

Report from the Arch2030 Visioning Workshop: Where are Computer Architects headed and what does it mean for GreenMetrics?

Thomas F. Wenisch

*(Special acknowledgements to Luis Ceze, Mark Hill,
& 40+ members of the architecture community)*

Arch2030 was supported by the Computing Community Consortium (CRA).

The (quick) backstory... (1/2)

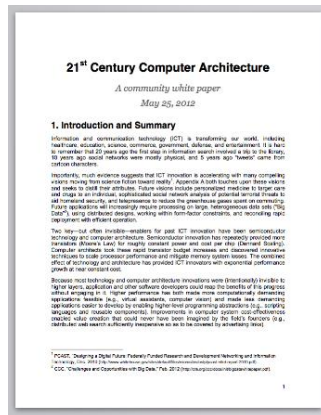
Moore's Law is ending. For real this time.



In 2011, Architects (via the Computing Community Consortium) sent a white paper to NSF on “21st Century Computer Architecture”

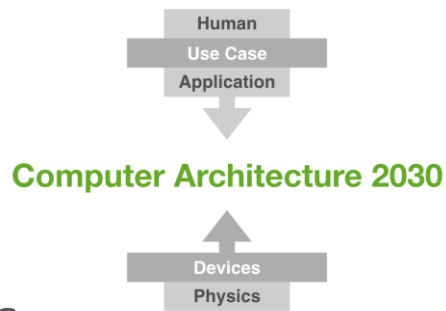
Contributed to launch of NSF eXploiting Parallelism and Scalability program

**But, the world has changed in 5 years;
the community did not foresee some critical trends...**



The (quick) backstory... (2/2)

To update the architecture research vision,
CCC sponsored **Arch2030** workshop with ISCA 2016.



This talk:

- 1) Context: What did architects say in 2011, and what is different now?
- 2) Summarize the technical story & recommendations of these whitepapers
- 3) Opine on where it intersects with GreenMetrics

20th century ICT set up:

Information & Communication Technology (ICT) has changed our world

<long list omitted>

Required innovations in algorithms, applications, programming languages, ... ,
& system software

Key (invisible) enablers to (cost-)performance gains:

Semiconductor technology (“Moore’s Law”)

Computer architecture (~80x per Danowitz et al.)

21st century ICT promises more



Data-centric personalized health care



Computation-driven scientific discovery



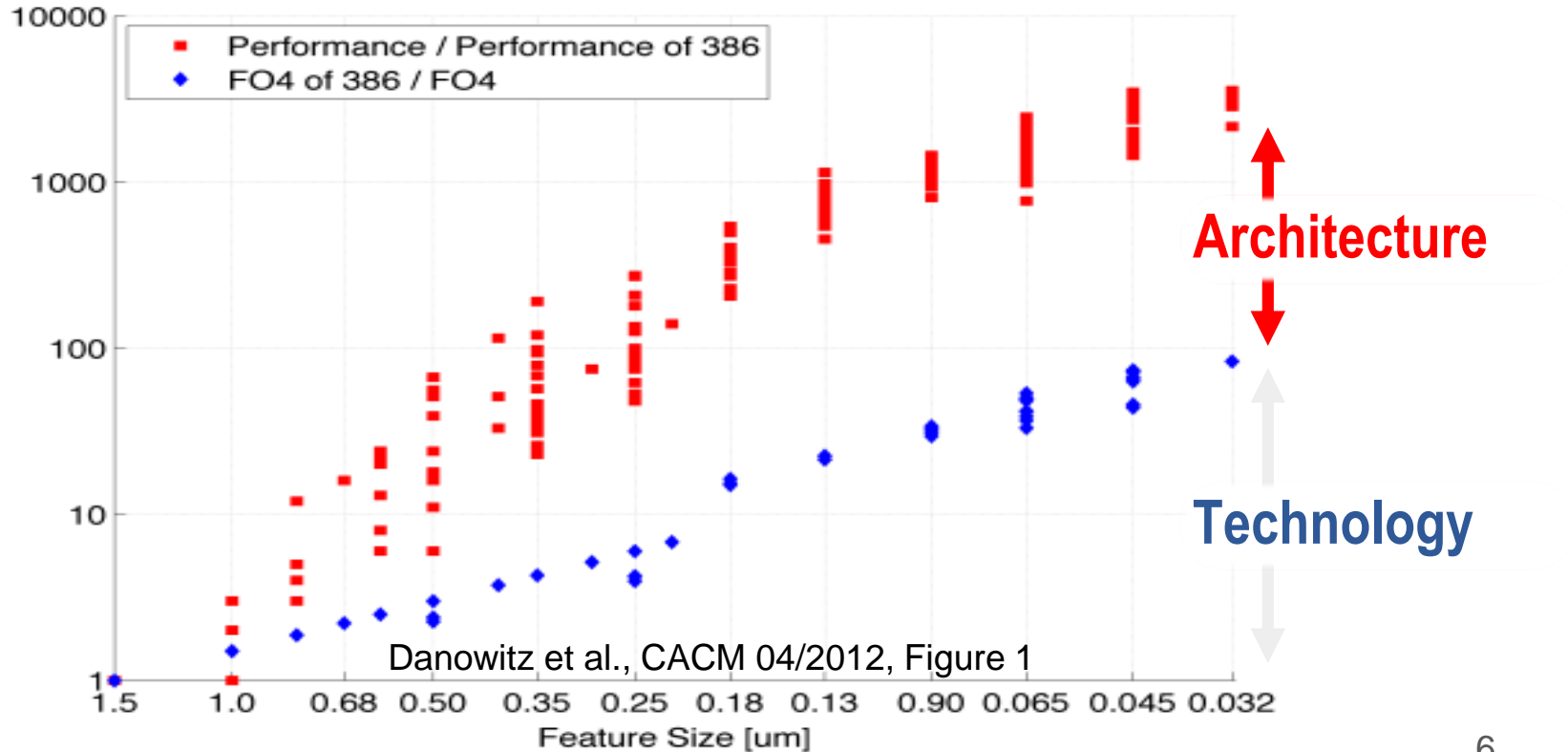
"You never call, and the federal government will back me up on that."

Human network analysis

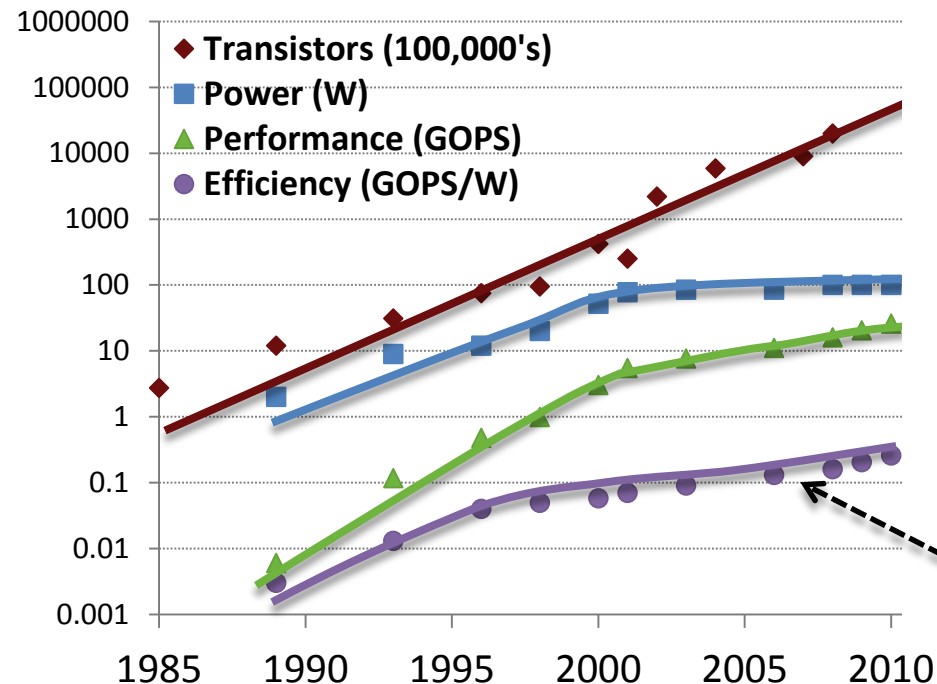


Much more: known & unknown

Enablers: Technology + Architecture



Technology: a paradigm shift in the 2000s...

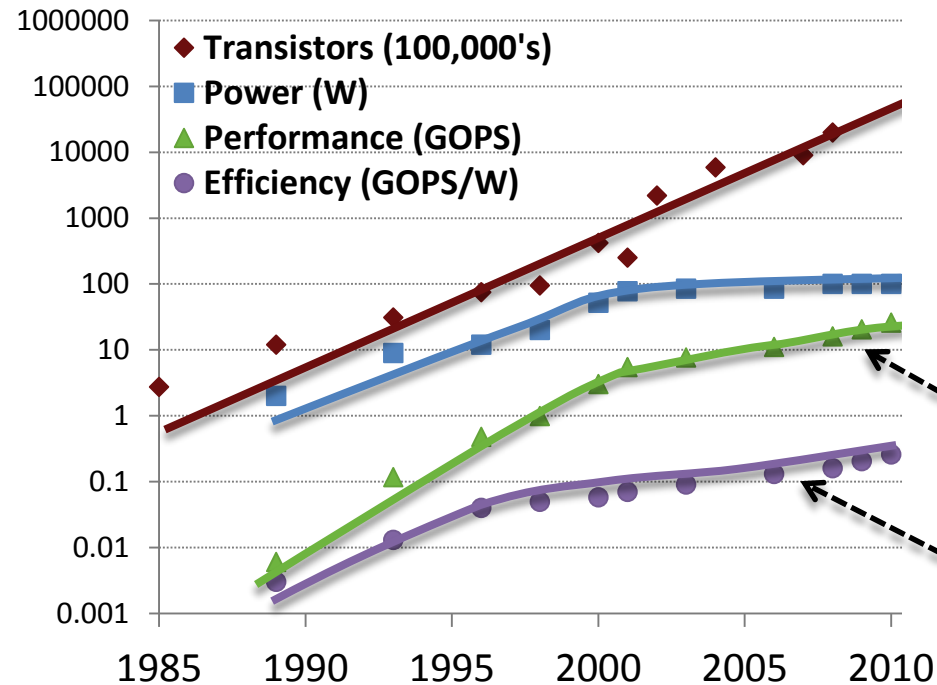


IEEE Computer—April 2001
T. Mudge

Limits on heat extraction

Limits on energy-efficiency of operations

Technology: a paradigm shift in the 2000s...



IEEE Computer—April 2001
T. Mudge

Limits on heat extraction

Stagnates performance growth

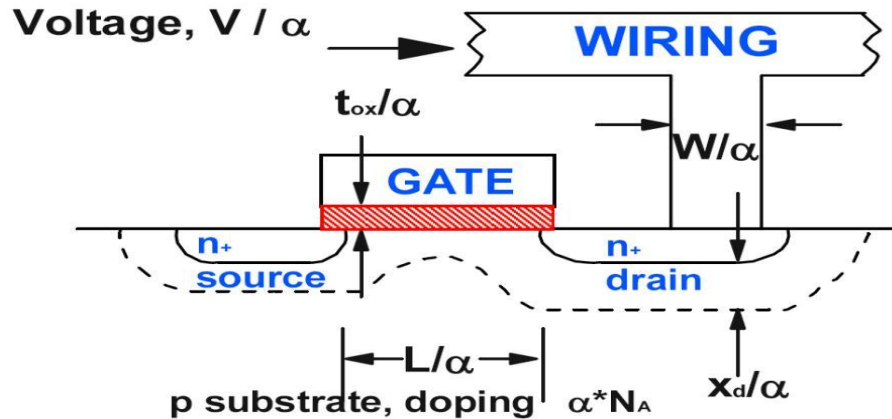
Limits on energy-efficiency of operations

← Era of Delay-Constrained Computing

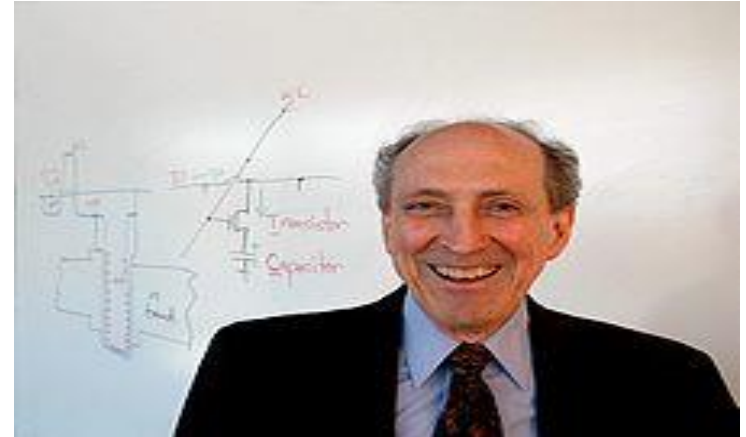
→ Era of Power-Constrained Computing

c. 2000

The magic underlying technology gains: Dennard Scaling



Dennard et. al., 1974



Robert H. Dennard, picture from Wikipedia

Dennard scaling in a nutshell:

Power = Capacitance \times Voltage² \times frequency

Scale transistor by $\alpha \rightarrow$ density grows by $\alpha^2 \rightarrow$ power nominally grows by α^2

To compensate: lower voltage by $\alpha \rightarrow$ more transistors; constant power!

So what happened? Leakage killed Dennard Scaling

To distinguish zeros and ones, supply voltage must be about triple the transistor switching (threshold): $V_{dd}/V_{th} > 3$

So, scaling down supply requires scaling down threshold

But, transistor leakage power is exponential in V_{th}

→ V_{dd} can't go down anymore

No more free lunch...

(2011 edition)

Dark Silicon: can't use all transistors all the time

Need system-level approaches to...

...turn increasing transistor counts into customer value

...without exceeding thermal limits

Energy efficiency is the new performance

21st Century Arch. – Key Challenges

Late 20 th Century	The New Reality
Moore's Law — 2× transistors/chip	Transistor count still 2× BUT...
Dennard Scaling —~constant power/chip	Gone. Can't repeatedly double power/chip
Modest (hidden) transistor unreliability	Increasing transistor unreliability can't be hidden
Focus on computation over communication	Communication (energy) more expensive than computation
1-time costs amortized via mass market	One-time cost much worse & want specialized platforms

How should architects step up as technology falters?

Recommendations from 2011

20 th Century	21 st Century	
Single-chip in stand-alone computer	Architecture as Infrastructure: Spanning sensors to clouds Performance + security, privacy, availability, programmability, ...	Cross-Cutting: Break current layers with new interfaces
Performance via invisible instr.-level parallelism	Energy First <ul style="list-style-type: none">• Parallelism• Specialization• Cross-layer design	
Predictable technologies: CMOS, DRAM, & disks	New technologies (non-volatile memory, near-threshold, 3D, photonics, ...) Rethink: memory & storage, reliability, communication	

Six years elapse...
... and some new realities emerge

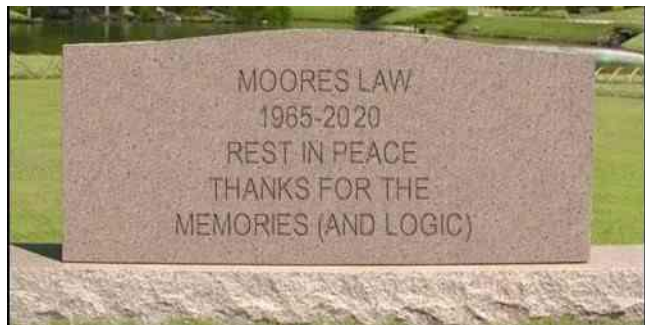
What changed?

Machine learning is a key workload

Specialization already happening at scale

Cloud is truly ubiquitous

Wide acceptance Moore's Law is really ending



Holographic Processing Unit



Arch2030 Visioning Workshop: Process

Reached out to prior efforts

21st Century CA, Rebooting Computing

Reached out to community for input

Invited experts (devices, applications)

Held with ISCA: 120+ participants

Sent out report for comments

40+ endorsers

What are the most important challenges to be addressed by architecture research in the next 15 years? *

Long answer text

What are some of these challenges/trends that you are *not* working on? *

Long answer text

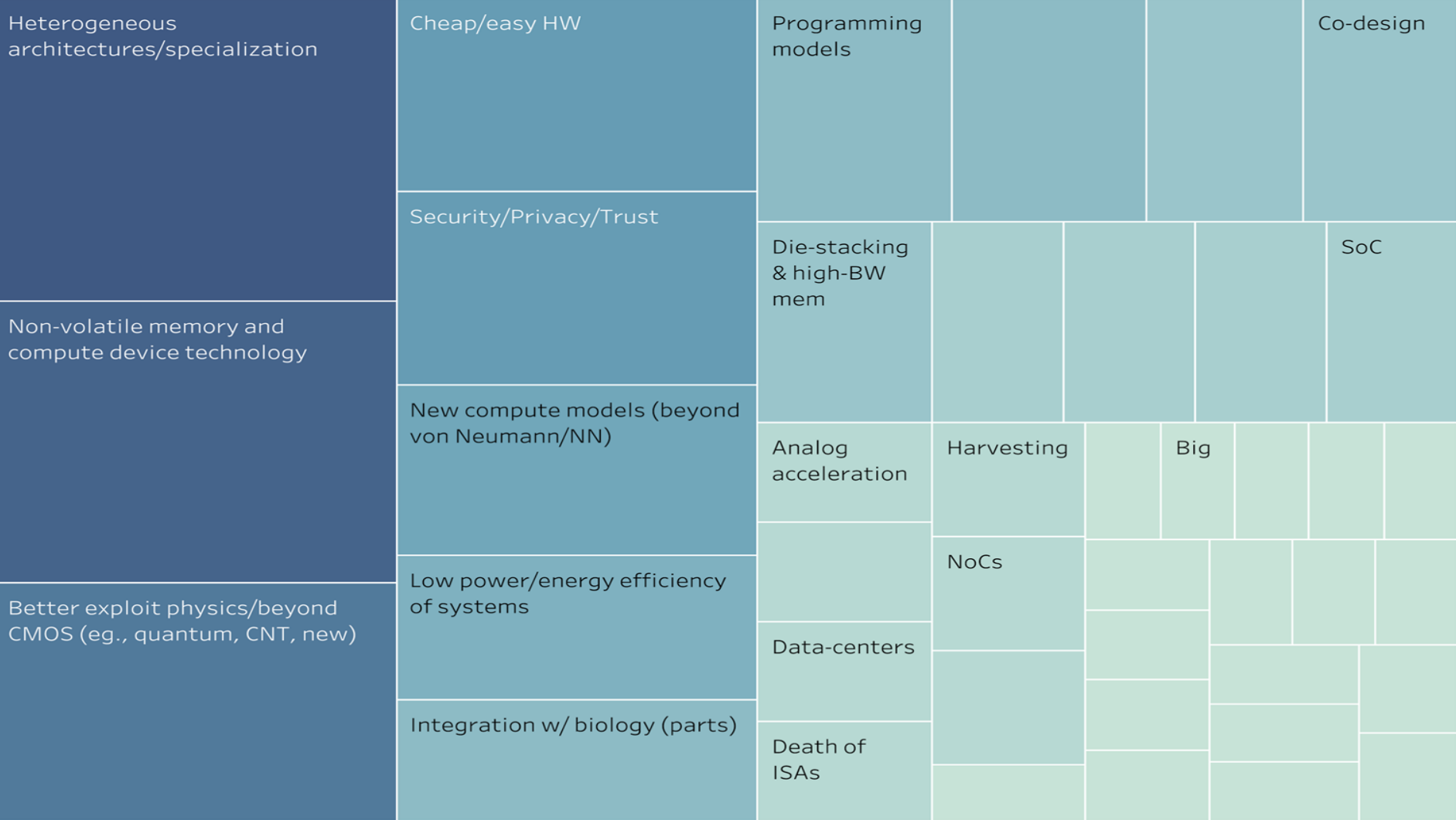
What are the main technology opportunities that the community should take advantage of? *

Long answer text

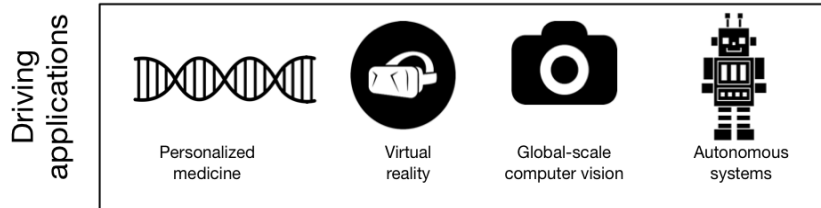
How well do you think the architecture community anticipated today's challenges/opportunities? *

1 2 3 4 5

Bad ☐ ☐ ☐ ☐ ☐ Great!



Arch2030: The upshot



<i>Observation</i>	<i>Implications for next 15 year</i>
1. Specialization gap	Democratize HW design: tools and open source designs
2. Ubiquitous cloud: innovation abstraction	Cloud model provides practical deployment path for new architectures
3. 3D stacking is real	Opportunities for new architectures and integration models
4. Getting “closer to physics”	Need for more adventurous architectures
5. Machine learning as key app. component	New architectures are enablers: need real collaboration with core ML community

1. Specialization Gap: Democratizing HW Design

Performance gap: many applications aren't possible without specialization

AR/VR, autonomous vehicles, large-scale AI/ML

General purpose processors aren't efficient enough

Design cost/effort gap:

HW design costs growing too fast

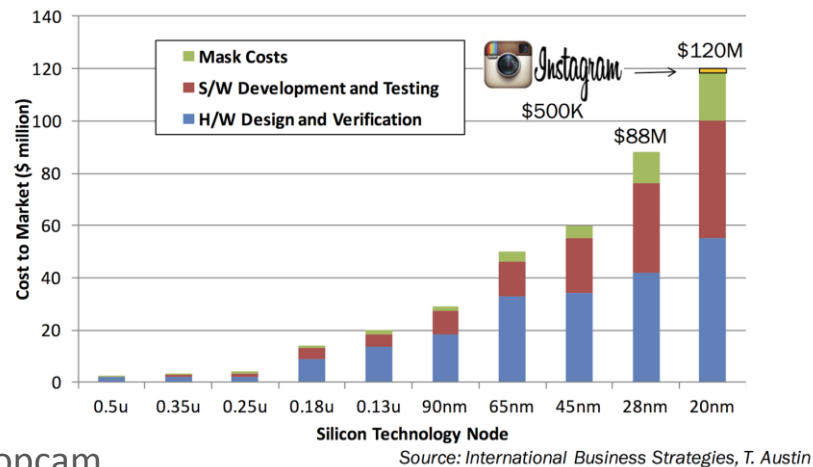
Need: better models, tools, open-source design

Can create new business/innovation forces

Emerging “HW” companies: fitbit, Oculus, Pebble, Dropcam, ...

Open source can create agility for ASIC-based startups

**Developing specialized hardware must become as easy,
inexpensive, and agile as developing software**



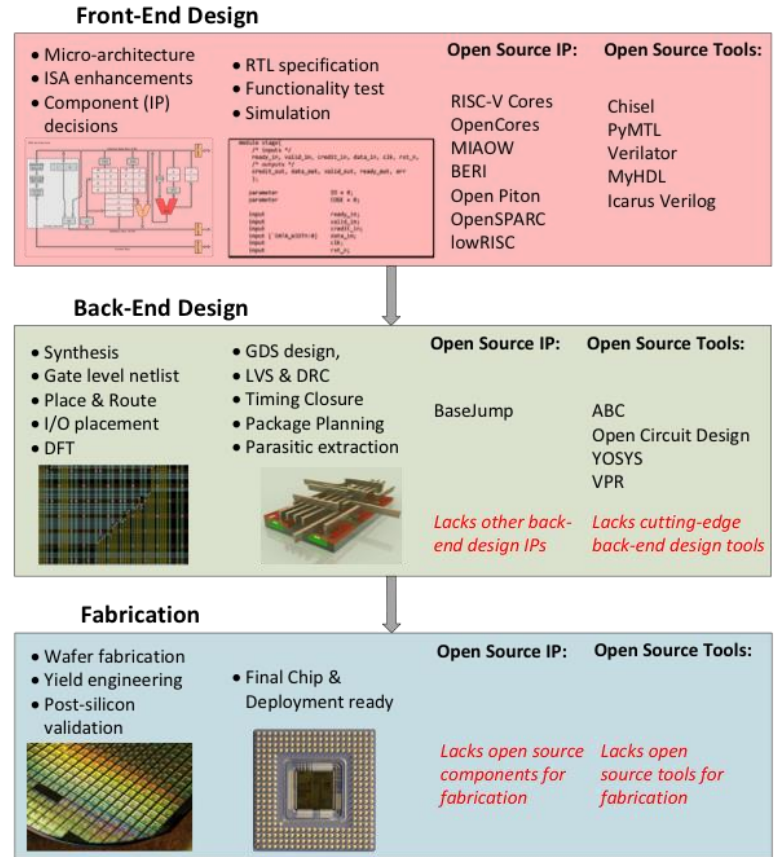
Opportunity:

Need infrastructure to reduce barrier-to-entry for custom ASICs

Faster impact via tightly integrated FPGAs

Need open/reusable IP cores and tools

Investigate “chiplet” / post-fab integration



Sankaralingam et al.

2. Cloud as Abstraction for Architectural Innovation

Ubiquitous public cloud infrastructure (Microsoft, Google, Amazon)

More than just software - entry point for new hardware

Clean service/microservice interfaces

Can hide exotic HW/devices

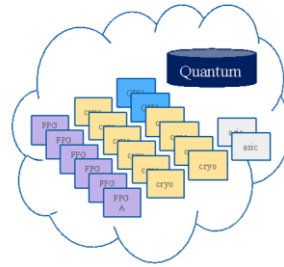
ASICs, FPGAs, quantum computers?



Yesterday



Today



Tomorrow

[Doug Carmean, ISCA'16 Keynote]

Through scale and virtualization, clouds can offer deep HW innovations transparently and at low cost

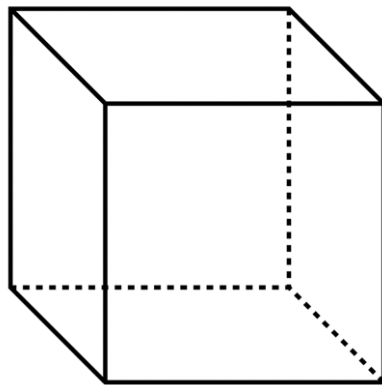
3. Going Vertical with 3D Integration

Denser memories, higher bandwidth

Capacity/bandwidth grows

Fundamental need for processing+memory integration

Integration of “chiplets” in 3D substrate a promising design/business model



capacity $\propto L^3$

bandwidth $\propto L^2$

3D integration provides a new dimension of scalability

4. Getting Closer to Physics

New memories and devices

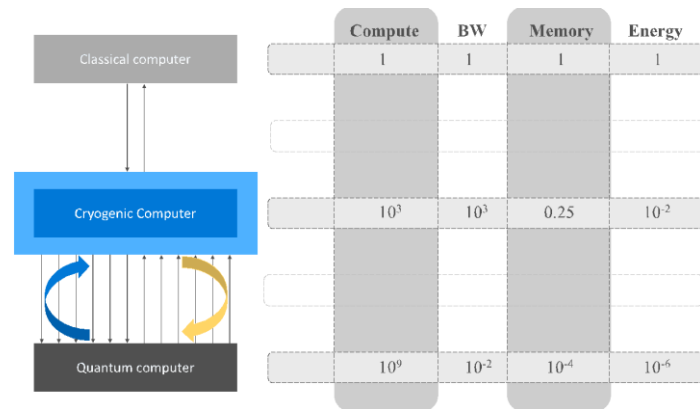
Carbon nanotubes

Quantum computing and superconducting logic

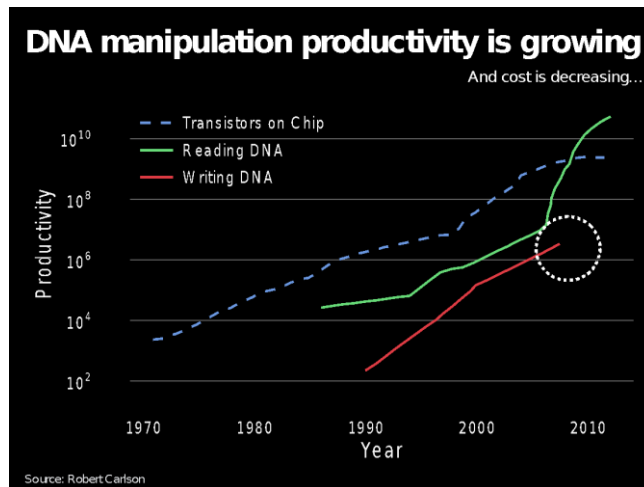
Borrowing from biology

	Access Time	Capacity	Durability
Flash	μs -ms	TBs	~5 yrs
HDD	10s ms	100s TBs	~5 yrs
Tape	minutes	PBs	~10s yrs
DNA-based Archival	hours	ZBs	~1000s yrs

[Bornholt et al.]



[Doug Carmean, ISCA'16 Keynote]



5. Machine Learning as a Key Workload

Training: HPC-like systems, turn-around time matters to evolution

Inference: Low latency, low power

Strong driver for architecture and systems innovation

Tensor flow, TPUs, MS CNTK, ...



Google's Tensor Processing Unit

Hardware advancement enables machine learning over “bigger data”

The Future: Architecture + *X*

Application and technology driven

It's clear we are beyond a processor + memory centric world

Examples: sensor/compute fusion, intelligent networks, intelligent storage systems

Critical to reach out to other CS areas and fields

What does it mean for GreenMetrics?

What does it mean for GreenMetrics? (1 of 2)

All of computing needs to think about “Democratizing HW Design”

What does source hardware mean for sustainability?

Can we foster a “Github movement” for HW?

What does the tech transfer pipeline from idea to system look like?

“The Cloud” does not mean commodity servers anymore

Tensor Processing Unit

Project Catapult at Microsoft

FPGA instances at Amazon

What does it mean for GreenMetrics? (1 of 2)

Machine Learning is everywhere

How do we enable sustainable training & inference?

Unsustainable to do all this compute in centralized data centers;
it needs to be pushed to the client/edge

Exotic hardware is coming

3D transistors? Memristors? Carbon Nanotubes?

Quantum? Biology inspired? DNA storage?