# The "Turing Tax" - reducing the interpretive bottleneck in computer architecture

## UK Design Forum, Manchester, June 2025

Paul H J Kelly

Group Leader, Software Performance Optimisation
Department of Computing, Imperial College London

**This talk includes work done by or influenced by:** David Ham (Imperial Maths), Andy Davison (Imperial Computing), Lawrence Mitchell (NVIDIA),Gerard Gorman (Imperial Earth Science Engineering), Peter Vincent (Imperial Aeronautics), Wayne Luk (Imperial Computing), Piotr Dudek (Manchester), Jim Whittaker (Huawei), Vassilios Chouliaras (Huawei), Martin Berger (Montanarius), Sam Ainsworth (Edinburgh) and more

**And affiliated postdocs/PhD students:** Luke Panayi, Jacky Wong, Nicholas Fry, Shinjeong Kim, Edward Stow, Yuncheng Lu, Shuang Liang, Riku Murai, George Bisbas

**And by group alumni:** Eduardo Carvalho (Valantis Labs), Marius Koch (NVIDIA), Navjot Kukreja (Senapt), Philippos Papaphilippou (U. Southampton), Edward Stow (TBC), Matthew Taylor, Sophia Vorderwuelbecke, Fabio Luporini (DevitoCodes), Graham Markall (NVIDIA), Florian Rathgeber (Google), Francis Russell (Quaisr), George Rokos (Intel), Tianjiao Sun (Orbit Markets), Thanasis Konstantinides (Cerebras), Carlo Bertolli (AMD), Doru Bercea (AMD), Michael Lange (ECMWF), Sajad Saeedi (UCL), Luigi Nardi (DBTune/Lund University), Navjot Kukreja (Senapt), Emanuele Vespa (Magic Leap), Miklos Homolya and more

**And by collaborators:** Mike Giles (Oxford), Gihan Mudalige (Warwick), Istvan Reguly (Pazmany Peter, Hungary), Matt Piggott (Imperial ESE), Spencer Sherwin, Chris Cantwell (Imperial Aero), Michelle Mills Strout (University of Arizona/HPE), Chris Krieger (Maryland), Cathie Olschanowsky (Colorado State University), Bruno Bodin (Yale-NUS), Richard Veras, Ram Ramanujam (Louisiana State University), Doru Thom Popovici (LBL), Franz Franchetti (CMU), Jan Hückelheim (Argonne), Freddie Witherden (Texas A&M), Chris–Kriton Skylaris (Southampton) Ridgway Scott (U Chicago), Lluis Guasch (Imperial Earth Science Engineering)

# Feynmann: plenty of room at the bottom

"…I do know that computing machines are very large; they fill rooms. Why can't we make them very small, make them of little wires, little elements – and by little, I mean little? For instance, the wires should be 10 or 100 atoms in diameter, and the circuits should be a few thousand angstroms across"

- >60 years of exponential progress since then
- We're much closer to such limits
- Much debate about where they really lie
- What is clear is that we're a lot closer
- We are confronted more and more with fundamental physical concerns
- Particularly wrt communication latency, bandwidth and energy.

**(1959, talk at the American Physical Society)**

Ca.10"

Ca.8.5"

$1.2 \times 10^{12}$ *transistors*

*400k cores*

*18GB SRAM*

*17-20KW TDP*

■ Moore's Law: "circuit density doubles every 18 months"

■ 60 years =40x18months

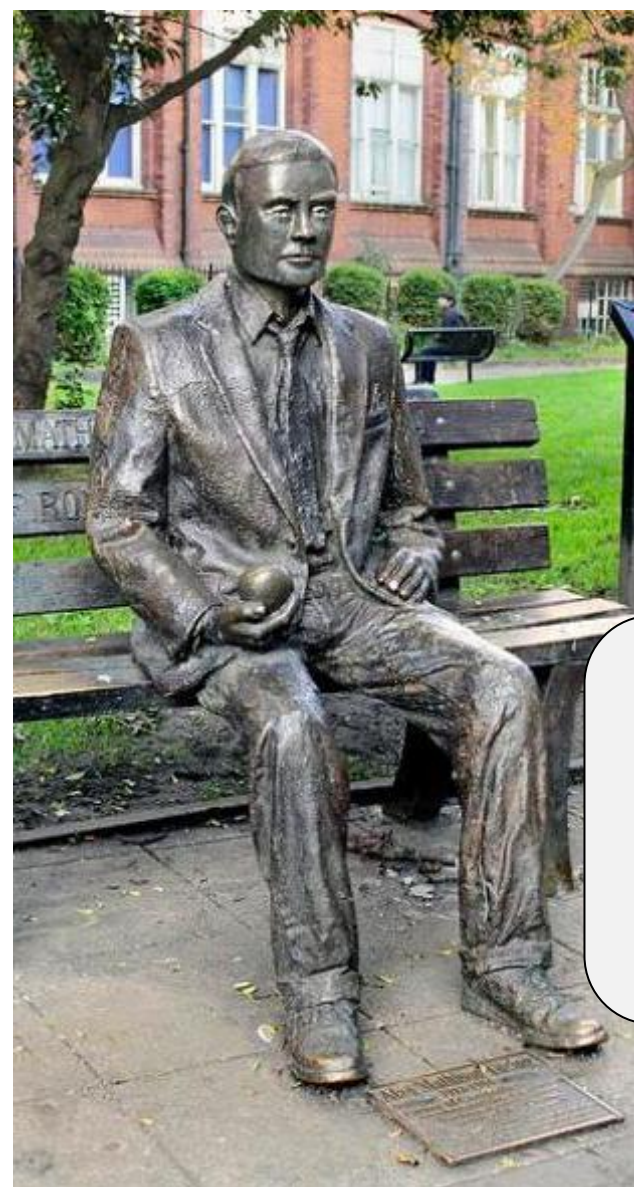■ So Moore's Law would predict $2^{40}$= $10^{12}$ increase

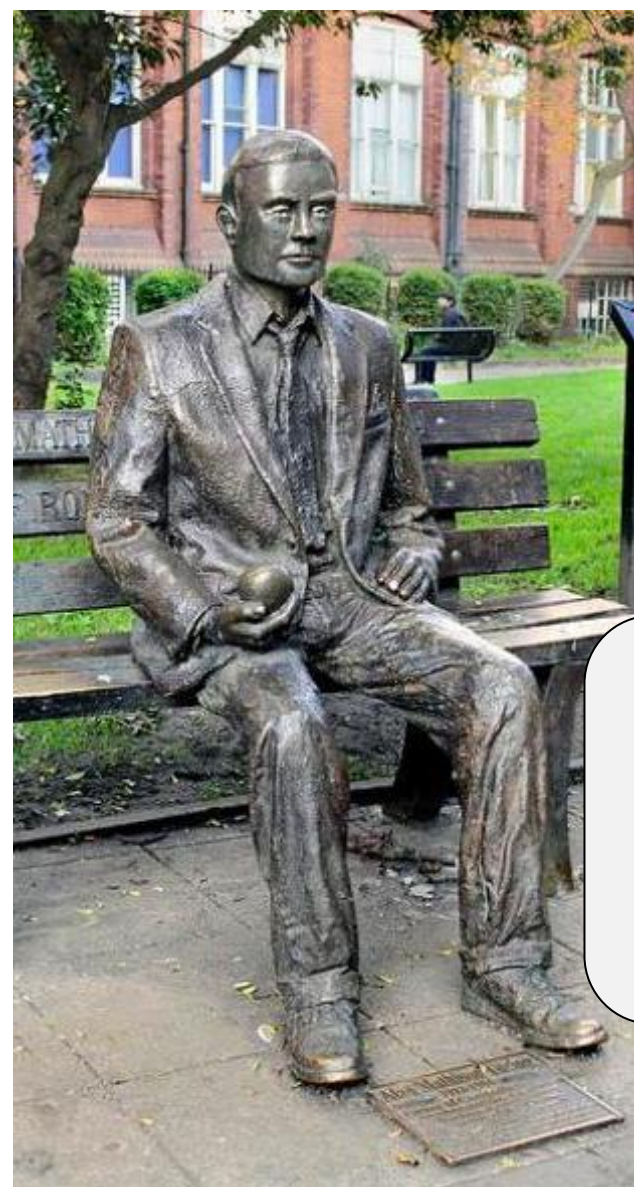Ferranti Computer Systems Ltd Pegasus valve computer circuit board, c. 1964
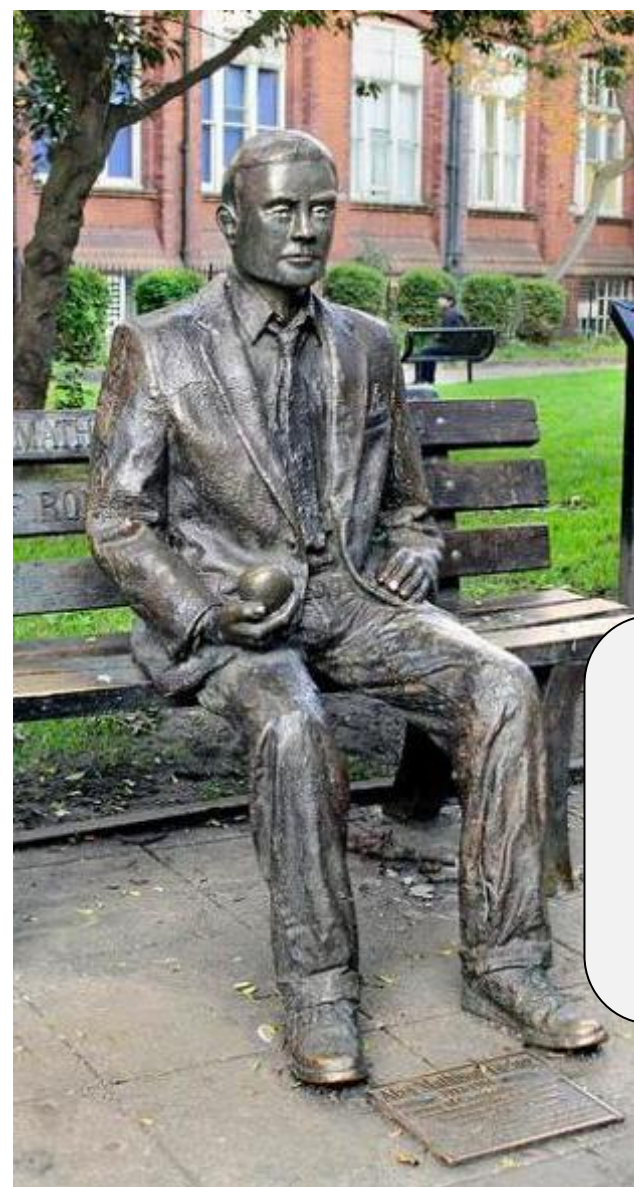*https://blog.sciencemuseum.org.uk/the-pegasus-computer/*

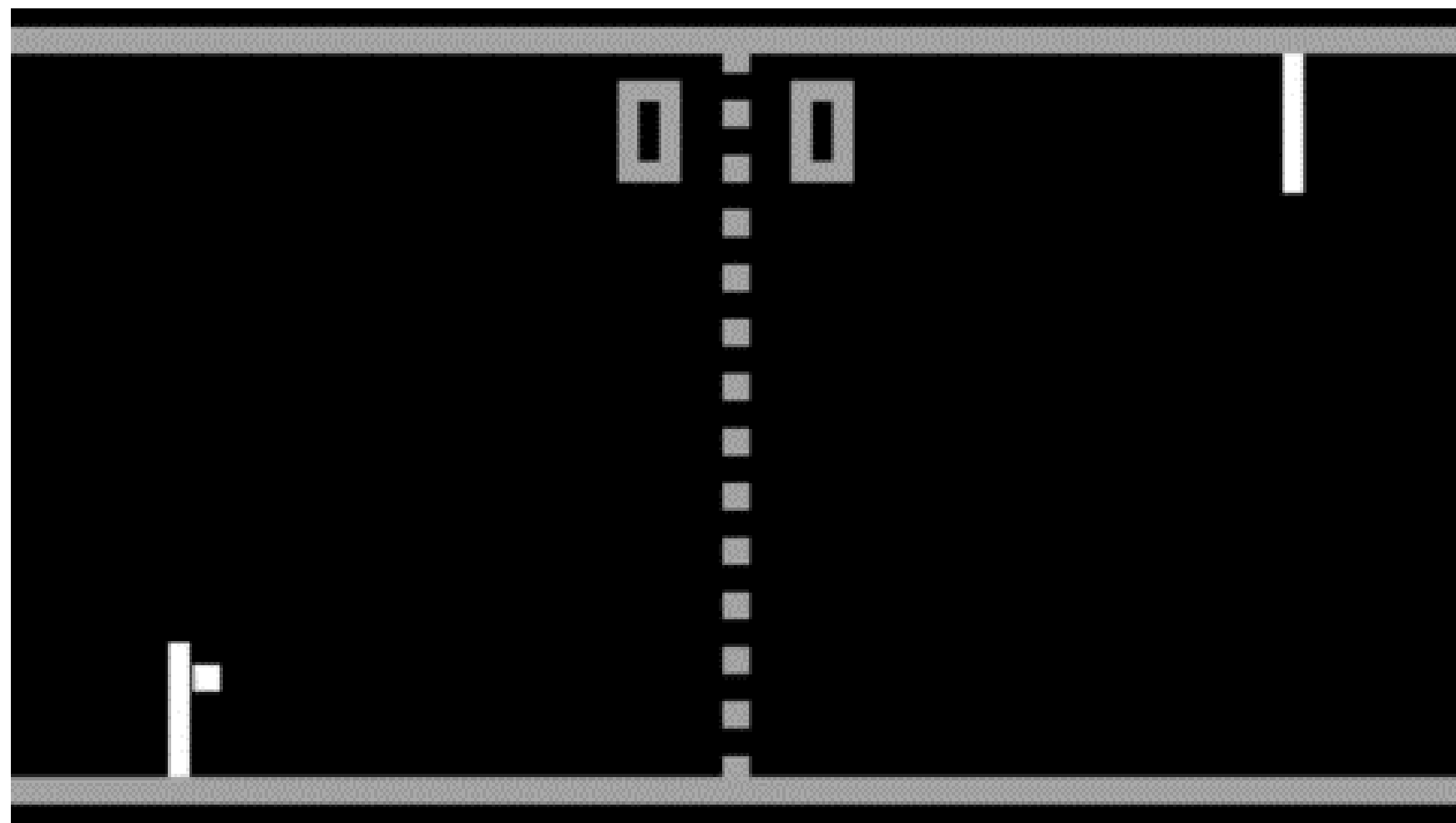Cerebras co-founder Sean Lie holding the Wafer Scale Engine. Image: Cerebras Systems

- Alan Turing realised we could use digital technology to implement any computable function

- He then proposed the idea of a "universal" computing device – a *single* device which, with the right program, can implement any computable function *without further configuration*

- "**Turing Tax**", **or** "**Turing Tariffs**": the overhead (performance, cost, or energy) of universality in this sense

- The performance (time/area/energy) difference between a **special-purpose** device and a **general-purpose** one

- **One of the fundamental questions of computer architecture is to how to reduce the Turing Tax**

**Imperial College London**

- Alan Turing realised we could use digital technology to implement any computable function

- He then proposed the idea of a "universal" computing device – a *single* device which, with the right program, can implement any computable function *without further configuration*

- **"Turing Tax", or "Turing Tariffs"**: the overhead (performance, cost, or energy) of universality in this sense

- The performance (time/area/energy) difference between a **special-purpose** device and a **general-purpose** one

- **One of the fundamental questions of computer architecture is to how to evade the Turing Tax**

- Alan Turing realised we could use digital technology to implement any computable function

- He then proposed the idea of a "universal" computing device – a *single* device which, with the right program, can implement any computable function *without further configuration*

- **"Turing Tax", or "Turing Tariffs"**: the overhead (performance, cost, or energy) of universality in this sense

- The performance (time/area/energy) difference between a **special-purpose** device and a **general-purpose** one

- **One of the fundamental questions of computer architecture is to how to reduce the Turing Tariffs**
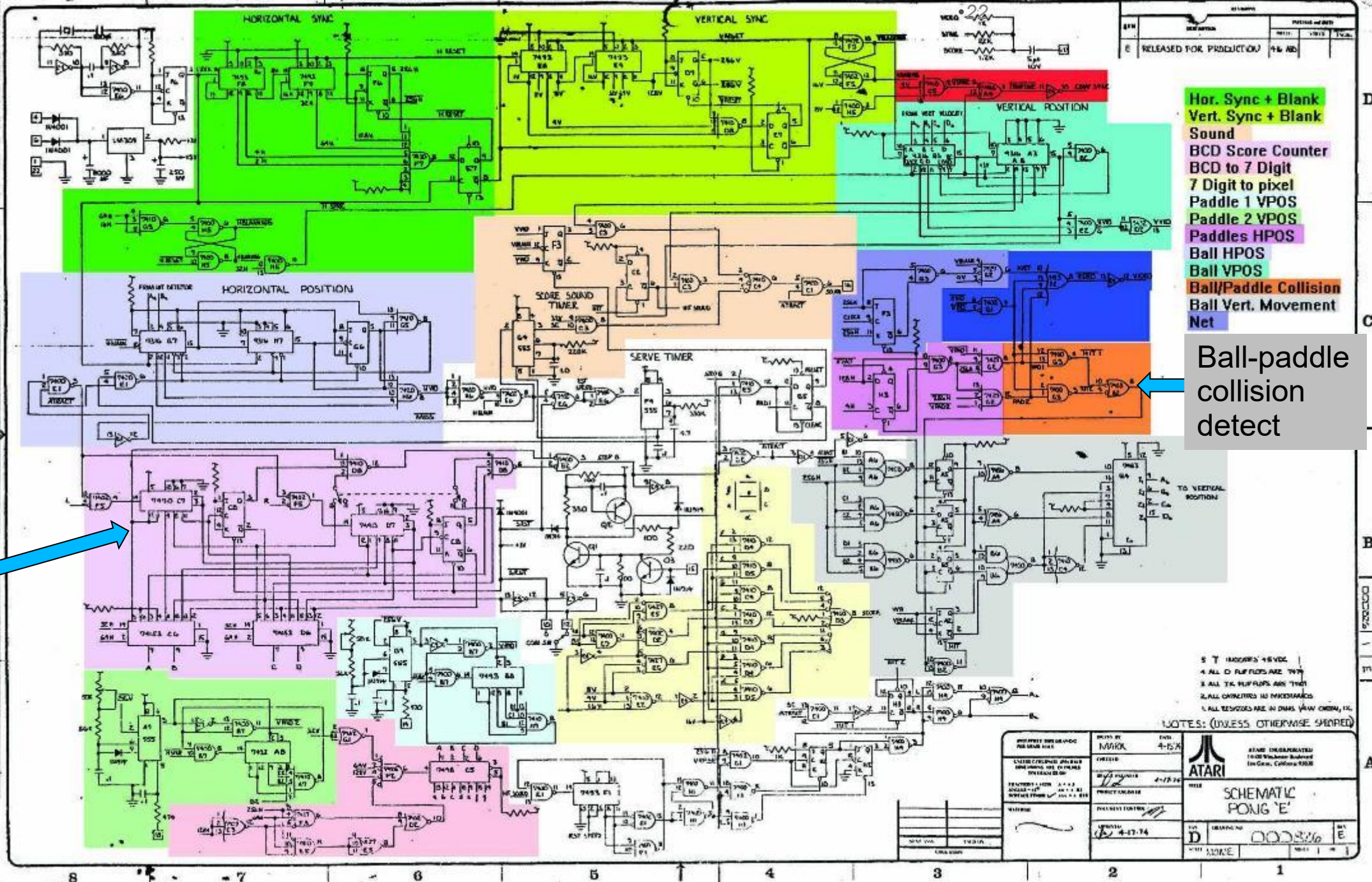
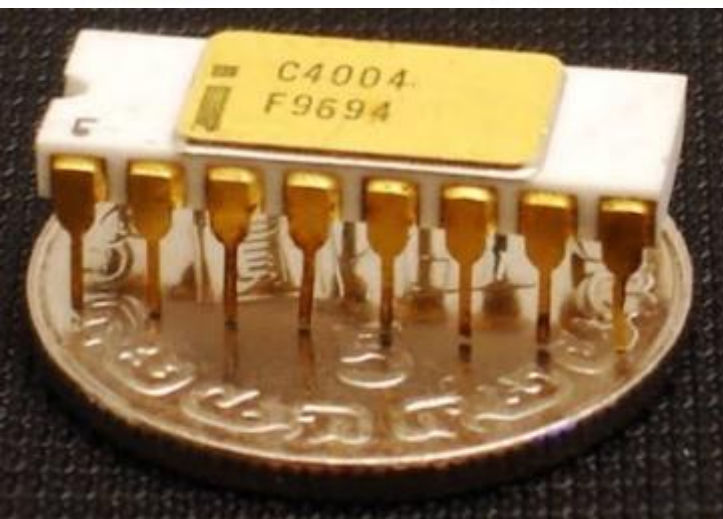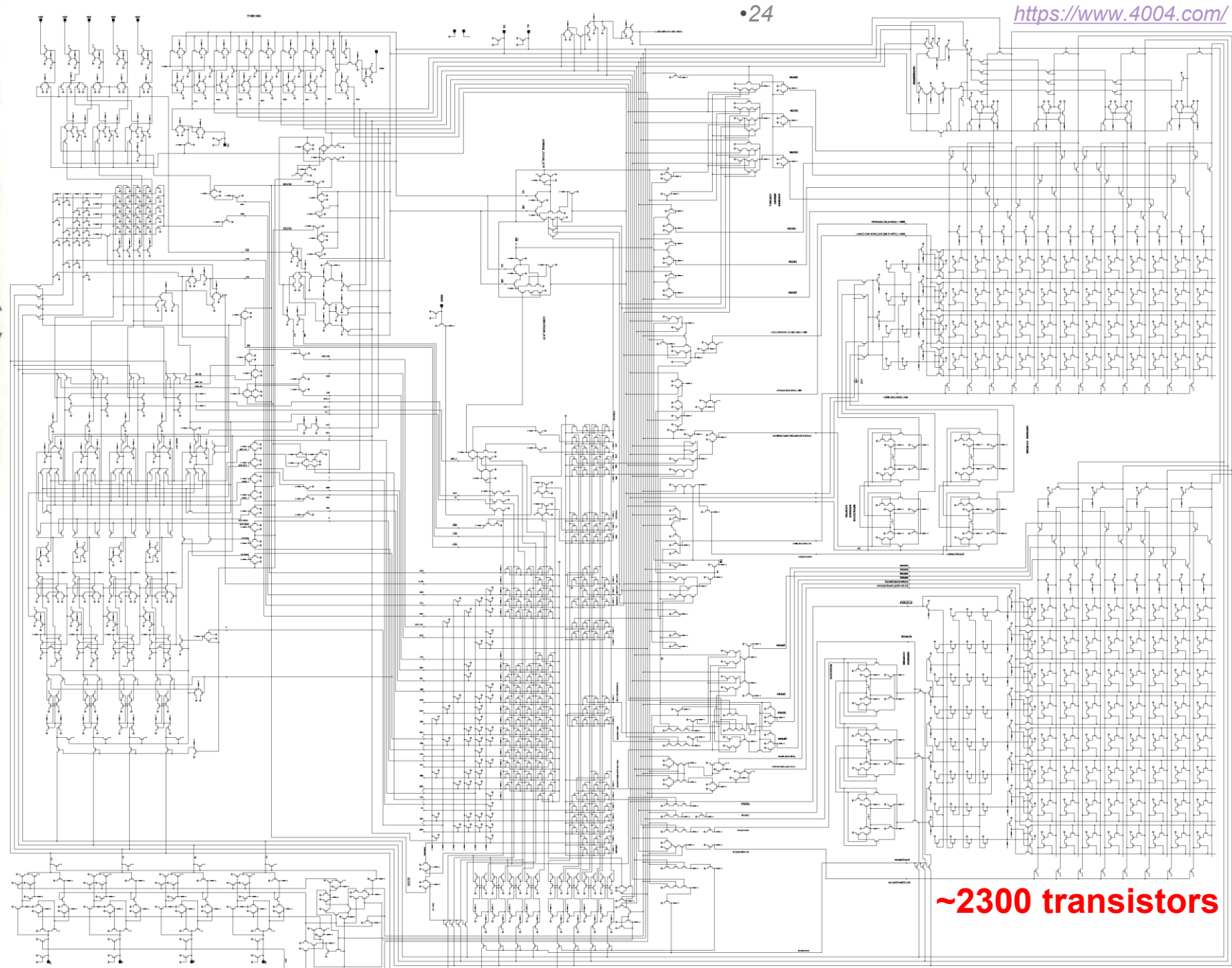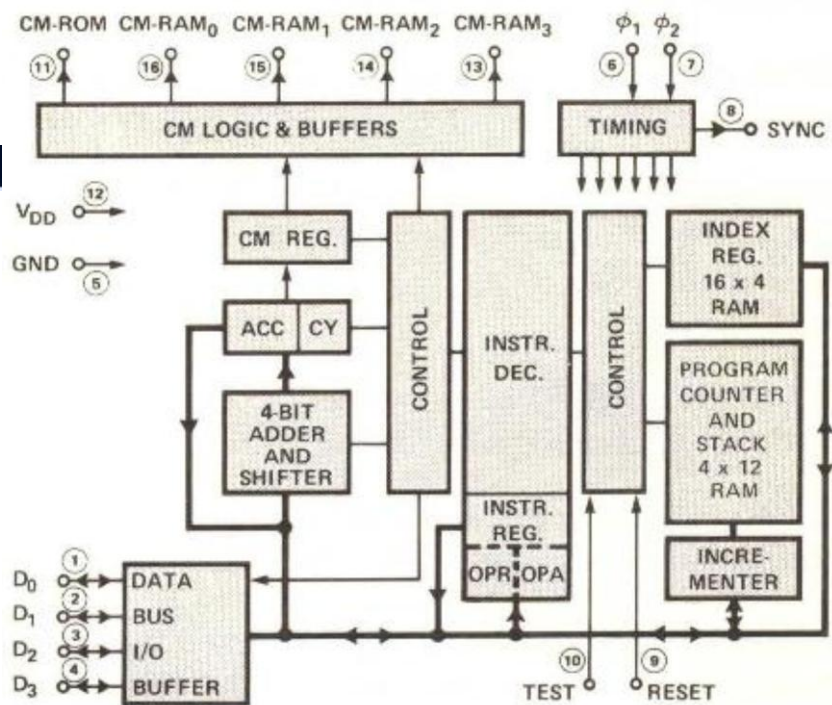# Circuit diagram for Atari's Pong game

Designed by Allan Alcorn in around 1972

Alcorn hired Steve Jobs in 1974

7490 TTL IC: "decade counter" for tracking the score (ca.45 transistors)

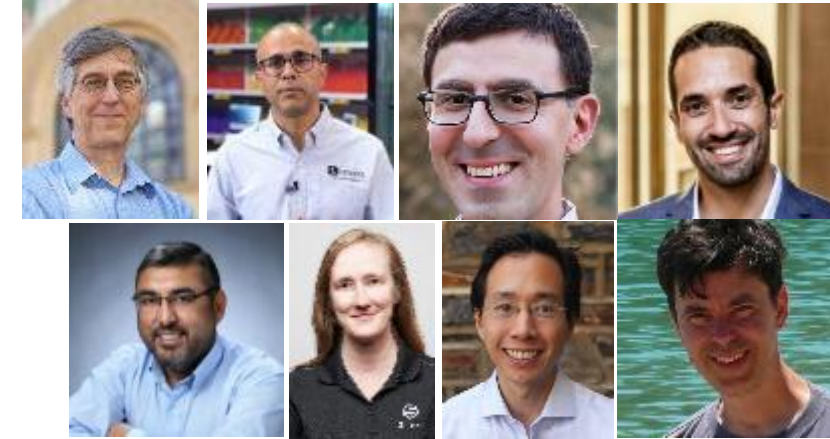https://www.pong-story.com/LAWN_TENNIS.pdf

Ball-paddle collision detect

**Legend:**
- Hor. Sync + Blank
- Vert. Sync + Blank
- Sound
- BCD Score Counter
- BCD to 7 Digit
- 7 Digit to pixel
- Paddle 1 VPOS
- Paddle 2 VPOS
- Paddles HPOS
- Ball HPOS
- Ball VPOS
- Ball/Paddle Collision
- Ball Vert. Movement
- Net

~2300 transistors

*24  Circuit diagram for the first commercially-available microprocessor, Intel's 4004 (1971)*
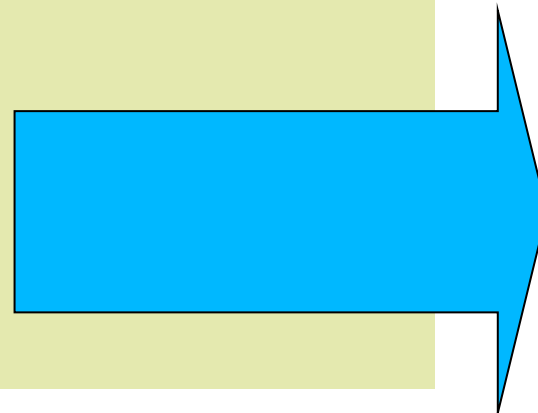
**So how big is the Turing Tax?**

Suppose you have a job that you can do on a single CPU core

But you invest in a fully-specialised ASIC instead

**How much faster?**

**How much smaller?**

**How much less energy?**

- 2010 article: Hameed, Qadeer, Wachs, Azizi, Solomatnikov, Lee, Richardson, Kozyrakis & Horowitz: **Understanding sources of inefficiency in general-purpose chips**. ISCA'10

- Intel's hand-coded implementation of H.264 encoding for Pentium 4 for 1280x720 HD:
  - 11fps
  - $122mm^2$
  - 2023mJ/frame

- ASIC implementation of H.264 encoding for 1280x720 HD:
  - 30fps
  - $8mm^2$
  - 4mJ/frame

(Chen, Chien, Huang, Tsai, Chen, Chen, & Chen: Analysis and architecture design of an HDTV720p 30 frames/s H.264/AVC encoder. IEEE Trans. Cir. and Sys. for Video Technol. 16, 6 (September 2006))

27

# Turing tariffs

- Fetch-execute is the original Turing tariff

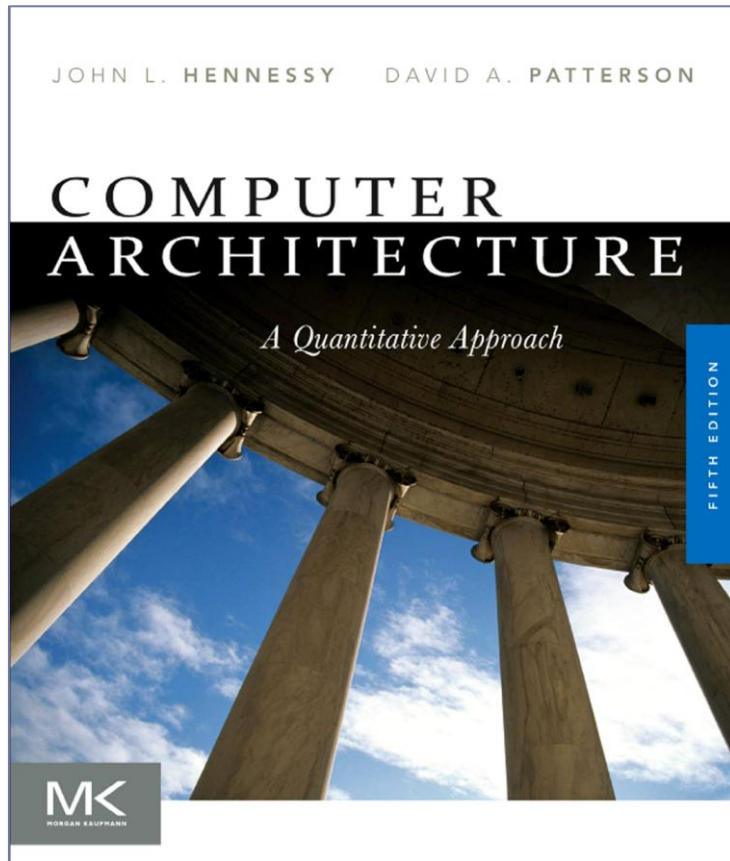- FPGAs pay Turing tariffs in the reconfigurable fabric

- Registers are a Turing Tariff
  - Because if we know the program's dataflow, we can use wires and latches to pass data from functional unit to functional unit

- Memory
  - But if we can stream data from where it's produced to where it's used, maybe we don't need so much RAM?

- Cache
  - If we know exactly when the reuse will occur, we can program movement to and from local fast memory explicitly

- Floating-point arithmetic:
  - If we know the dynamic range of expected values…

- Fetch-execute, decode

- Registers, forwarding

- Dynamic instruction scheduling, cracking, packing, renaming

- Cache tags

- Cache blocks

*Basically the whole computer architecture textbook*

- Cache coherency

- Prefetching

- Branch prediction

- Speculative execution

- Address translation

- Store-to-load forwarding, write combining, address decoding, ECC, DRAM refresh

- Mis-provisioning: unused bandwidth, unusable FLOPs, under-used accelerators

JOHN L. HENNESSY    DAVID A. PATTERSON

COMPUTER ARCHITECTURE

*A Quantitative Approach*

FIFTH EDITION

MK
MORGAN KAUFMANN

- SIMD: amortise fetch-execute over a vector or matrix of operands
- Predication, VLIW, EPIC, register rotation
- Macro-instructions: FMA, crypto, conflict-detect, custom ISAs

- Streaming dataflow: FPGAs, CGRAs, systolic arrays
- Circuit switching instead of packet switching
- DMA, long cache lines: move a lot of data with one instruction

- Non-temporal loads/stores, explicit prefetch instructions
- Scratchpads ("shared memory" in CUDA)
- Message passing, instead of cache-coherent shared memory

- Texture caches, interpolation, decompression

- Fine-grain multithreading to avoid pipeline hazards, hide latency

- **And many more ideas…. *Your ideas?***

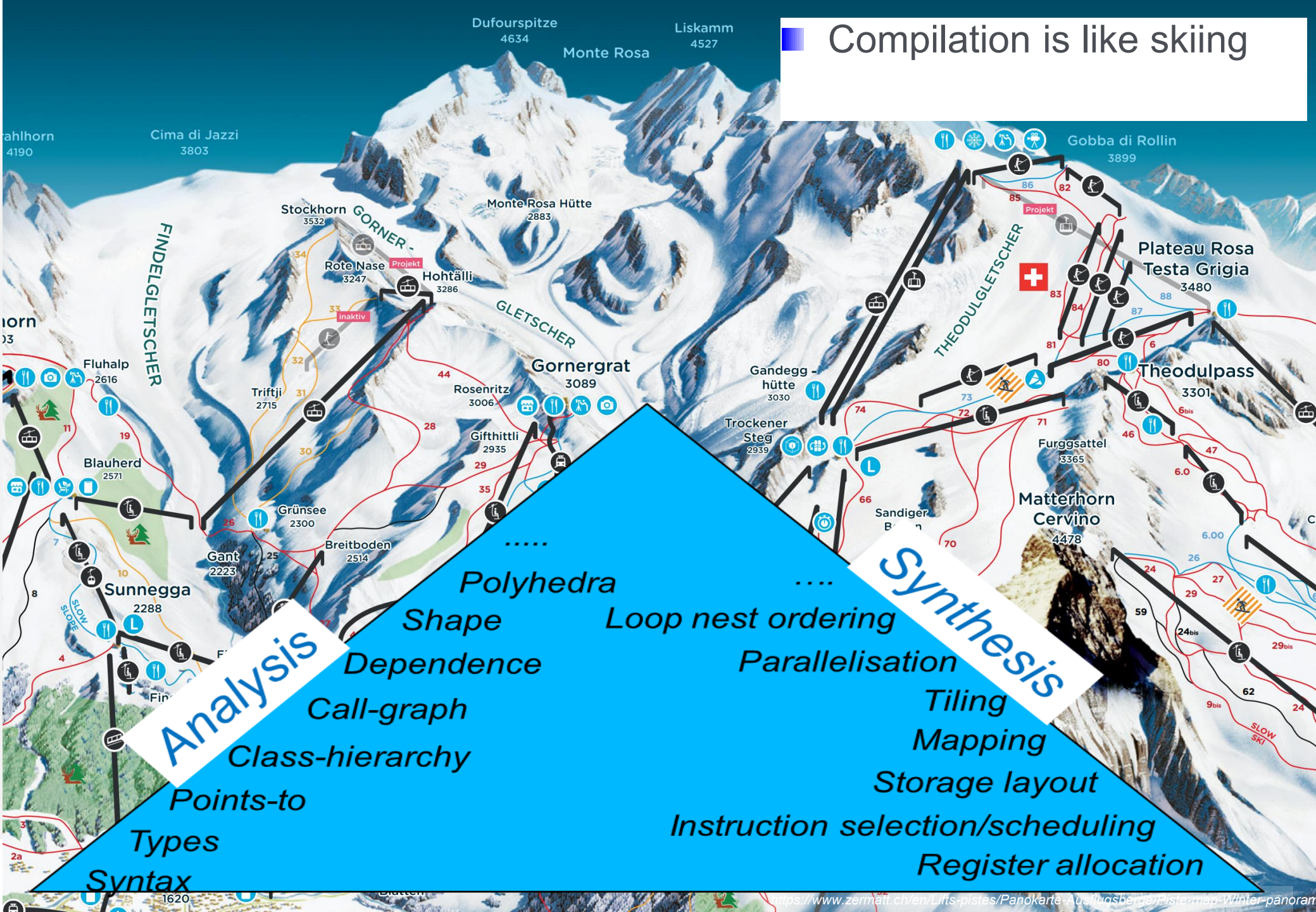Generating code to avoid the need for interpretive mechanisms in hardware:

- Vectorisation
- Static instruction scheduling
- Offloading
- Predication
- Message aggregation
- Synchronisation minimization

Generating code that is specialized for a specific purpose:

- Function inlining, type disambiguation, object inlining
- Specialisation: metaprogramming, JIT, metatracing

Compilation is like skiing

What about **compilers**?

**Is there a Turing Tax for compilers too?**

Analysis

Syntax
Types
Points-to
Class-hierarchy
Call-graph
Dependence
Shape
Polyhedra
.....

Synthesis

.....
Loop nest ordering
Parallelisation
Tiling
Mapping
Storage layout
Instruction selection/scheduling
Register allocation

https://www.zermatt.ch/en/Lifts-pistes/Panokarte-Ausflugsberge/Piste-map-Winter-panorama

Compilation is like skiing

What about **compilers?**

**The price you pay for coding in a general-purpose language**

**When you could have used a DSL**

Carrying your skis up the mountain isn't the best bit
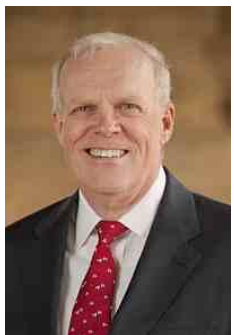
**All of this is Turing Tariff!**

Synthesis

https://www.zermatt.ch/en/Lifts-pistes/Panokarte-Ausflugsberge/Piste-map-Winter-panoram

# Computer architecture – the book

**Computer Architecture: A Quantitative Approach**

**Six editions since 1990**

**Revolutionary landmark book brought experimental discipline to processor design**

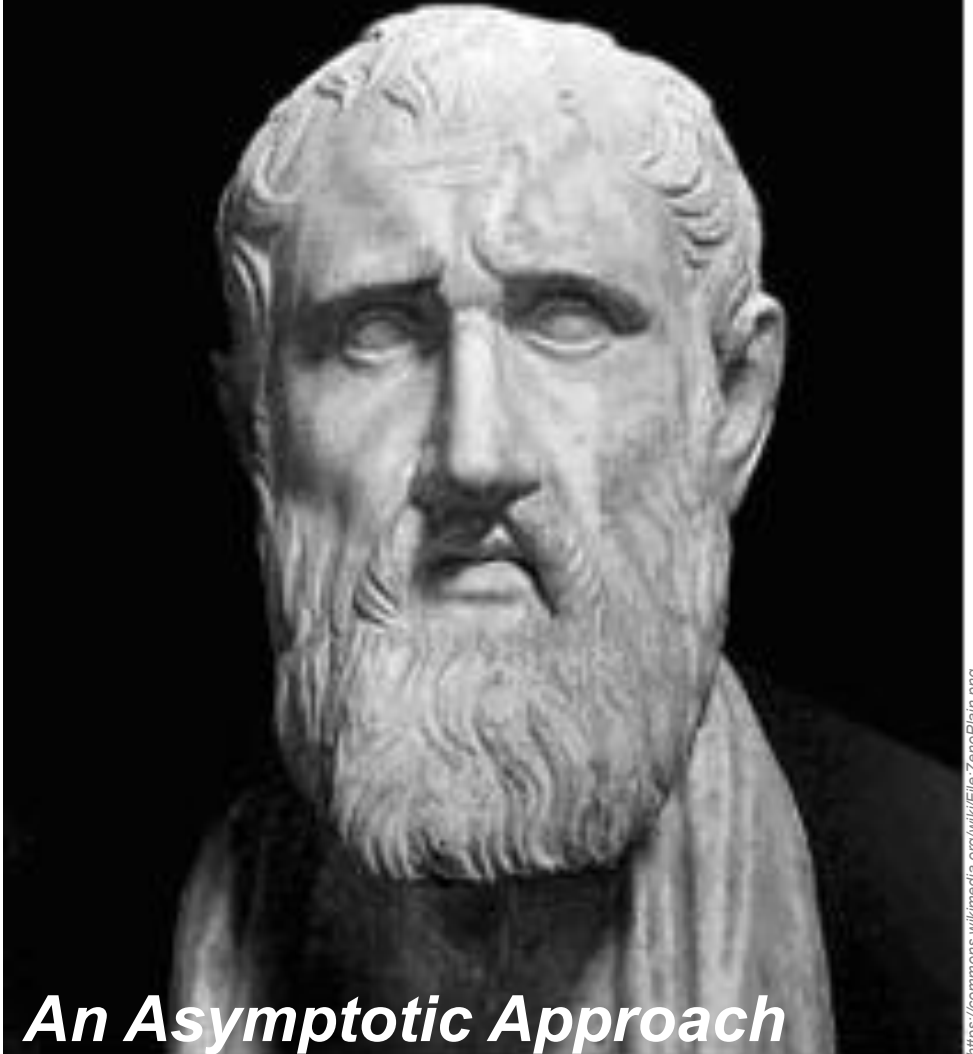**Almost entirely devoid of theory**

*David Patterson*

*John Hennessy*

# Computer architecture – the future?

***Computer Architecture***

- **A manifesto**
- **For computer architecture at the end of Moore's Law**
- **Where we confront fundamental physical constraints**
- **Where we have to account for fundamental costs**
- <span style="color:red">**Where architectural efficiency is paramount**</span>

***An Asymptotic Approach***

*(Statue is of Zeno of Citium, perhaps the patron saint of asymptotics)*

58

https://commons.wikimedia.org/wiki/File:ZenoPlain.png