Energy, power and other trends

Topics we've covered

- Data representations and digital computation
- Basic CPU execution
- Pipelined execution
- Caches
- Virtual memory
- Dynamic ILP
- Static ILP/Compiler considerations
- DLP
- GPUs
- ISA design
- Security

What sorts of tradeoffs (implicit or explicit) have we seen?

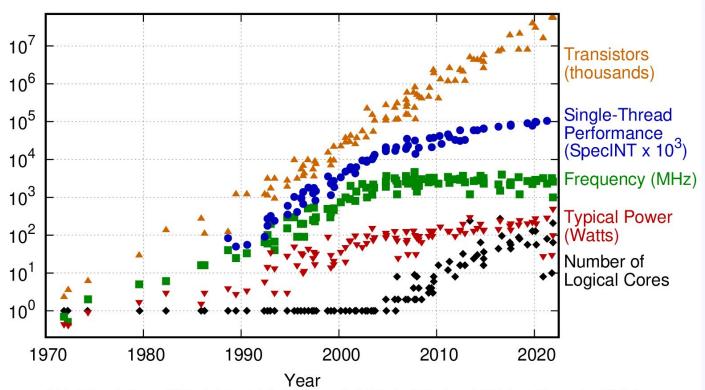
Ex: small/fast/expensive vs. large/slow/cheap memory

Ex: developer effort vs. program efficiency



Source

50 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 2010-2021 by K. Rupp

Top500 performance



source



source

<u>Green500 data</u>

Green500 performance

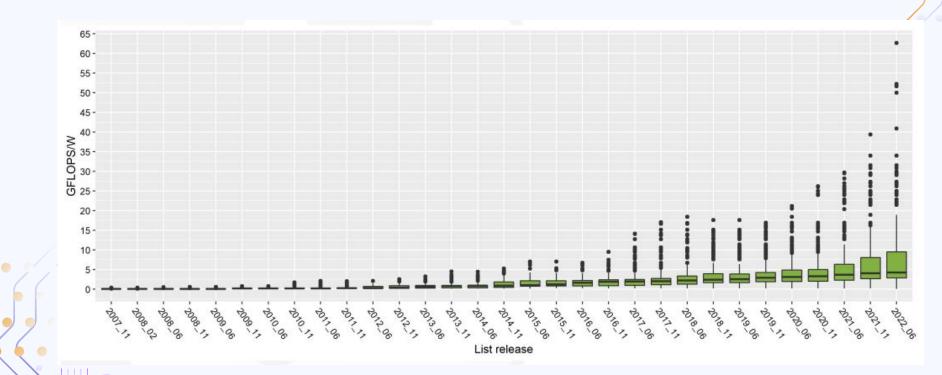


image source

The power wall

Limits to how far we can push a processor (clock frequency, transistor density) based on power limitations such as TDP (Thermal Dissipation Power)

Note: there are many opinions on what the limits of computer architecture are and when we'll reach them. So far, we keep seeing growth in one area or another



Adjusting clock frequency

TASK slack time

What is the energy expenditure compared to above (assuming voltage stays constant)?

If CPU gets put in sleep mode for slack time, we halve the energy usage in this time period

TASK (at half frequency)

But note that reducing frequency *also* allows us to reduce voltage!

This allows us to save energy (*why?*)

DVFS

Dynamic voltage and frequency scaling

Adjusts voltage/frequency based on workload

Modern systems manage this at the OS level

(CPUFreq governors on Linux)

Based on coarse samples of system performance

Require hardware interface

Other approaches: "Dark Silicon"

Low-power modes when hardware isn't being used

For example: I/O interrupts instead of polling

But there is a latency cost to coming out of low-power mode

Clock gating

Turn off clock to idle unit (reduces useless switching)

IBM Power5 claim: 25% reduction in switching power w/o reduction in perf

Detecting narrow-width operands

If buses are width 64 but using only 32 bits, disable those wires



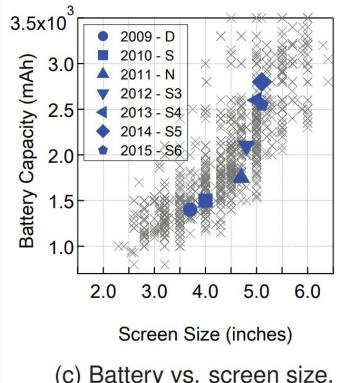
Mobile device challenges

"No Moore's Law for Batteries"

Unless battery technology sees drastic innovation, focus needs to be on energy efficiency

(Good news: mobile devices have different usage patterns than laptops/desktops)

M. Halpern, Y. Zhu and V. J. Reddi, "Mobile CPU's rise to power: Quantifying the impact of generational mobile CPU design trends on performance, energy, and user satisfaction," 2016 IEEE International Symposium on High Performance Computer Architecture (HPCA), Barcelona, Spain, 2016, pp. 64-76, doi: 10.1109/HPCA.2016.7446054. **IEEE link**



(c) Battery vs. screen size.

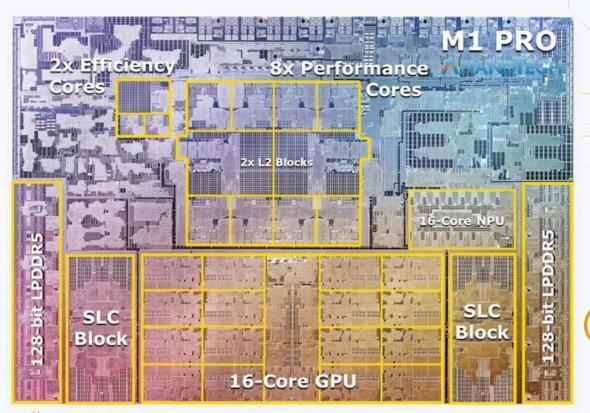
Heterogeneous architectures

Image source

Use different types of processors (also called asymmetric) or processing units for different tasks

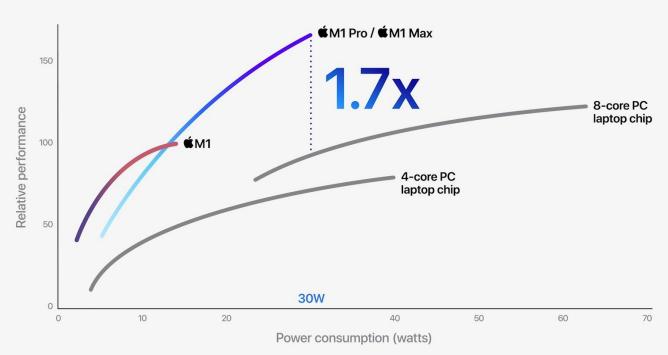
Allows system to use efficient hardware for specific applications

Challenges: coordination, scheduling, design



Claims by Apple, Inc (taken from AnandTech)

CPU performance vs. power



Domain-Specific Architecture

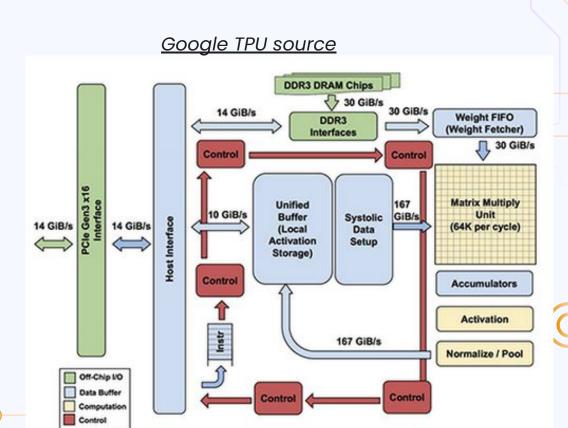
Specialized hardware for software domain

Can adapt memory, precision, parallelism to application

Increase both performance and energy efficiency!

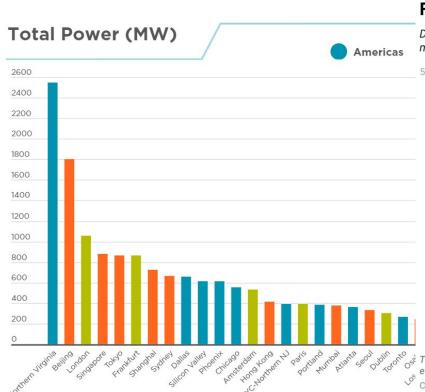
When/how to justify design cost for specialized hardware?

What implications does this have on longevity?



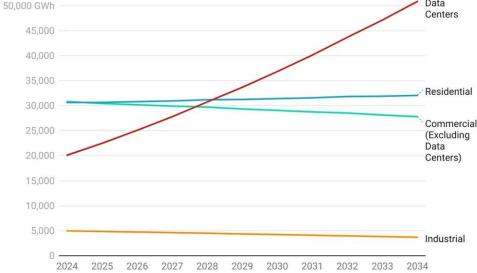
The other concerns with power

source



Forecasted Dominion Energy annual electricity sales

Data center electric sales will increase by 152% in the next decade, while others sectors remain mostly the same.



The overall increase in electricity sales is forecasted to be 32% over 10 years. That accounts for increased energy efficiency among other sectors. The forecast does not include projected electricity demand from electric vehicles.

Chart: Emily Richardson/VCU Capital News Service • Source: The Energy Transition Initiative at the Weldon Cooper Center for Public Service. • Created with Datawrapper

Related reading/viewing

<u>Jim Keller: Moore's Law is Not Dead</u> (argues that, whenever anyone says that some technology has reached its limits and cannot scale, people find ways to come up with a new technology)

<u>Charles E. Leiserson et al.</u>, There's plenty of room at the Top: What will drive computer performance after Moore's law? Science 368, eaam9744(2020).

<u>Hennessy, John L., and David A. Patterson.</u> "A new golden age for computer architecture." Communications of the ACM 62.2 (2019): 48-60.

<u>Kaxiras, Stefanos, and Margaret Martonosi</u>. Computer architecture techniques for power-efficiency. Morgan & Claypool, 2008.

A. Reuther, P. Michaleas, M. Jones, V. Gadepally, S. Samsi and J. Kepner, "Al and ML Accelerator Survey and Trends," 2022 IEEE High Performance Extreme Computing Conference (HPEC), Waltham, MA, USA, 2022, pp. 1-10, doi: 10.1109/HPEC55821.2022.9926331.