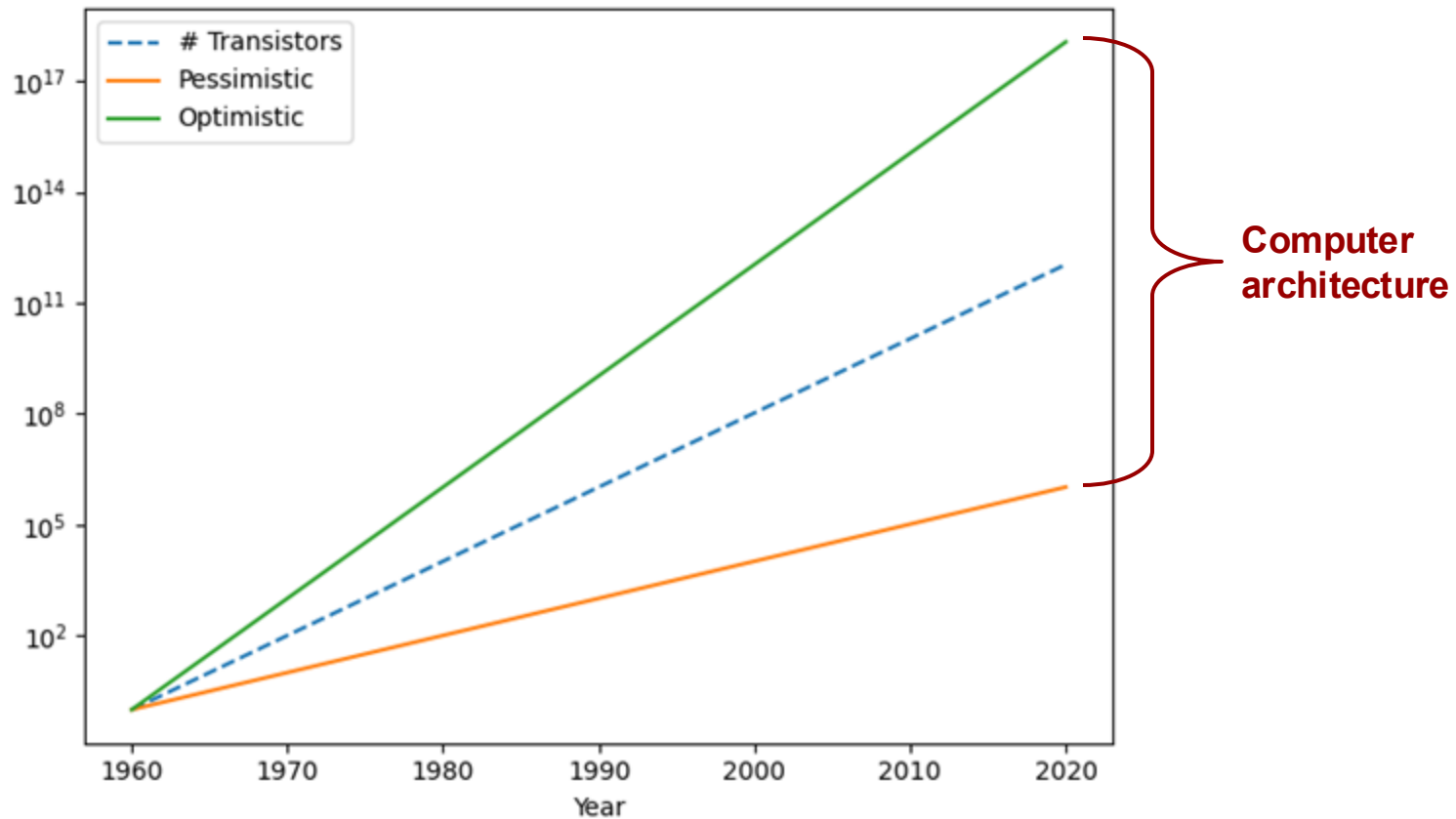


# How does performance scale?

Pessimistic performance = frequency =  $f$

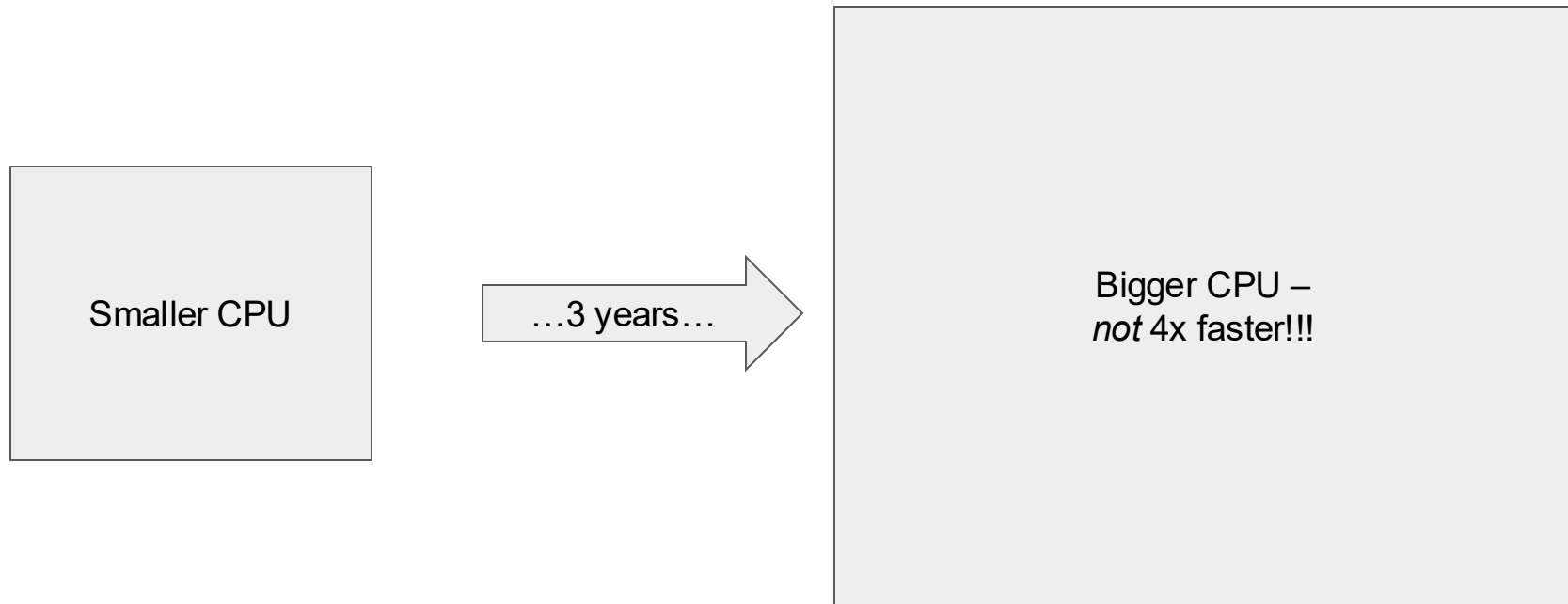
$$f' = \sqrt{2} f$$

⇒ Pessimistic scaling factor =  $\sqrt{2} \times$  per technology generation



Huge performance gap! *Need to put extra transistors to good use!*

# Better historical performance model



# Better historical performance model

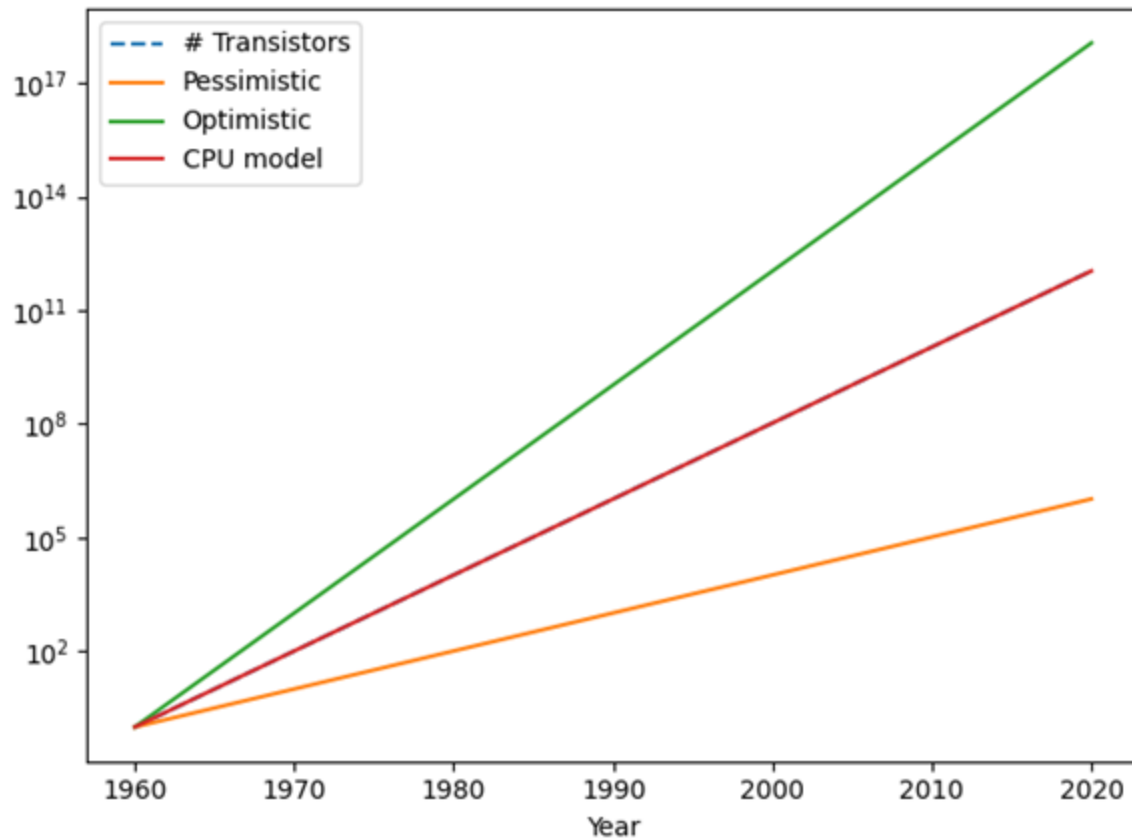
Realistic performance  $\approx \sqrt{N} f$

(slightly optimistic)

$$N' = 2N$$

$$f' = \sqrt{2}f$$

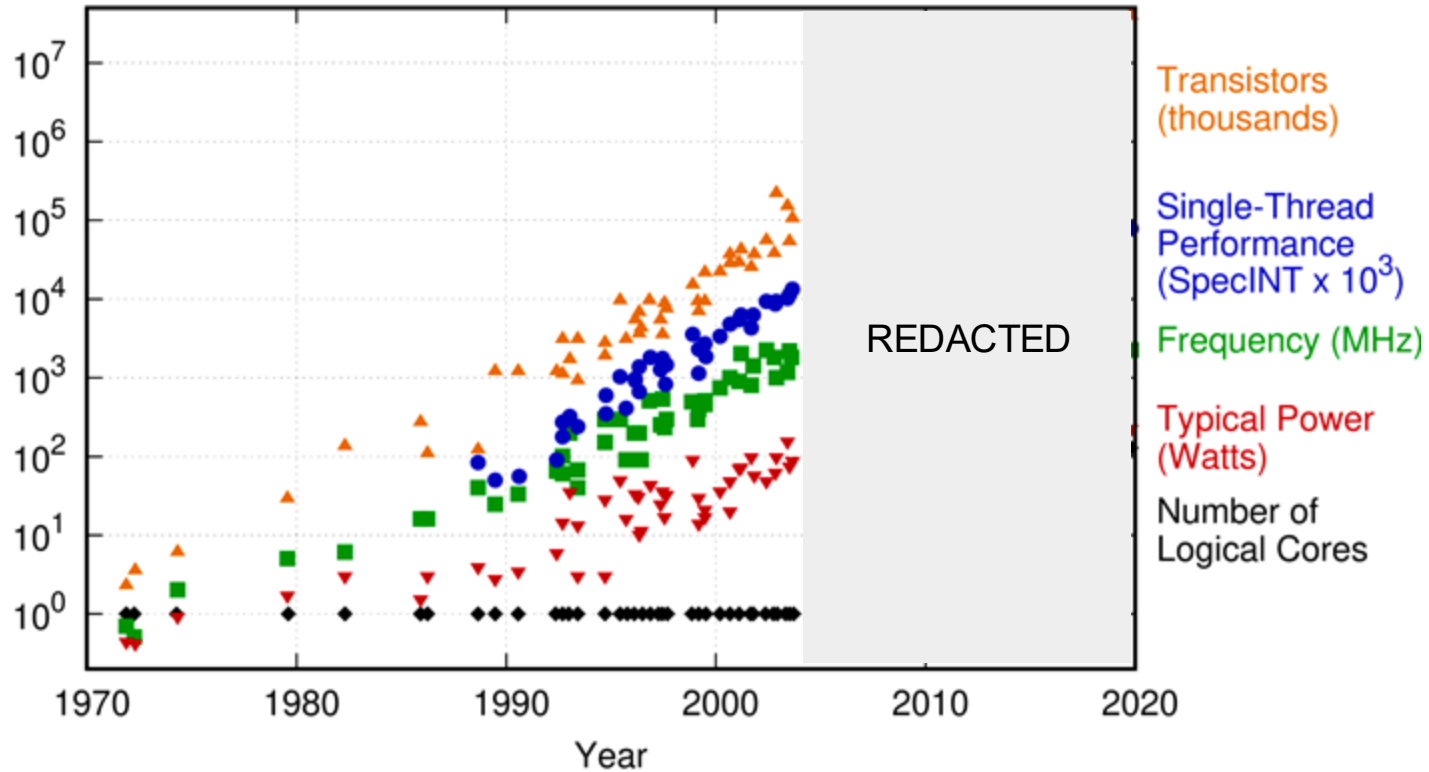
⇒ Scaling factor = 2× per technology generation



CPU's close ~half the gap & are easy to program



## 48 Years of Microprocessor Trend Data



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten  
New plot and data collected for 2010-2019 by K. Rupp

Cartoons aren't far from reality!

# The elephant in the room



# What about power & energy?

Power = Number of transistors  $\times$  Transistor power

# What about power & energy?

Power = Number of transistors  $\times$  Transistor power

$$\text{Transistor power} = C V^2 f$$

The diagram illustrates the scaling of transistor power with transistor dimensions. It features the equation  $\text{Transistor power} = C V^2 f$  at the top. Three blue arrows point from descriptive text below to the variables in the equation: one from 'Capacitance  $\approx$  Transistor width' to 'C', one from 'Voltage  $\approx$  Transistor length' to 'V', and one from 'Frequency  $\approx$  Transistor length' to 'f'.

Frequency  $\approx$  Transistor length

Capacitance  $\approx$  Transistor width

Voltage  $\approx$  Transistor length

# How do these parameters scale?

After one step of Moore's Law...

$$N' = 2 N$$

Area  $\times 1/2$

$$C' = C / \sqrt{2}$$

$$V' = V / \sqrt{2}$$

$$f' = \sqrt{2} f$$

Linear dimension  $\times 1/\sqrt{2}$

(All of the above just from the geometry of a smaller transistor!)

# Dennard scaling: Power is constant!

$$\begin{aligned}\text{Power}' &= N' C' V'^2 f' \\ &= 2N (C/\sqrt{2}) (V / \sqrt{2})^2 \sqrt{2}f \\ &= N C V^2 f \\ &= \text{Power}\end{aligned}$$

Power is **constant!**

(Even with  $2\sqrt{2}\times$  more transistor toggling)



# Dennard scaling ended in ~2003.

*(Voltage can only drop up to a point. Band gap, and all that.)*

# What does this mean for performance scaling?

Without Dennard scaling...

uh-oh! 😞

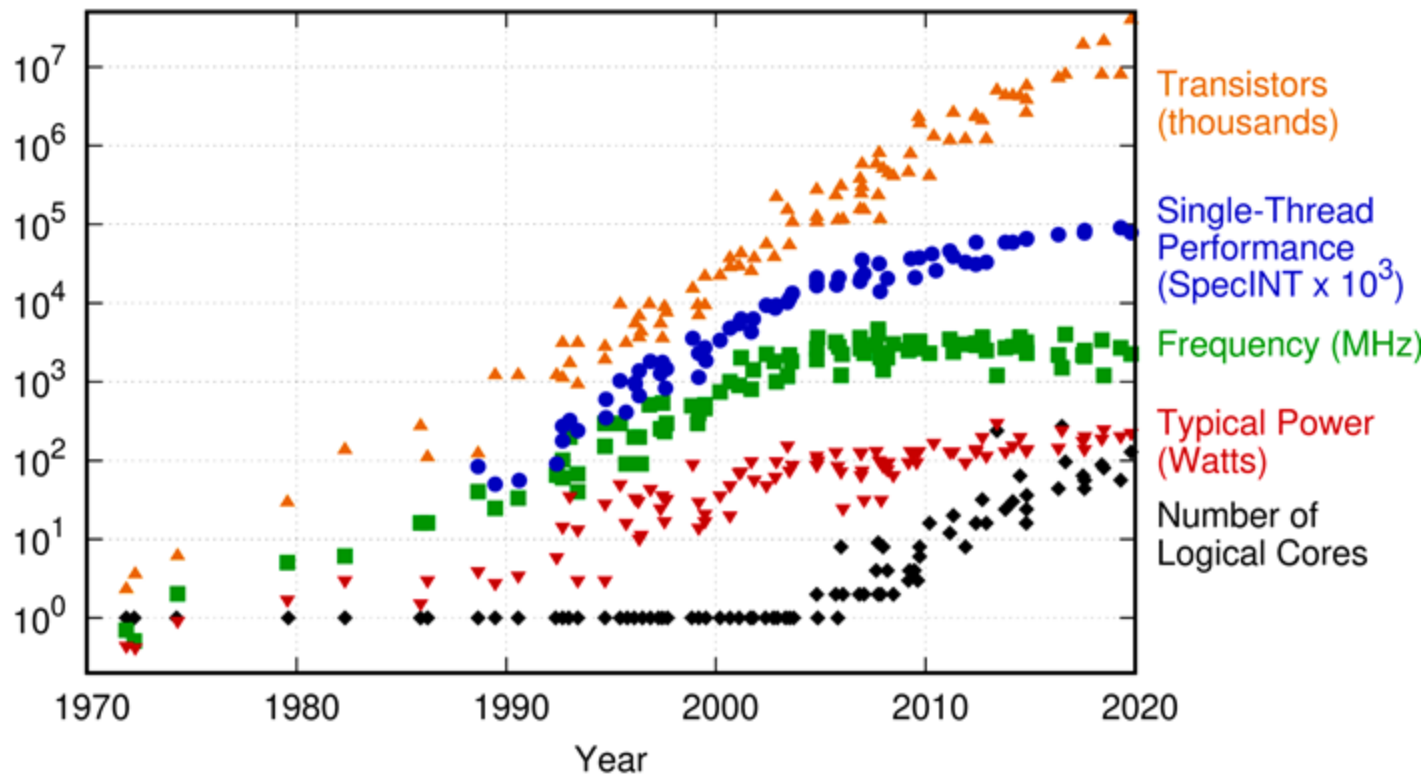
$$\text{Power}' = N' C' V'^2 f' = 2N (C/\sqrt{2}) V^2 \sqrt{2}f = \mathbf{2\times \text{Power}}$$

But power *must* remain ~constant! (heat sinks have limits)

$$(\text{Power}' / \text{Power}) = 1 \quad \Rightarrow \quad N' f' = \sqrt{2} N f$$

Steep decrease in performance scaling (lose 2× off optimistic)

## 48 Years of Microprocessor Trend Data



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten  
New plot and data collected for 2010-2019 by K. Rupp

# Technology scaling recap

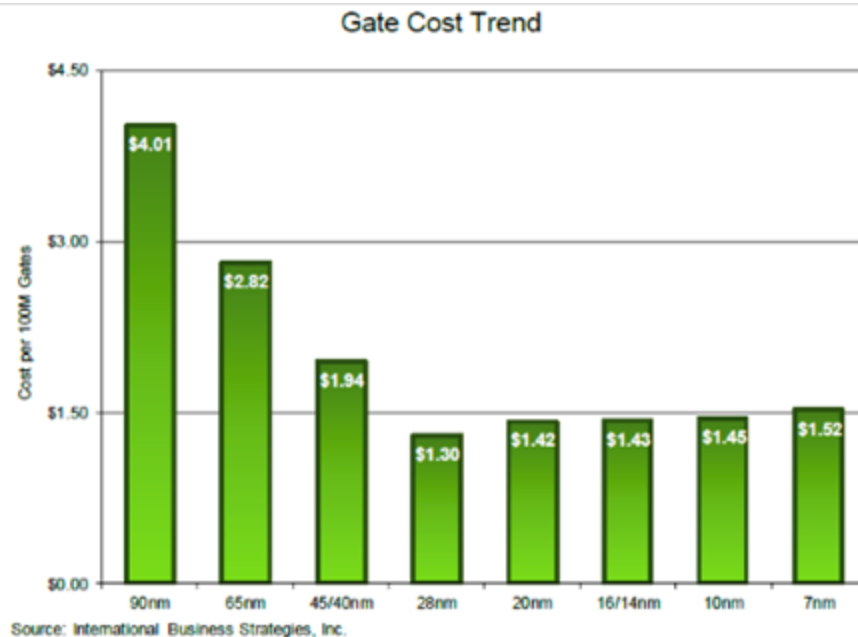
Moore's Law isn't dead

- At least not in Taiwan & Korea

**But!**

Cost scaling has stopped

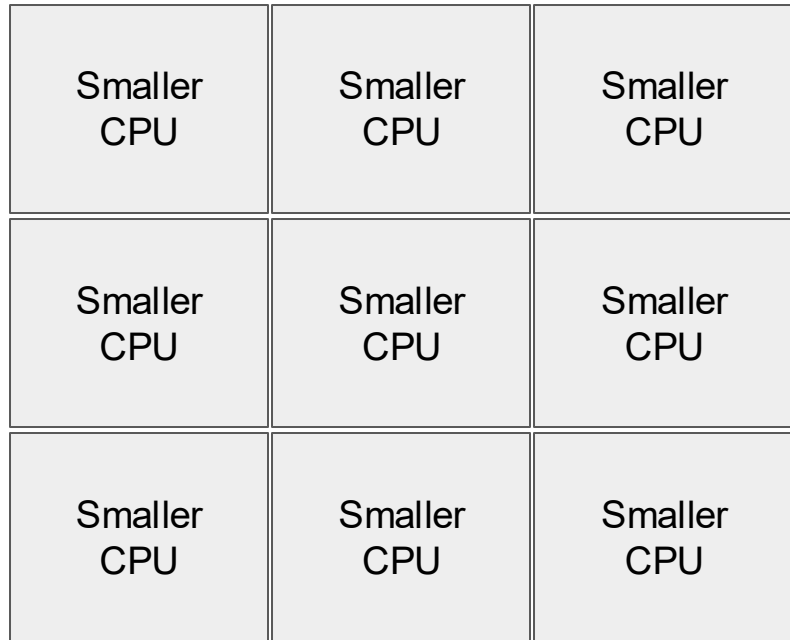
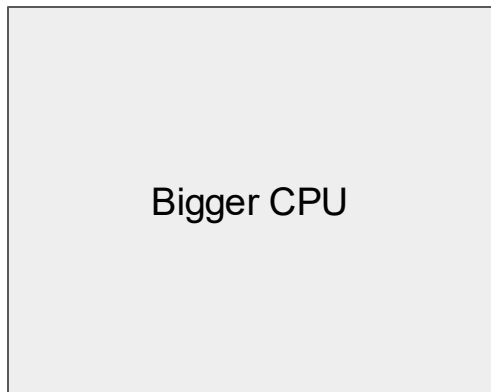
Dennard (power) scaling has stopped



# Back to the drawing board...

“Manycore” designs

- TRIPS, RAW, Tiler





# “Amdahl’s Law in the Multicore Era”

Model programs as sequential & parallel regions

Same CPU model ( $\text{perf} = \sqrt{N}$ ), but now splitting  $N$  across different #s of cores

Amdahl’s Law:  $\text{Speedup} = 1 / ( (1 - s) + s / P )$

- $s$  = fraction of program that can run in parallel
- $P$  = number of processors
- (assumes ideal speedup in parallel region)





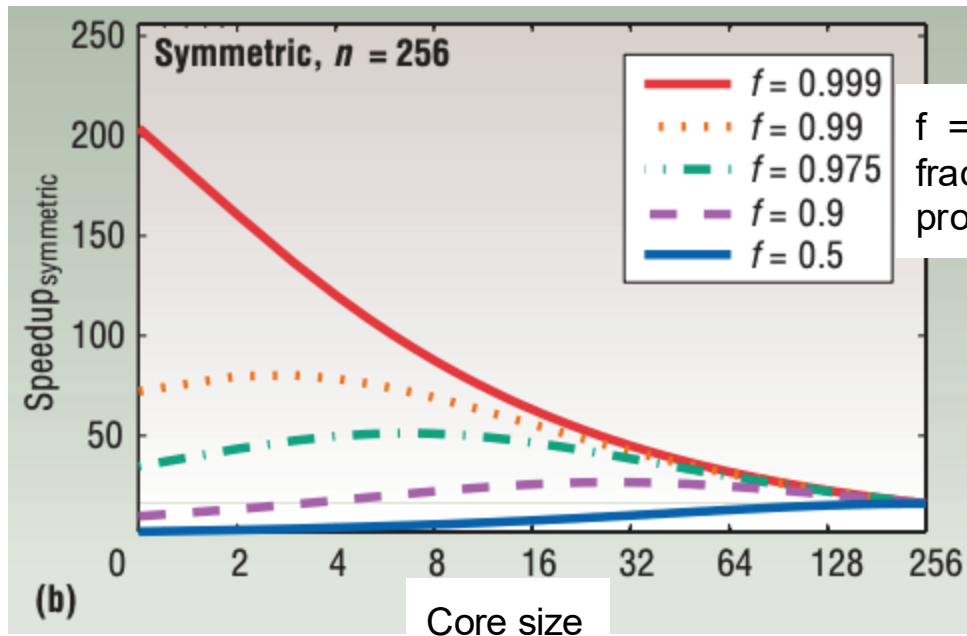
# “Amdahl’s Law in the Multicore Era”

Speedup from a manycore design:

$x = \text{CoreSize}$

$P = n / x$

Speedup =  
 $\sqrt{x} / ((1-s)+s/P)$



$f$  = Parallel  
fraction of  
program

**Speedup is  
surprisingly awful  
unless program is  
≥99.9% parallel!**

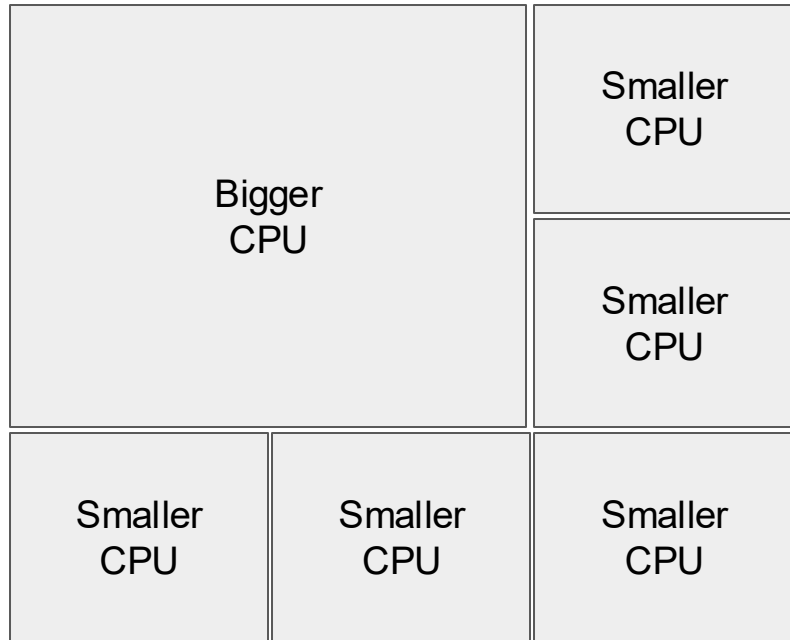
**No good single choice  
for core size!**

# Back to the drawing board again...

Problem is *sequential bottleneck*

⇒ Heterogeneous multicore design

Big CPUs are locally inefficient, but globally efficient





# “Amdahl’s Law in the Multicore Era”

Speedup from a heterogeneous/asymmetric design:

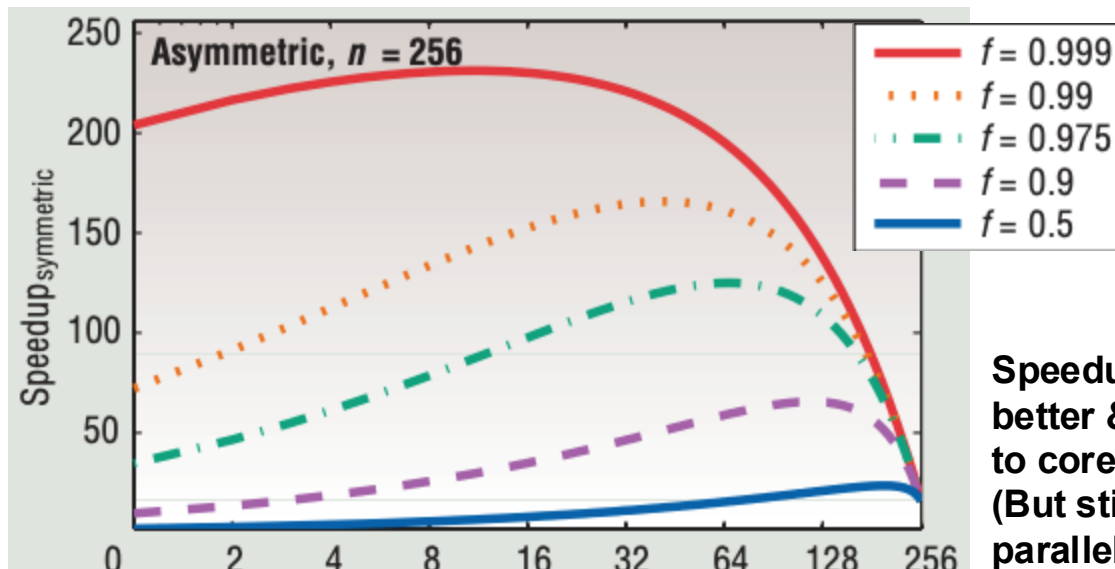
$x$  = BigCore

$PS_{\text{small}} = (n - x) + 1$

$T = PS_{\text{small}} + \sqrt{x}$

Speedup =

$$1 / ((1-s) / \sqrt{x} + s / T)$$



$f$  = Parallel  
fraction of  
program

**Speedup is much  
better & less sensitive  
to core size!  
(But still limited by  
parallel fraction)**

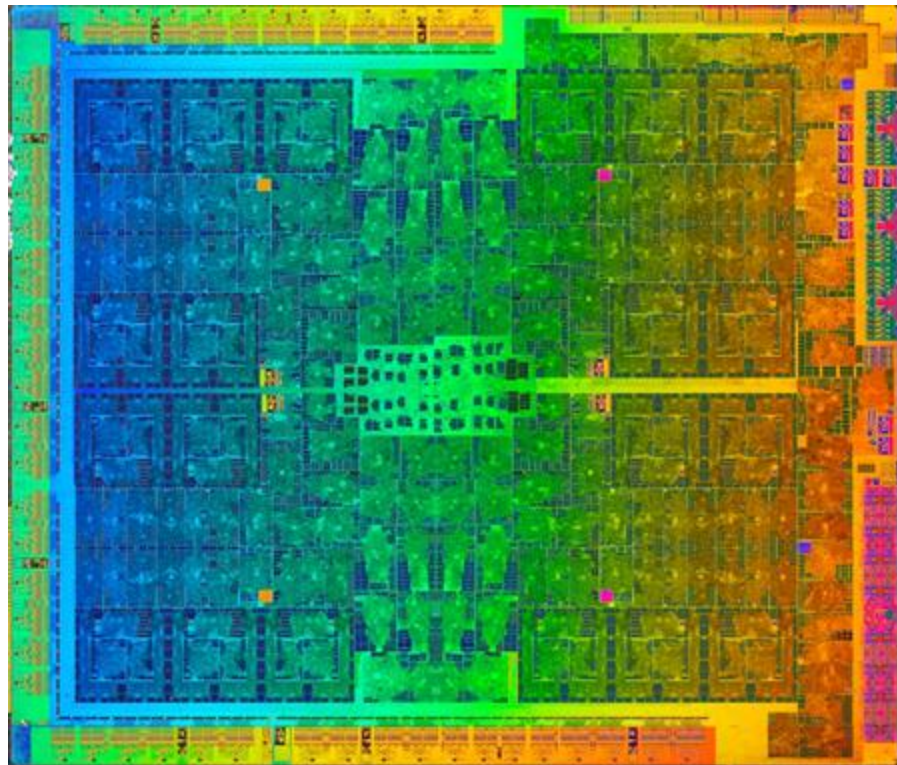
Big core size (small core size = 1)

# Heterogeneity & manycore in practice

Common in mobile (e.g., big.LITTLE)

Common in server systems:

- Intel Xeon is the “big CPU”
- Nvidia GPU is the “manycore” — design looks a little different, but the principle is the same
- Important: GPU is *not* an “accelerator” in this story; it is a programmable manycore (really, its somewhere in between)



Nvidia GP104

Why specialization?

## “Dark silicon”

$$N' = 2 N$$

*but*

$$N' f' = \sqrt{2} N f$$

⇒ We cannot use 30% of the new transistors each generation!

(Or we must run everything 30% slower...)

# Dark silicon $\Rightarrow$ Specialized hardware

Mounting pressure to use transistors efficiently

Conventional designs will have idle transistors

*Why not add custom hardware for important computations?*

“Dark silicon and the end of multicore scaling”, Esmailzadeh+ ISCA’11

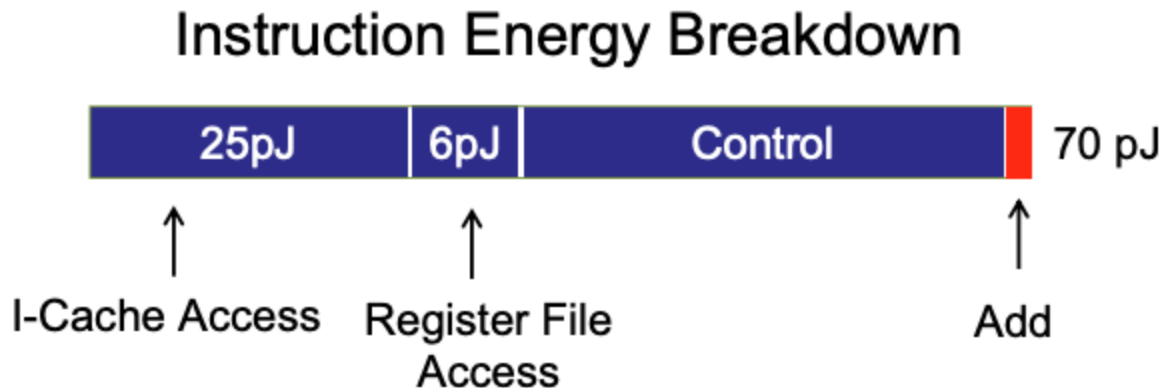
“Conservation cores: Reducing the energy of mature computations”, Venkatesh+ ASPLOS’10



# CPUs are very inefficient

CPU spends  $\approx 1\%$  energy on useful work

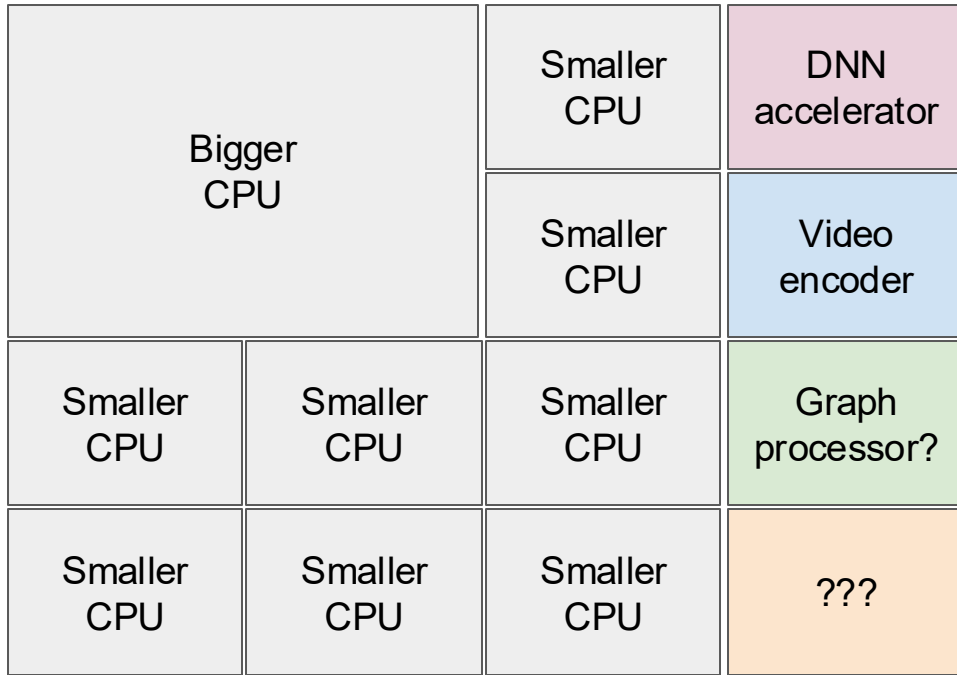
$\Rightarrow 100\times$  opportunity!





# Cartoon of a system on a chip (SoC)

Replaces some  
small CPUs in  
heterogeneous  
multicore with  
specialized  
accelerators that  
are ~100x more  
energy-efficient



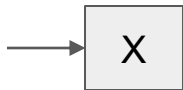
# Under the hood: Basics of accelerator design

Goal: Spend all energy on useful work

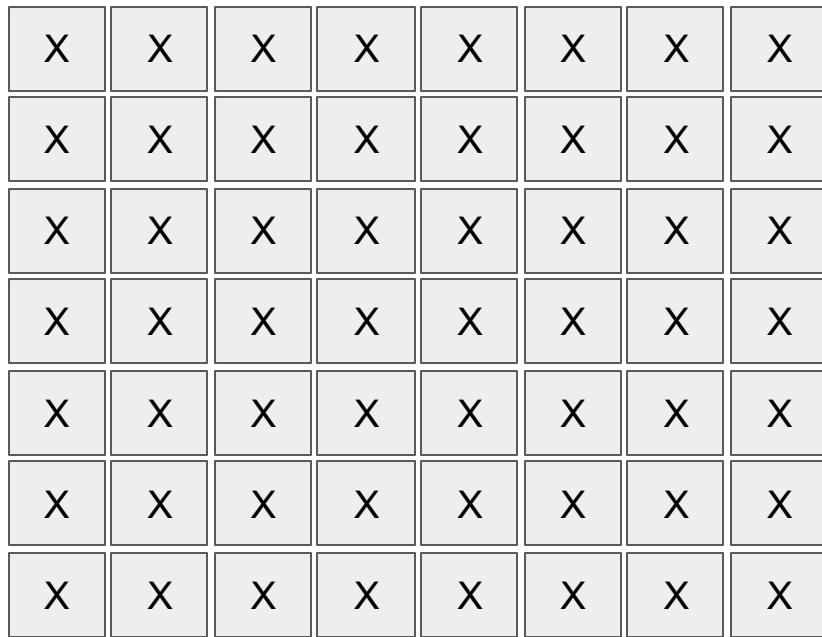
Replicate minimal, custom circuit many times

- Maximize parallelism
- Eliminate overhead
- Back to 'optimistic' scaling regime

Some useful  
computation



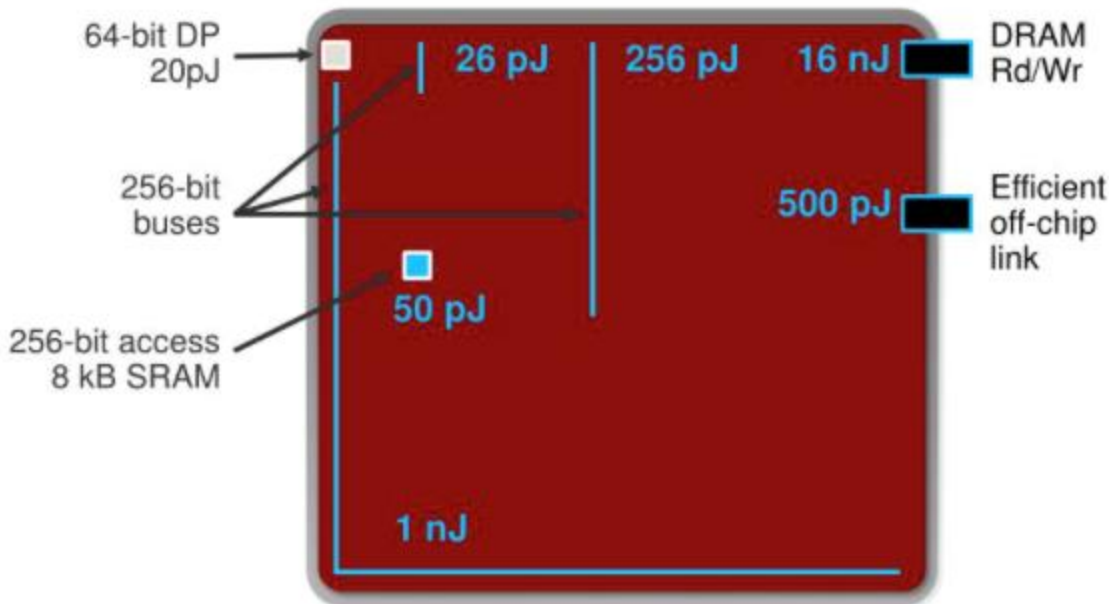
*Accelerate!*





# Accelerator design: Data movement

With CPU inefficiency removed, *data movement* consumes most energy





# Accelerators as algorithm design

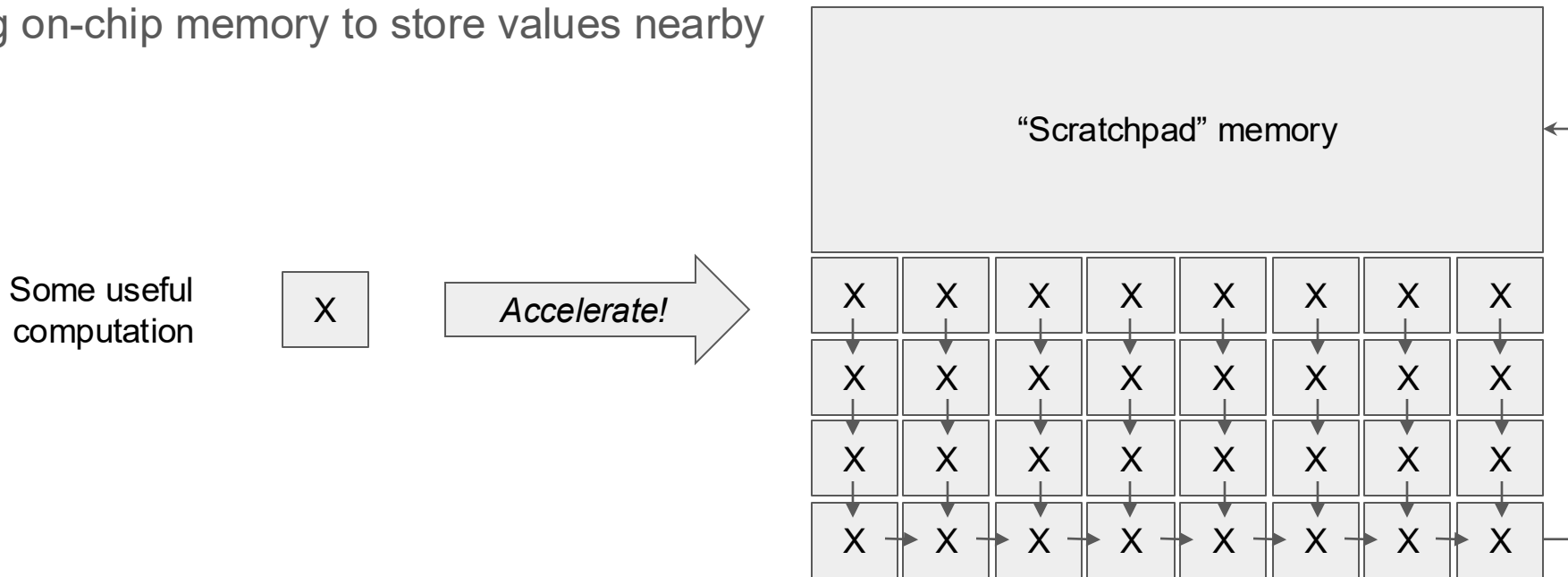
“Achieving high speedups and gains in efficiency from specialized hardware usually requires modifying the underlying algorithm...”

“**When logic is free, memory dominates.** Because the area and power of most accelerators are memory dominated, a reasonable first estimate of these costs can be made by considering only the memory...”

# Accelerator design: Data movement

Data values routed directly from one operation to another – no register file!

## Big on-chip memory to store values nearby

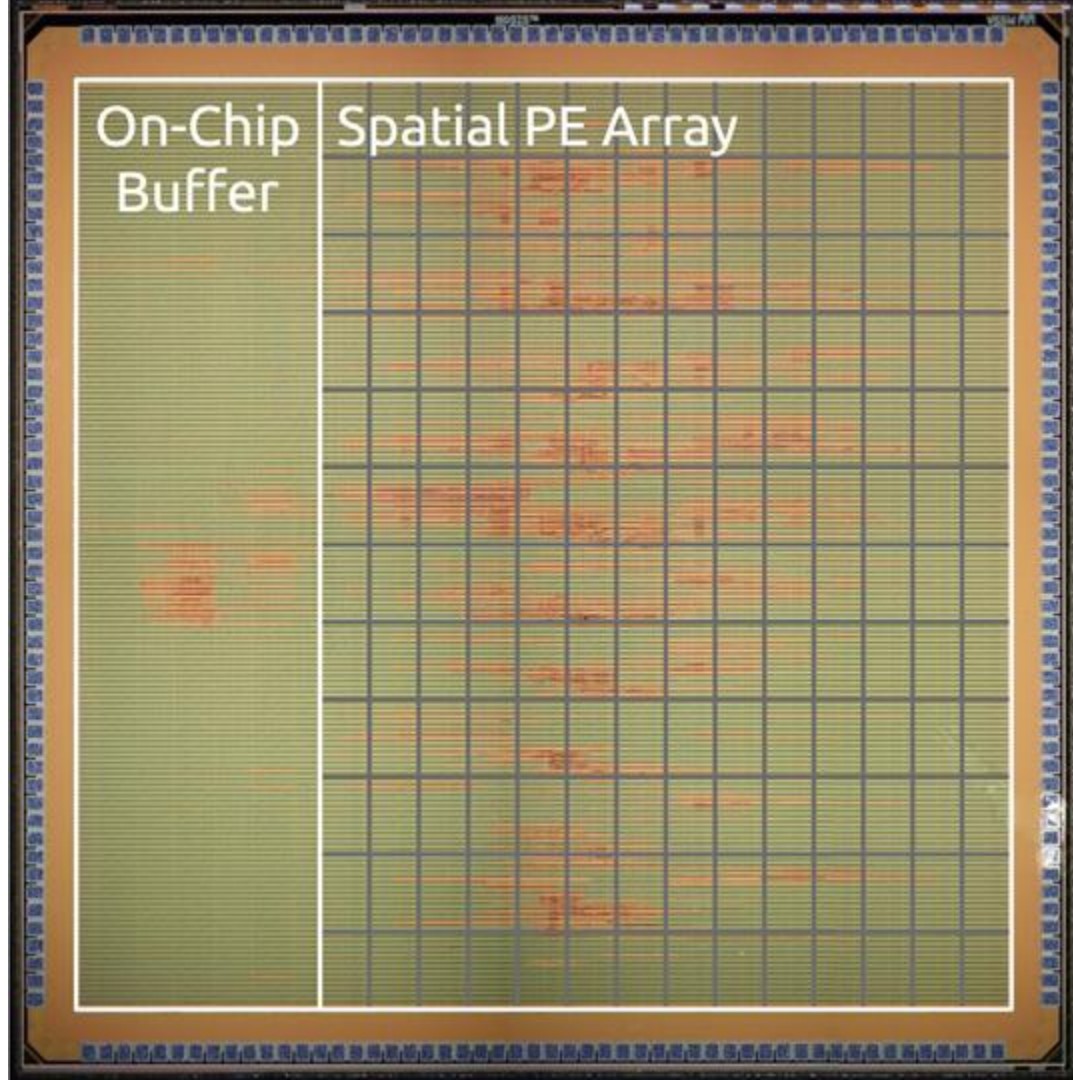


# Eyeriss

Groundbreaking CNN accelerator

- Big on-chip memory
- Custom on-chip network topology for CNNs

Yu-Hsin Chen (MIT), Tushar Krishna (MIT/GTech), Joel Emer (MIT/Nvidia), Vivienne Sze (MIT) at ISCA'16



# Challenges and open problems

# Challenges and open problems

## Amdahl's Law

- To scale,  $\geq 99.9\%$  of the program must run on the accelerator (deja vu...)
- How many computations are worth accelerating anyway?

## Programming

- How to code & compile to accelerators?
- Is there any hope for programmability?



# Amdahl's Law in the Accelerator Era

All the same problems as multicore, plus...

How many accelerators do we need?

- Rapidly diminishing returns
- Large up-front engineering costs

Are accelerators a scaling strategy or a one-off boost?

# Programmability vs specialization

How do you write programs for an accelerator?

- OS challenge: How to publish what accelerators are available?
- PL/compiler challenge: How to map code onto an accelerator?

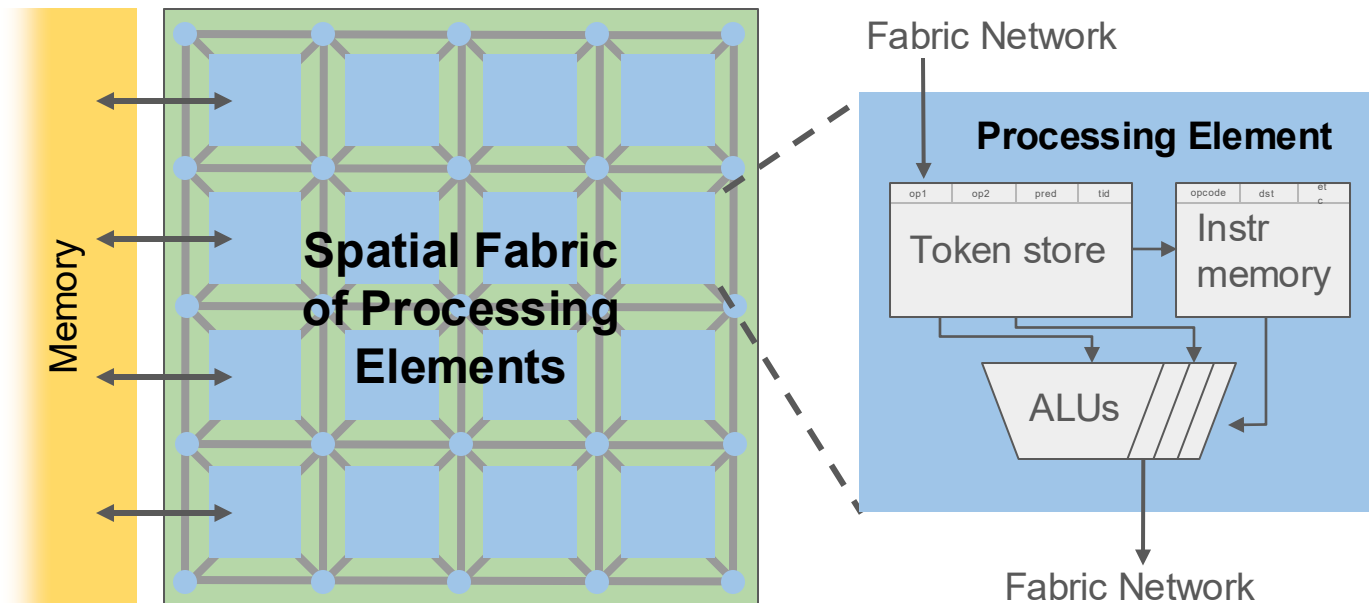
Common approach: Domain-specific languages & frameworks

- Restrict programmers to high-level operations that the framework understands
- Problem: Non-library code becomes Amdahl bottleneck
- Problem: What if an important operation is missing?

# General-purpose acceleration?

Rising interest in programmable, non-von Neumann architectures

E.g., coarse-grained reconfigurable arrays (CGRAs)



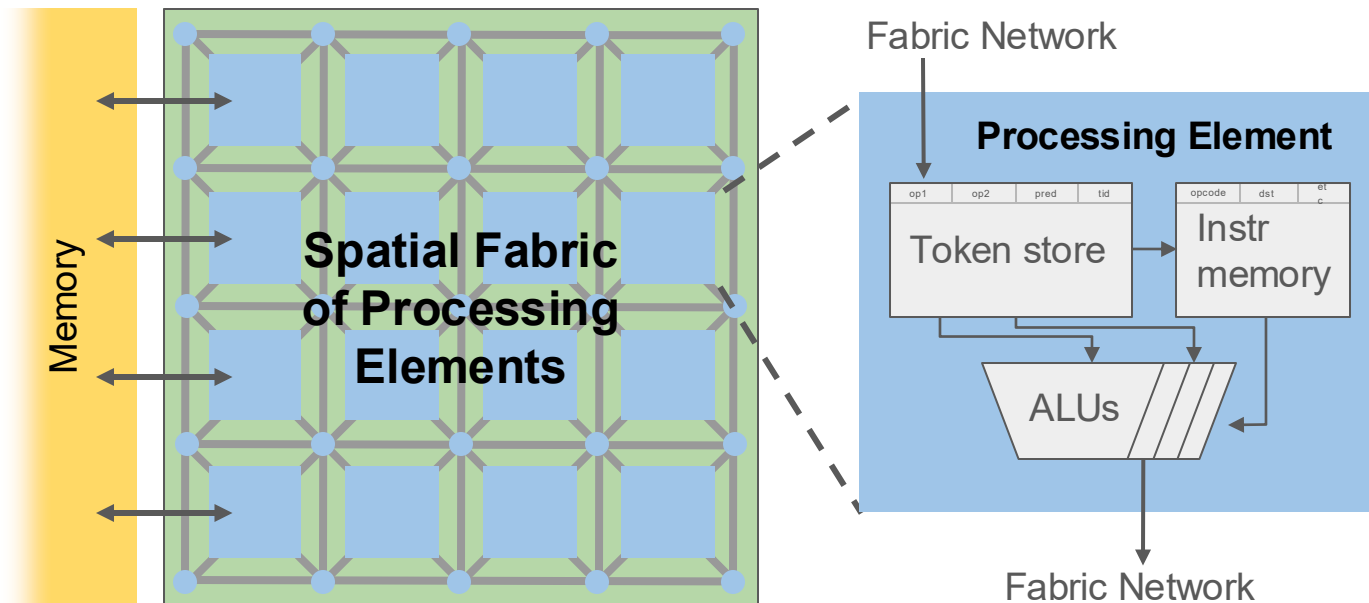
Processor consists of replicated “PEs” (necessary to get efficient scaling)

Compiler maps instructions to PEs, and they stay there!

- Directly encodes dataflow graph
- Can compile from high-level code; much simpler synthesis problem

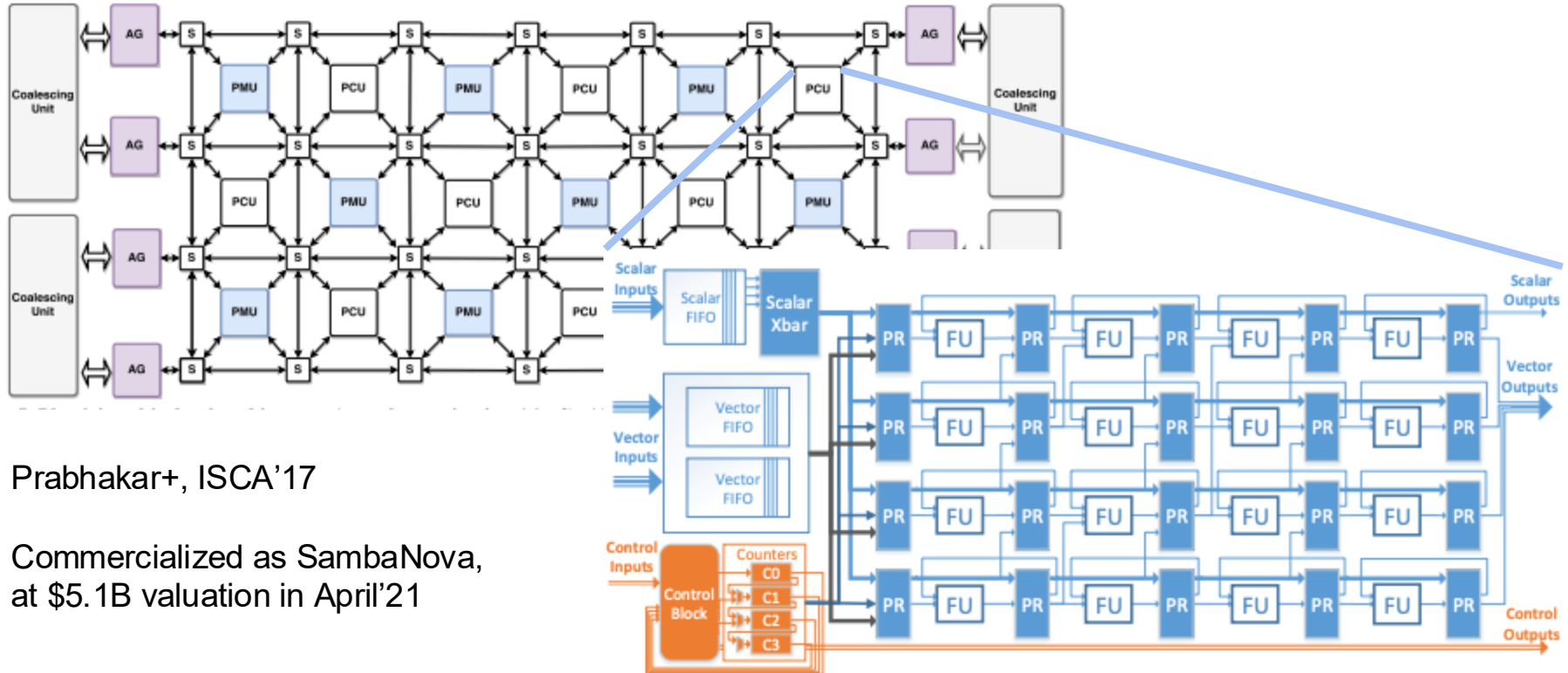
Old idea:

TRIPS, Garp,  
Piperech,  
Dyser, etc





# Example CGRA: Plasticine



Prabhakar+, ISCA'17

Commercialized as SambaNova,  
at \$5.1B valuation in April'21