

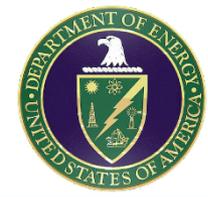
An Architect's Point of View of the Post Moore Era

Dr. George Michelogiannakis

Research scientist
Computer architecture group
Lawrence Berkeley National Laboratory

Work with Dr. Dilip Vasudevan

These are not DOE's or LBNL's official views



Poll: What Did Dr. Moore Say



- ◆ Transistor density will increase every **12** months
- ◆ Transistor density will increase every **18** months
- ◆ Transistor density will increase every **24** months

(may have multiple answers)



Poll: What Did Dr Moore Say



- ◆ Transistor density will increase every 12 months
 - In 1965
- ◆ Transistor density will increase every 18 months
- ◆ Transistor density will increase every 24 months
 - In 1975

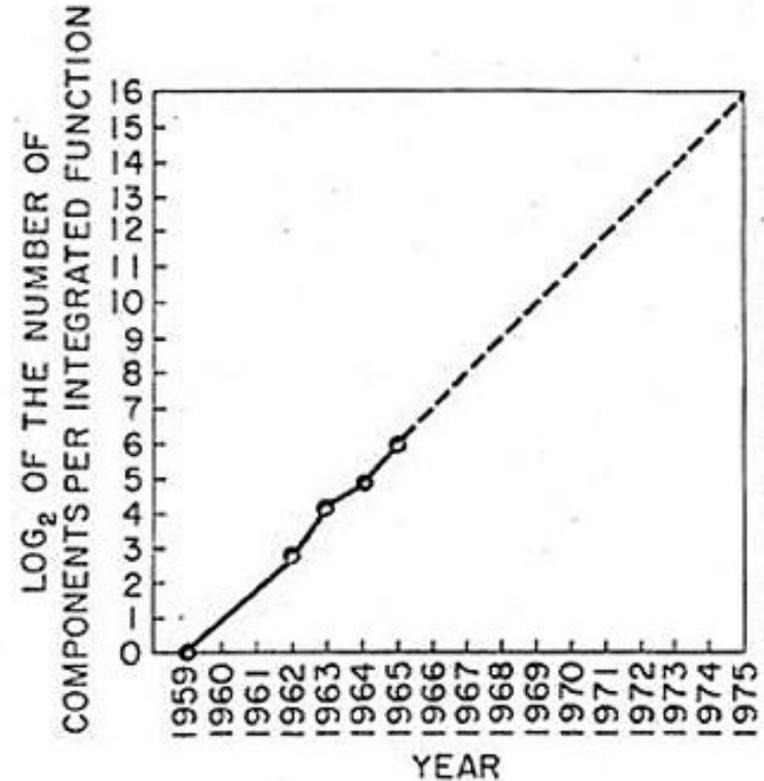
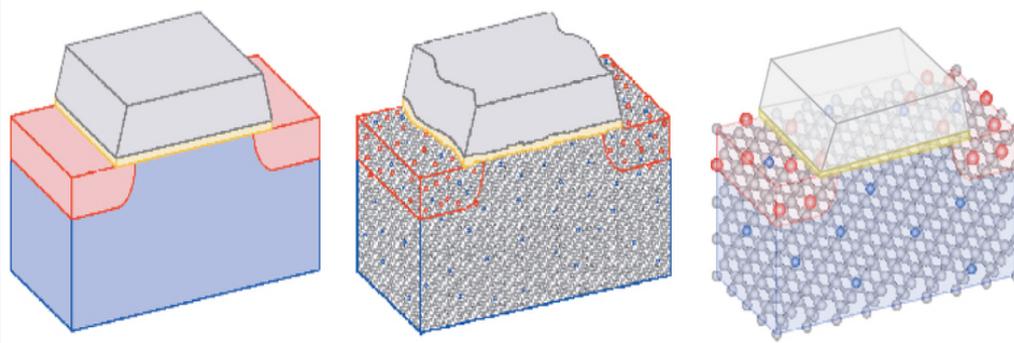


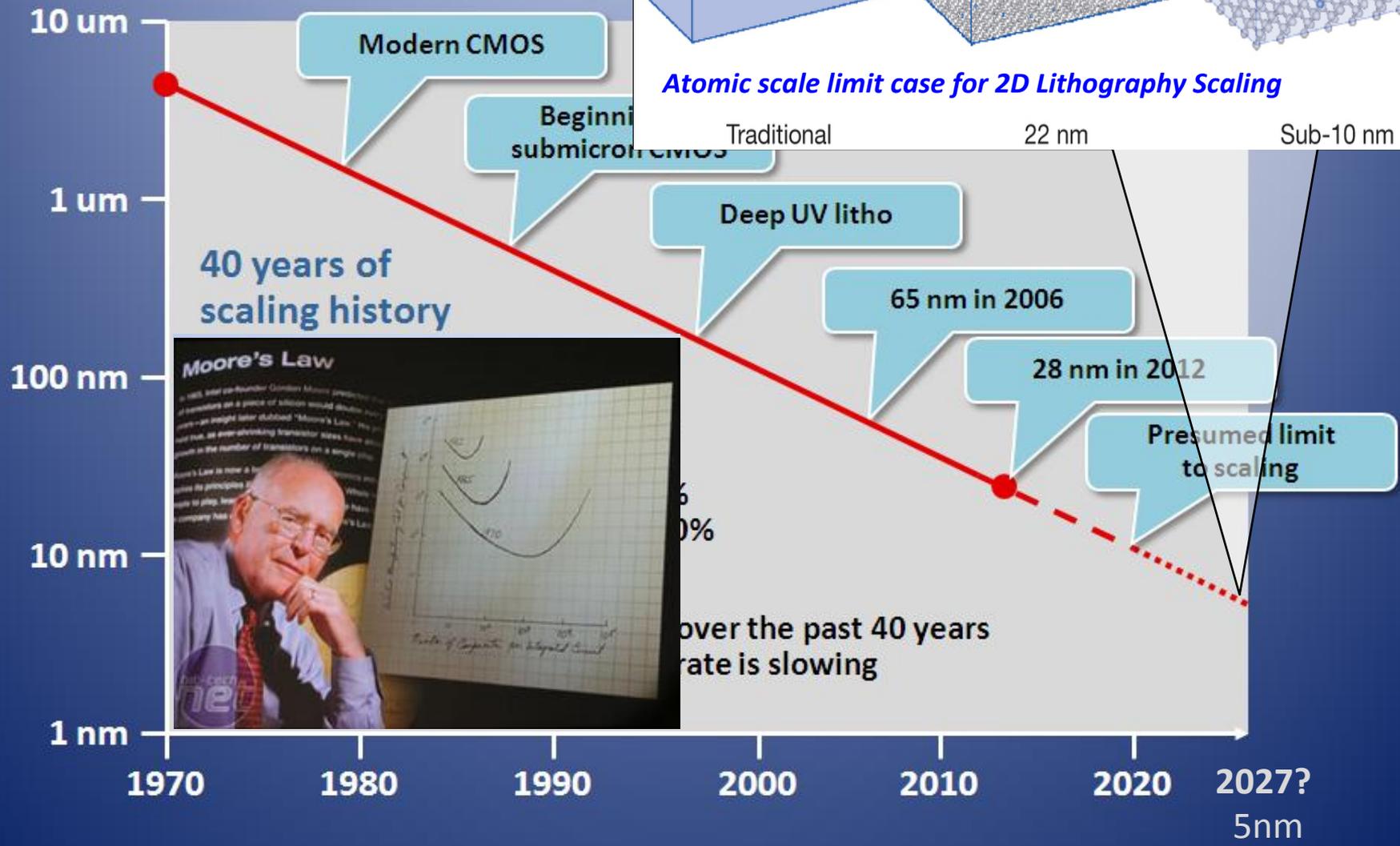
Fig. 2 Number of components per integrated function for minimum cost per component extrapolated vs time.

Dr. Moore's 1965 paper

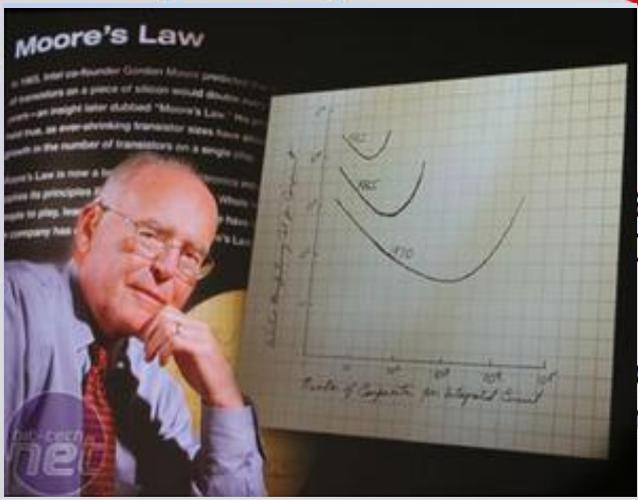
50 years of Sem



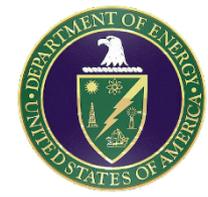
Atomic scale limit case for 2D Lithography Scaling



40 years of scaling history



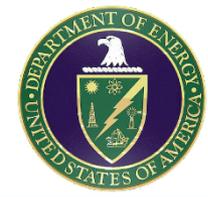
over the past 40 years rate is slowing



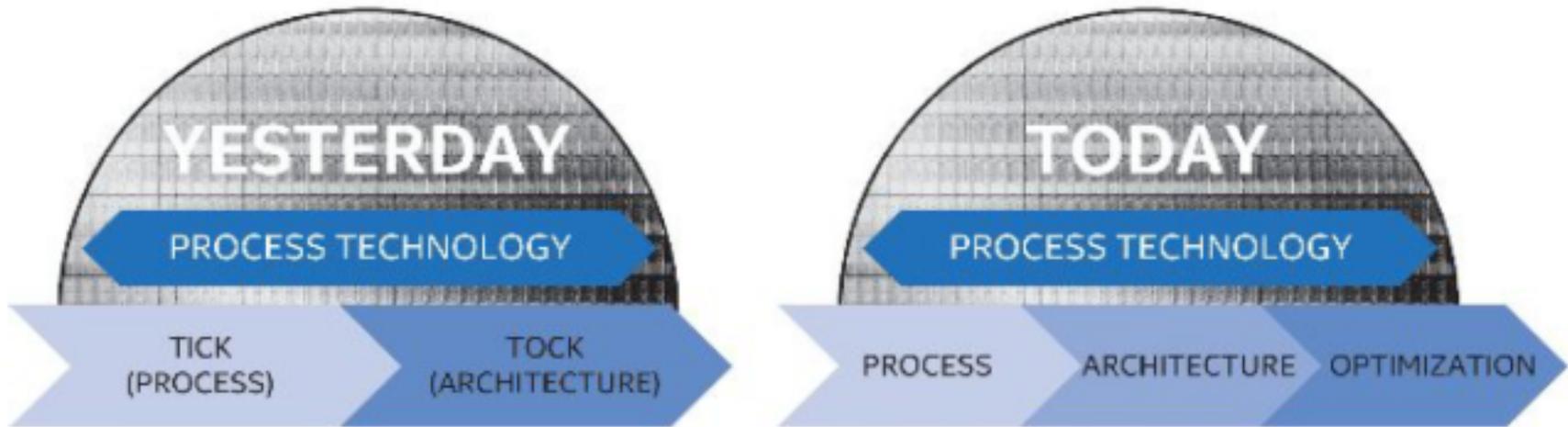
Moore's Law of Documentation



new "Moore's Law" on documentation volume
seen from the 14th floor at Fermilab perspective



Scaling Already Slowing Down



Peter Bright "Intel retires "tick-tock" development model, extending the life of each process", 2016



Preserve Performance Scaling With Emerging Technologies



Now – 2025

Moore's Law continues through ~5nm -- beyond which diminishing returns are expected.

2016

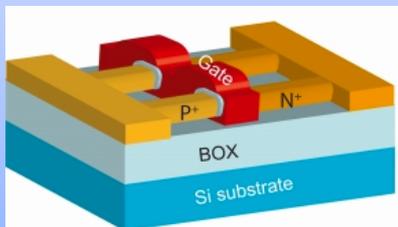
2016-2025

End of Moore's Law
2025-2030?

Post Moore Scaling

New materials and devices introduced to enable continued scaling of electronics performance and efficiency.

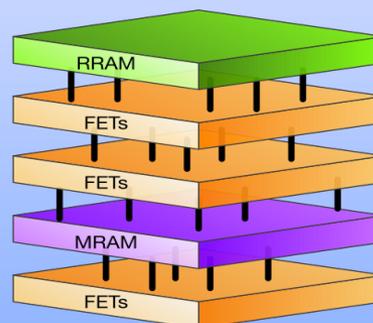
2025+



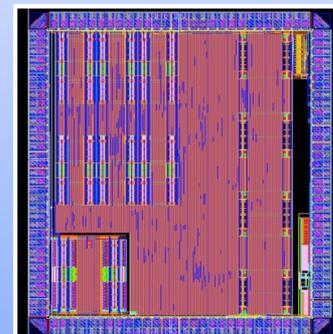
Emerging transistors



Emerging memories



3D integration



Specialized architectures

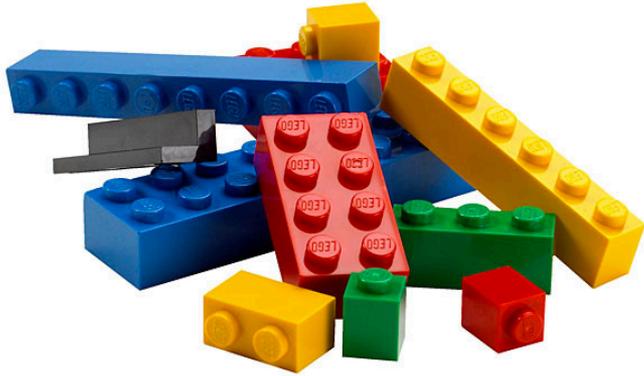
+ others



An Architect's Point of View

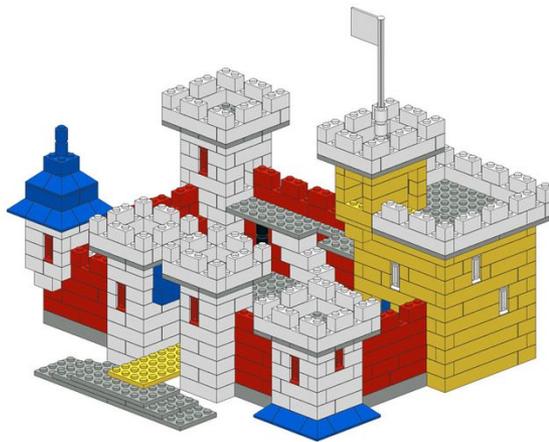


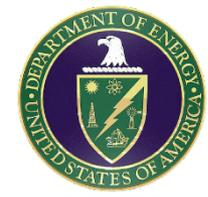
An Architect's Job





Lego Designs Have Been Getting Larger





New Lego Pieces



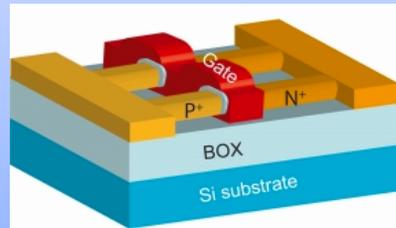
- ◆ Old designs can no longer become smaller with same strength
- ◆ Lego came up with new pieces:



- ◆ Which ones do we use?
 - What is their building-wide impact?
- ◆ How does each one change the optimal design?
- ◆ How does each piece interact with others?
- ◆ What feedback can we provide Lego to refine each piece?



Emerging Transistors



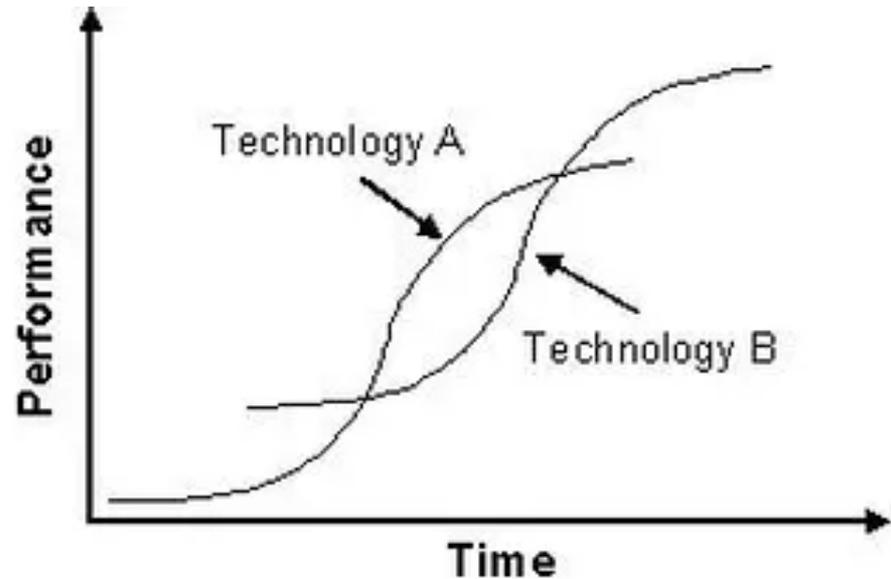
Emerging
transistors



New Devices



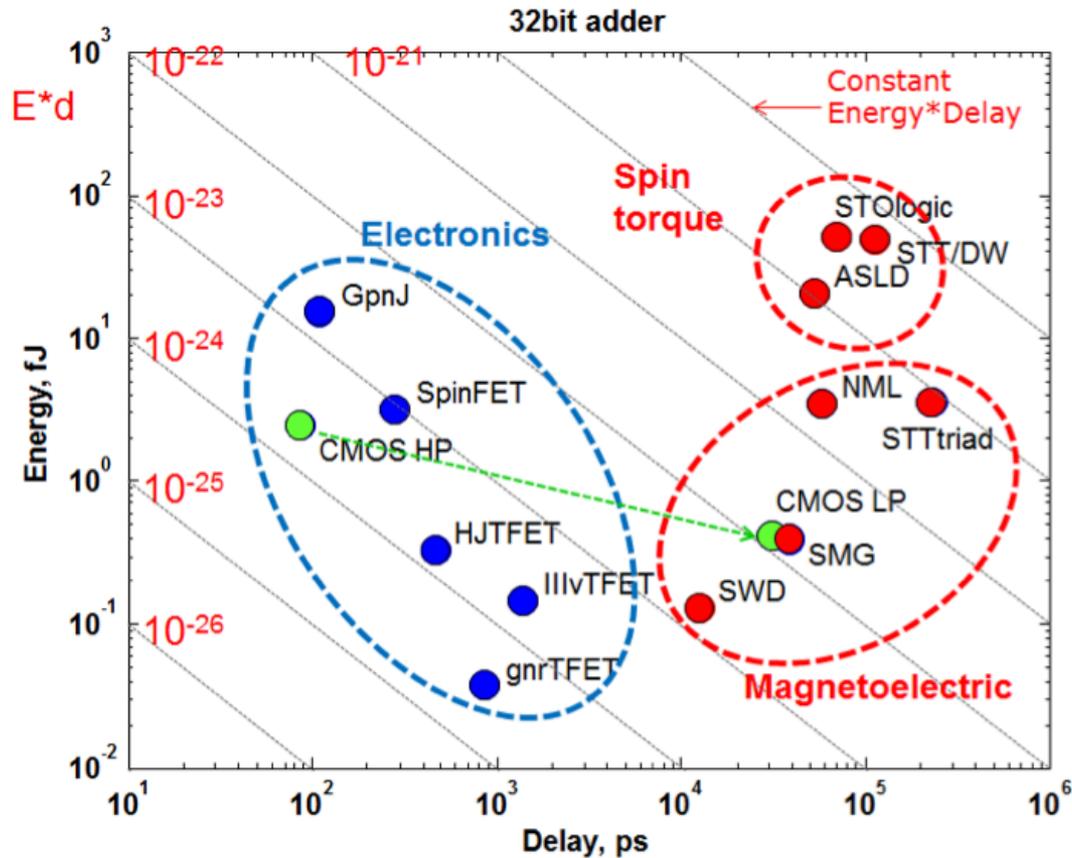
- ◆ New devices need time to show their potential
- ◆ Two broad categories:
 - New designs
 - New materials
- ◆ Maybe not a single replacement for MOSFETs



Rick Lindquist "3 Steps for Constructive Disruption"



Many More



Nikonov and Young, "Benchmarking of Beyond-CMOS Exploratory Devices for Logic Integrated Circuits", 2015

Each dot is a moving target. We have to judge the potential



Emerging Memories



Emerging
memories



Many Memories As Well



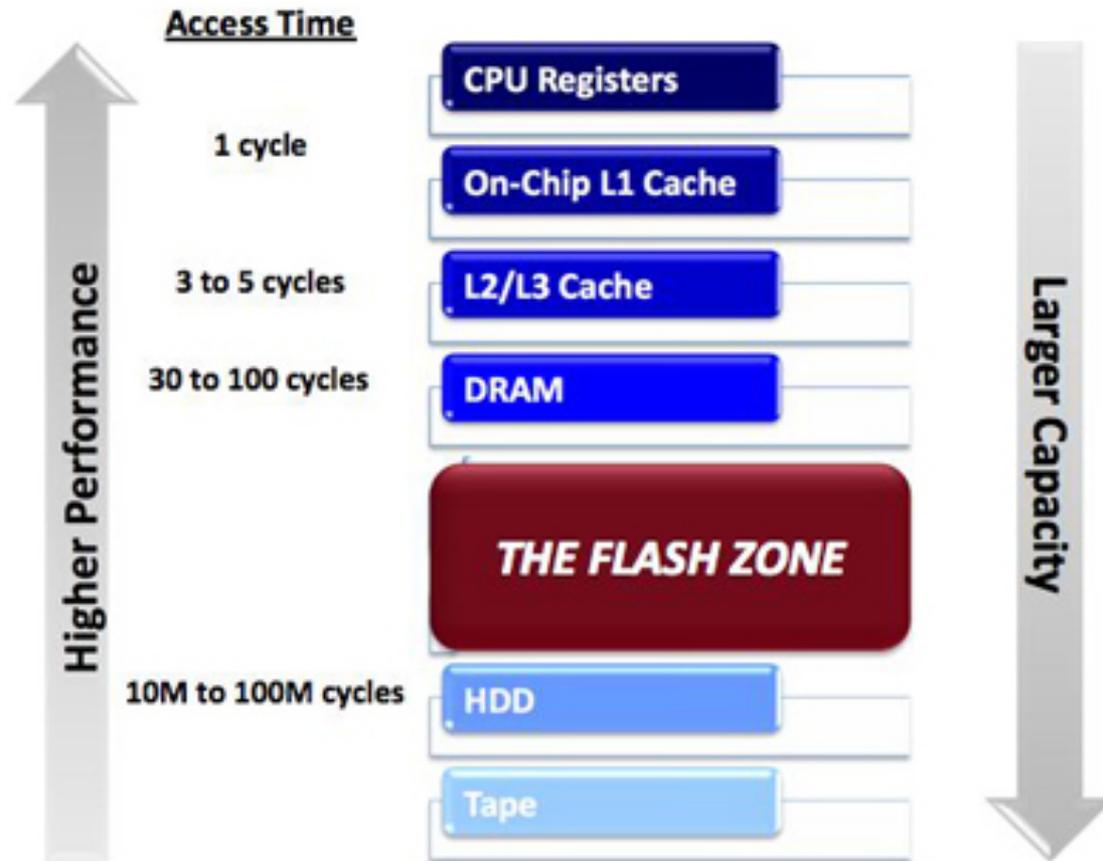
- ◆ Some of these are non-volatile

	SRAM	DRAM	eDRAM	2D NAND Flash	3D NAND Flash	PCRAM	STTRAM	2D ReRAM	3D ReRAM
Data Retention	N	N	N	Y	Y	Y	Y	Y	Y
Cell Size (F ²)	50-200	4-6	19-26	2-5	<1	4-10	8-40	4	<1
Minimum F demonstrated (nm)	14	25	22	16	64	20	28	27	24
Read Time (ns)	< 1	30	5	10 ⁴	10 ⁴	10-50	3-10	10-50	10-50
Write Time (ns)	< 1	50	5	10 ⁵	10 ⁵	100-300	3-10	10-50	10-50
Number of Rewrites	10 ¹⁶	10 ¹⁶	10 ¹⁶	10 ⁴ -10 ⁵	10 ⁴ -10 ⁵	10 ⁸ -10 ¹⁰	10 ¹⁵	10 ⁸ -10 ¹²	10 ⁸ -10 ¹²
Read Power	Low	Low	Low	High	High	Low	Medium	Medium	Medium
Write Power	Low	Low	Low	High	High	High	Medium	Medium	Medium
Power (other than R/W)	Leakage	Refresh	Refresh	None	None	None	None	Sneak	Sneak
Maturity									

J.S. Vetter and S. Mittal, "Opportunities for Nonvolatile Memory Systems in Extreme-Scale High Performance Computing," *CiSE*, 17(2):73-82, 2015.

- ◆ Non-volatility higher at the hierarchy
 - Challenge assumption that non-volatile storage is slow and distant
- ◆ New memories have different read, write, reliability constraints
- ◆ New memory hierarchy likely different

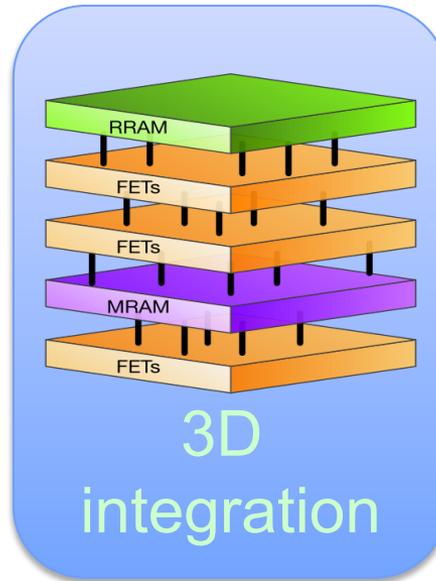
Flash Zone



AGIGARAM "The Flash Zone"



3D Integration



Enabled by Emerging Nanotechnologies

Massive Sensing

Data Storage
(NV memory)

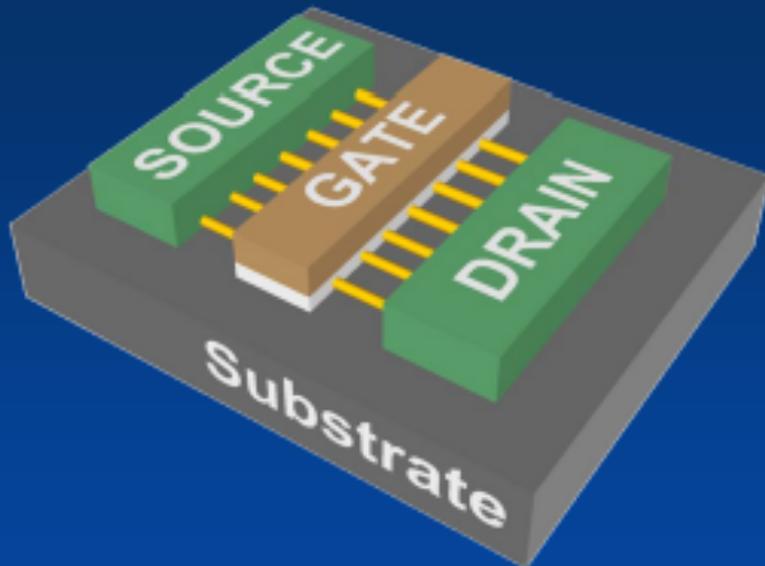
Computing Logic

Fine-grained
3D integration
(not TSVs)

- Low-temperature fabrication: **<400 °C**

Logic

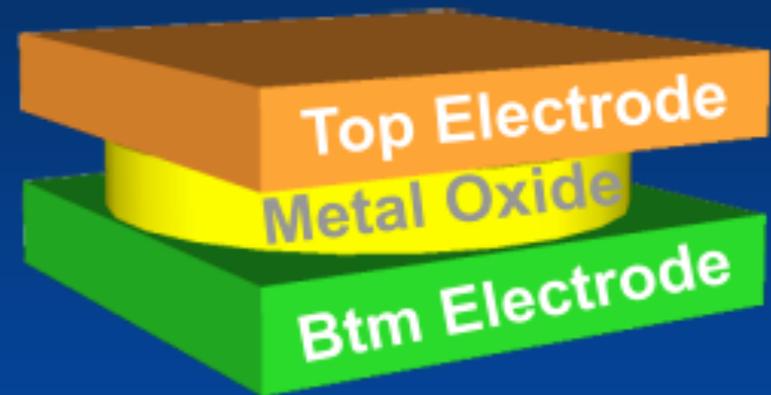
Carbon Nanotubes



<200 °C

Memory

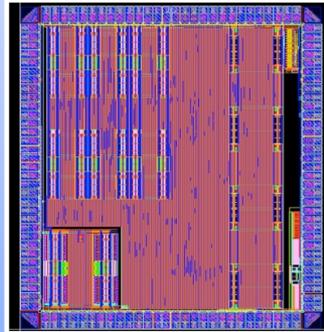
Resistive RAM



<200 °C

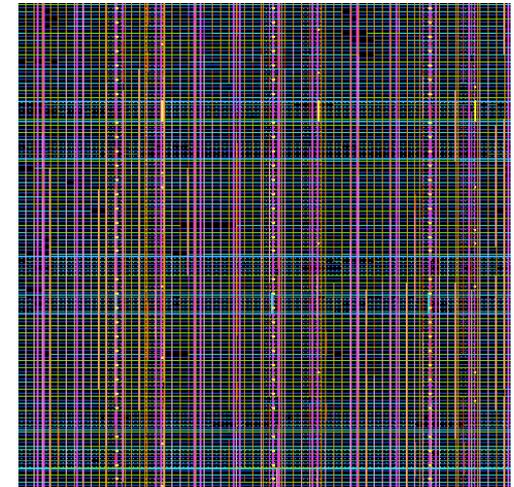
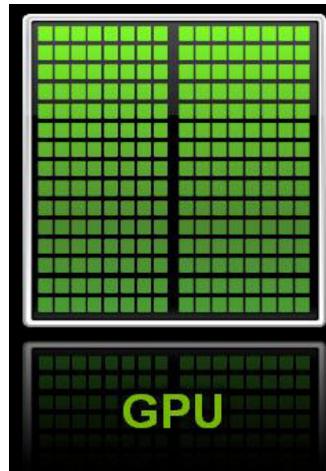
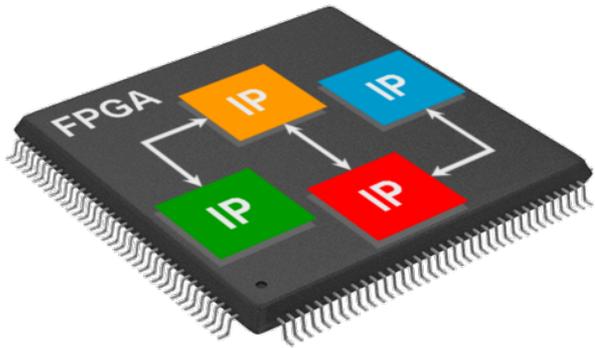


Specialization



Specialized architectures

- ◆ Hardware that is more suited for specific kinds of computation
 - Can also have accelerators for data transfer

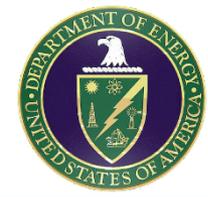


General
purpose

Accelerators

Fixed
function

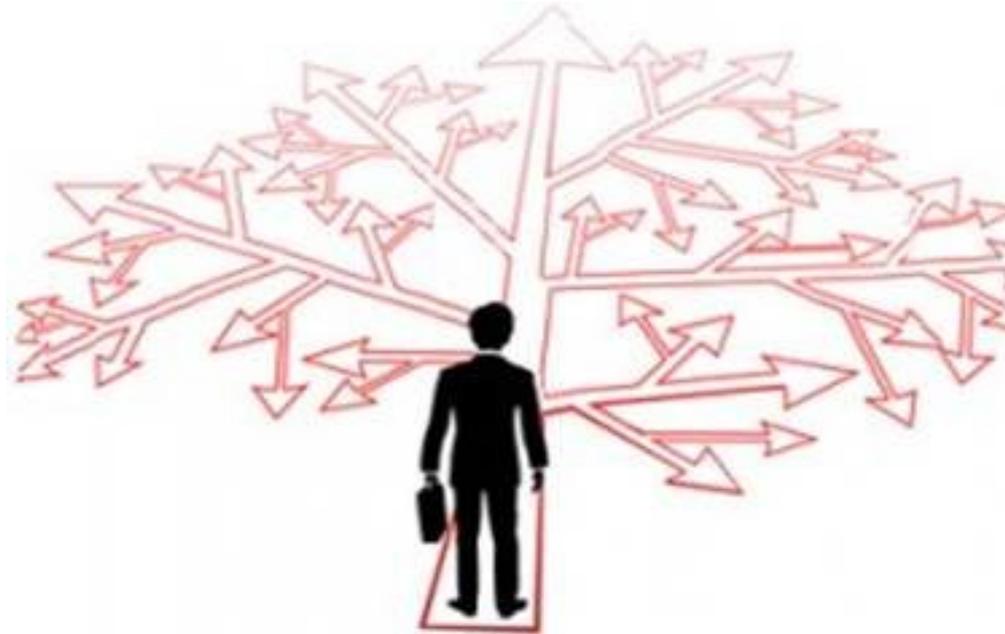




The Variety of Choices Is Overwhelming

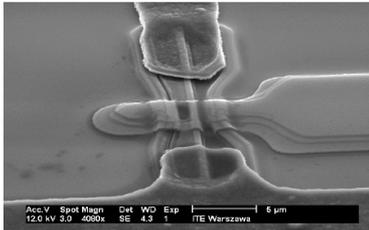


- ◆ The vast number of choices is a problem by itself
 - It makes finding a good design harder, especially when designing manually

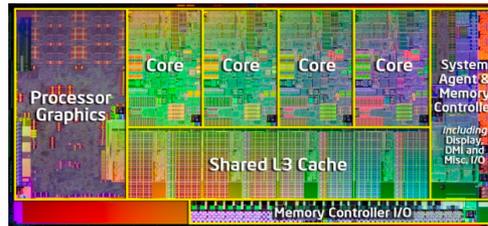


- ◆ Evaluating each option in isolation misses the big picture
 - Devices can be better designed with high-level metrics
 - Architects can figure out how to best use new technologies
 - Software experts can assess impact to programmability and compilers

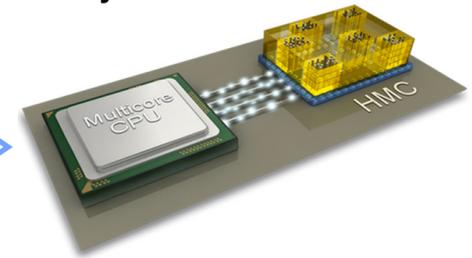
Transistor/Devices



Architecture



System

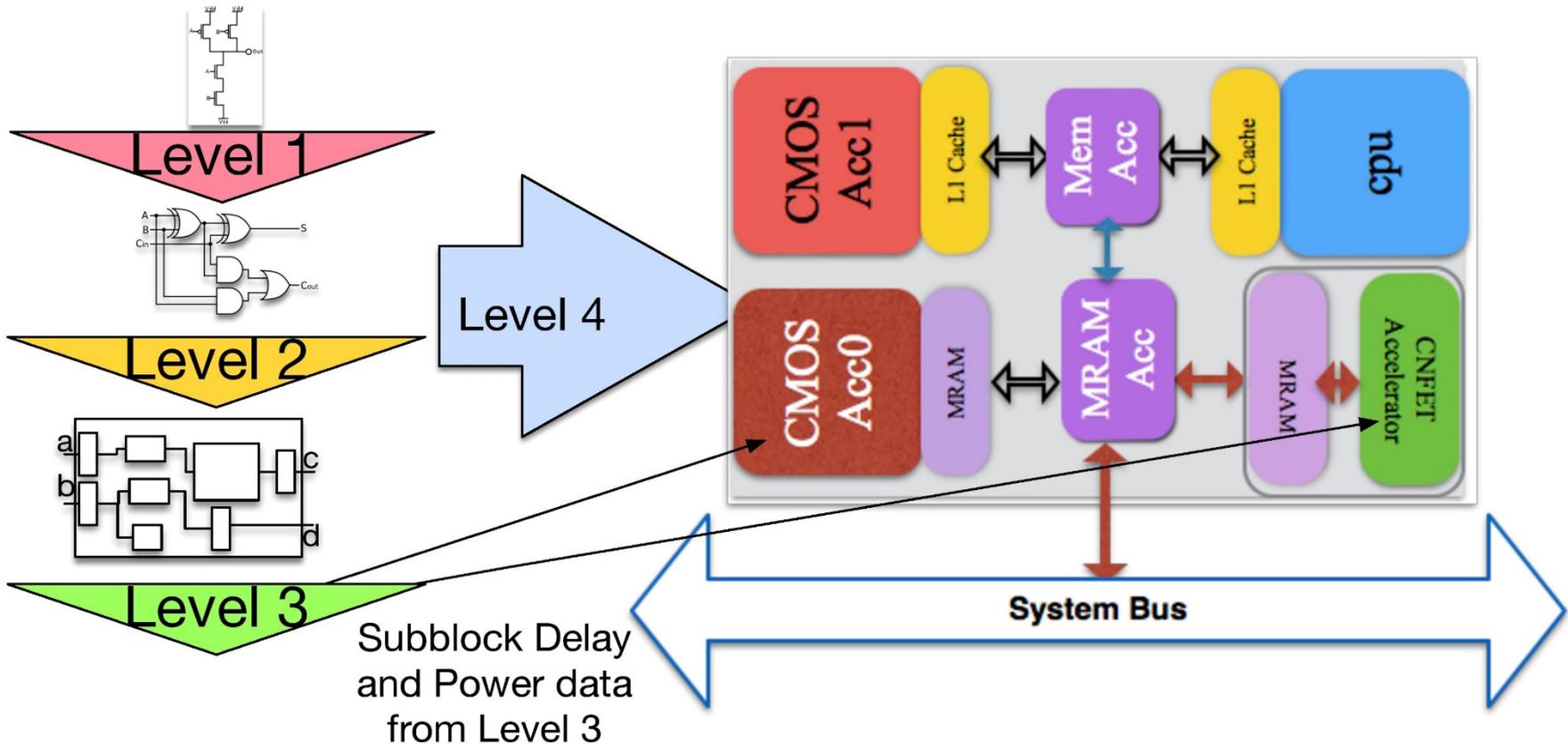


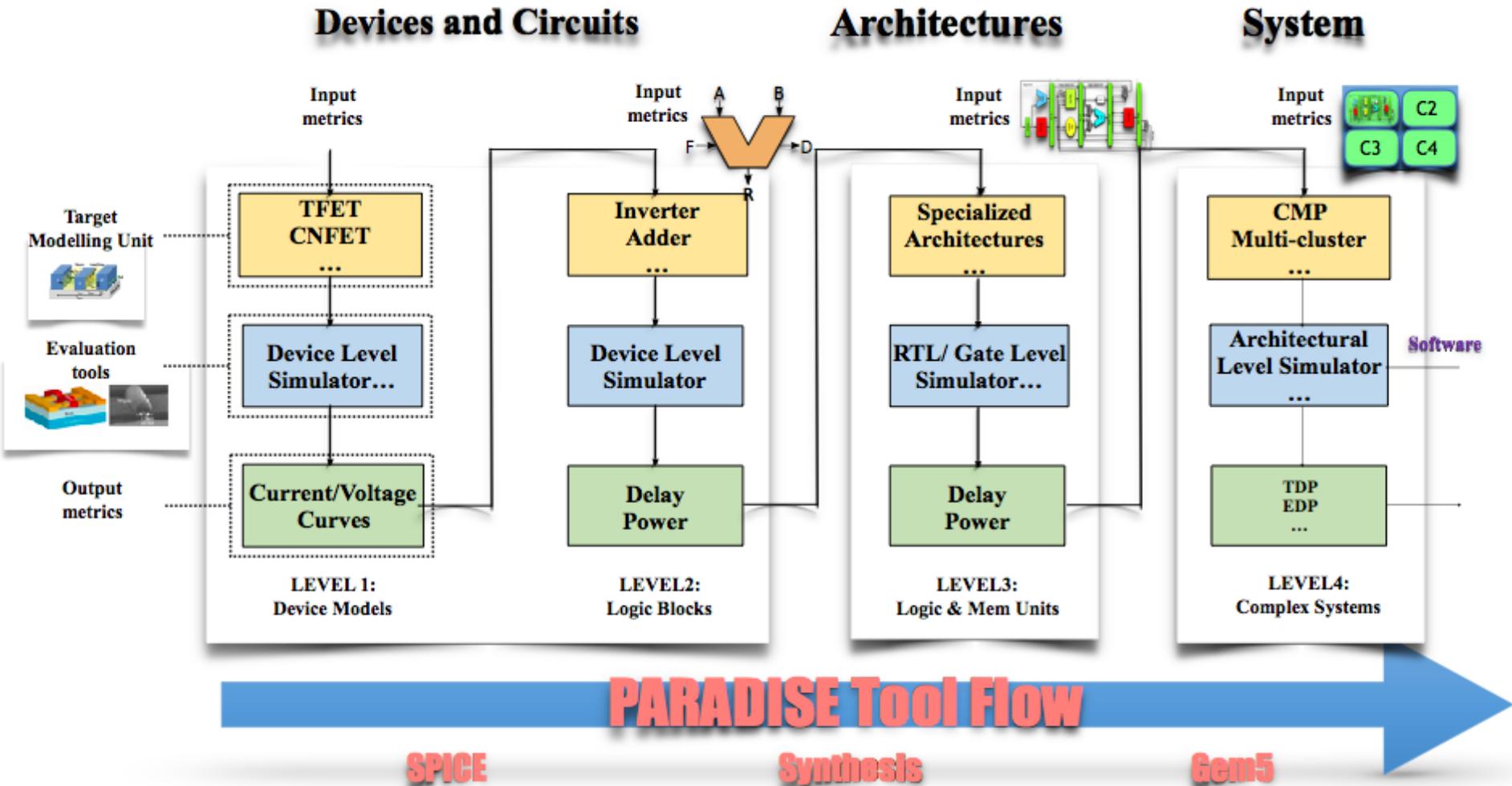
- ◆ But we lack the tools to do so systematically for many technologies



How To Make An Architect's Job Easier?

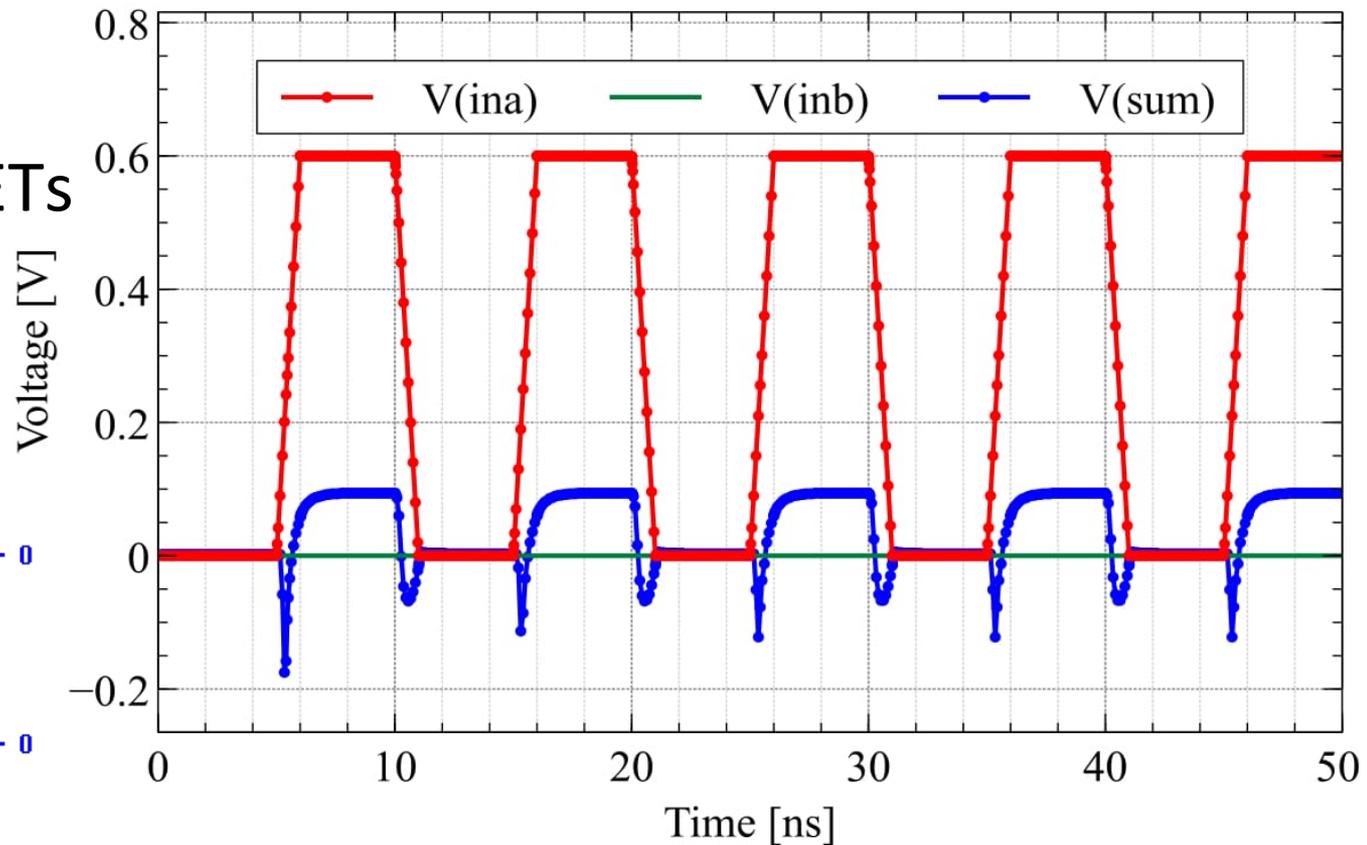
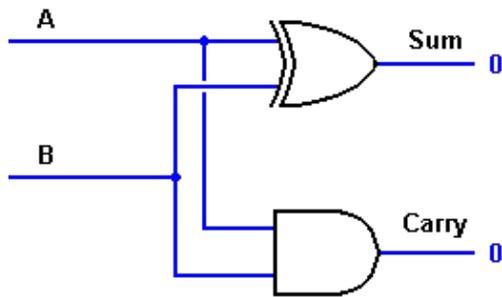
Beyond Moore Architectural Simulation from Emerging Device and Gate Level analysis

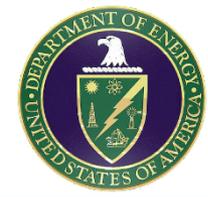




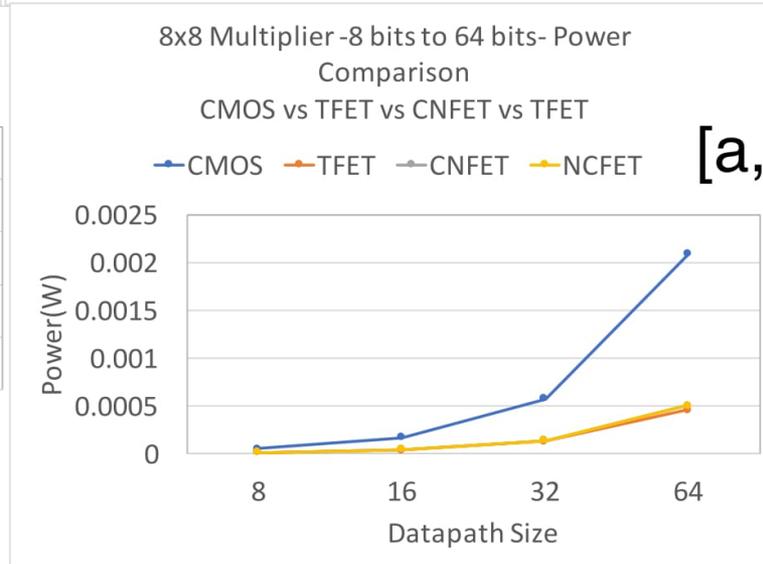
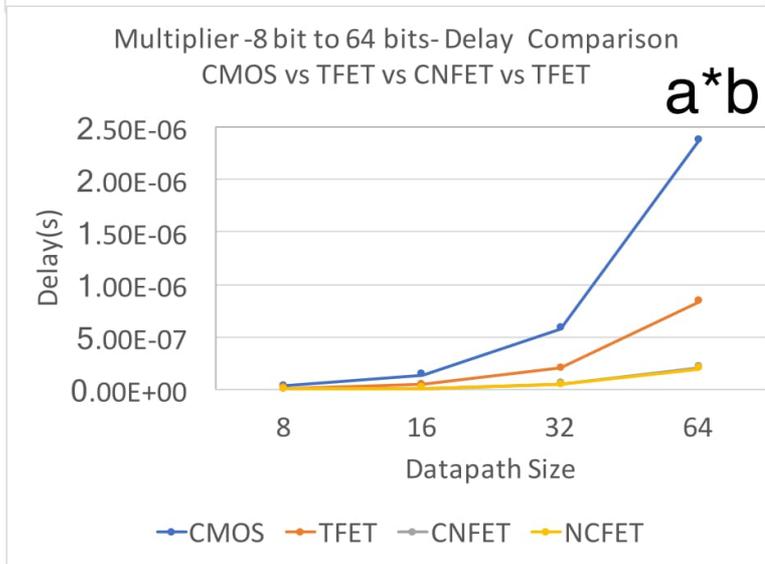
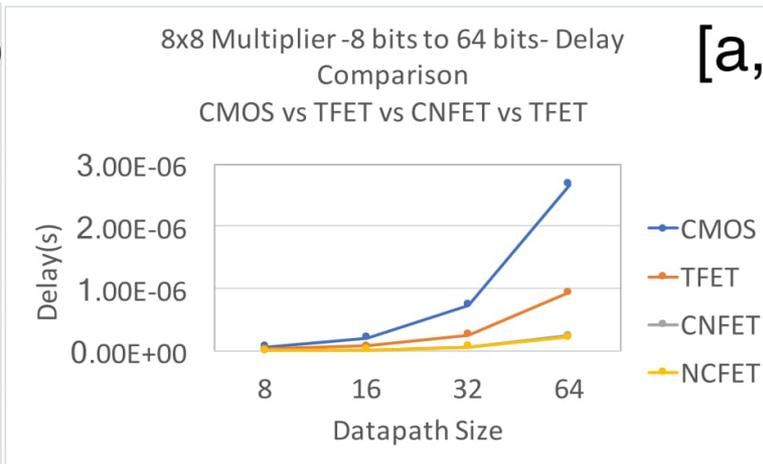
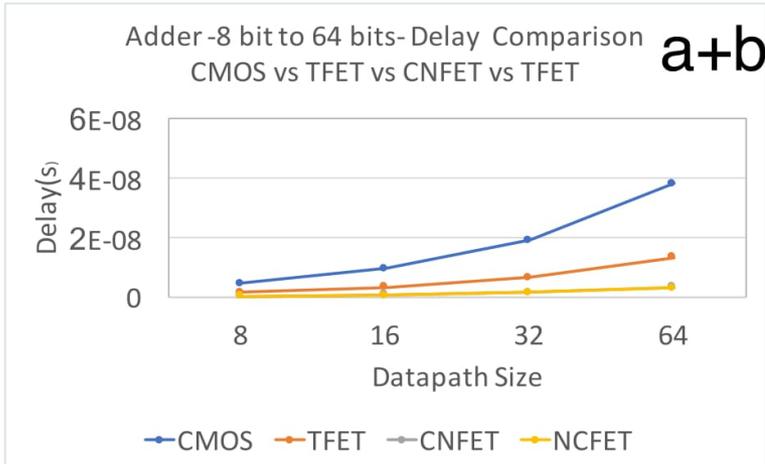
- ◆ Level 1 is the input for devices
- ◆ Xyce: open source parallel SPICE client

Adder using TFETs





Comparison Studies



Delay comparison for a multiplier

Power comparison for a multiplier

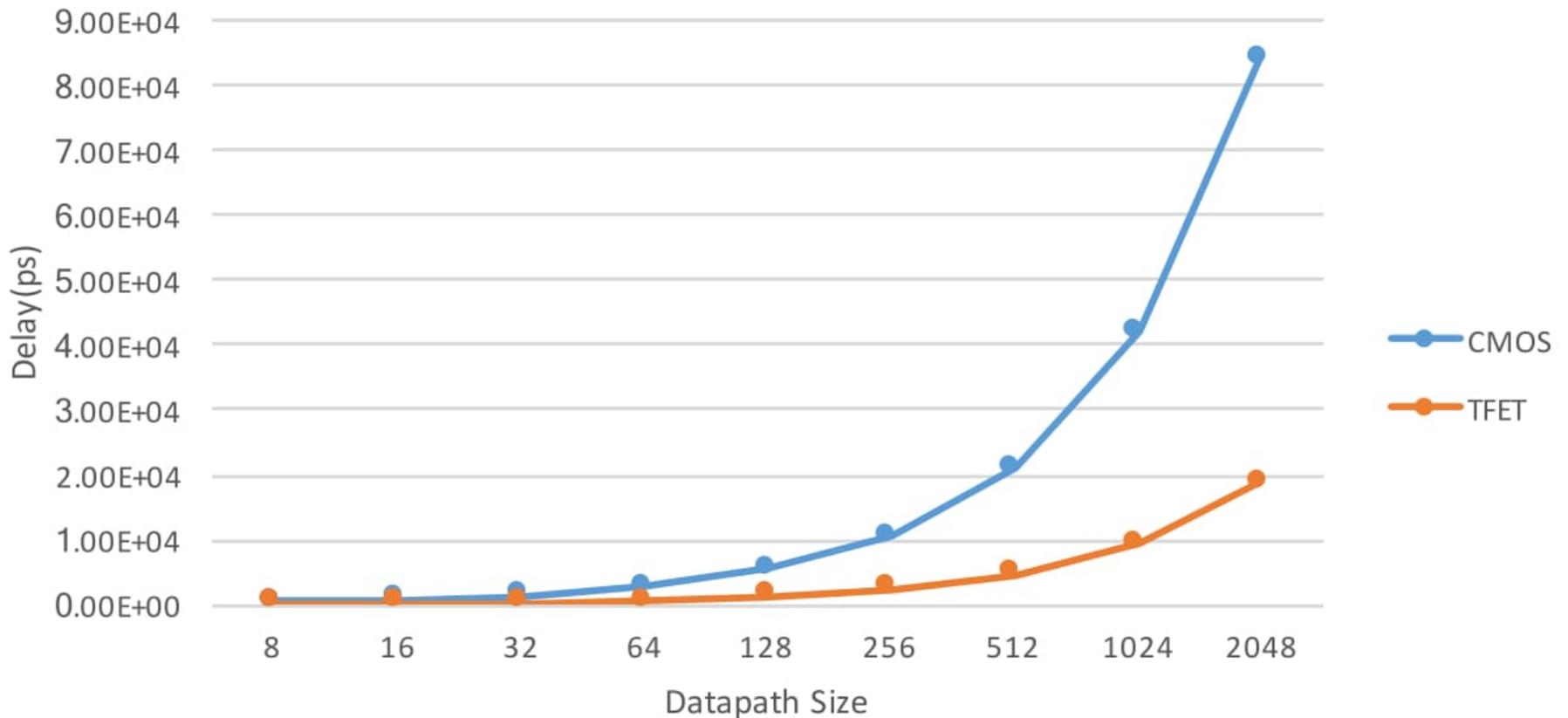


Level 3: RTL Synthesis



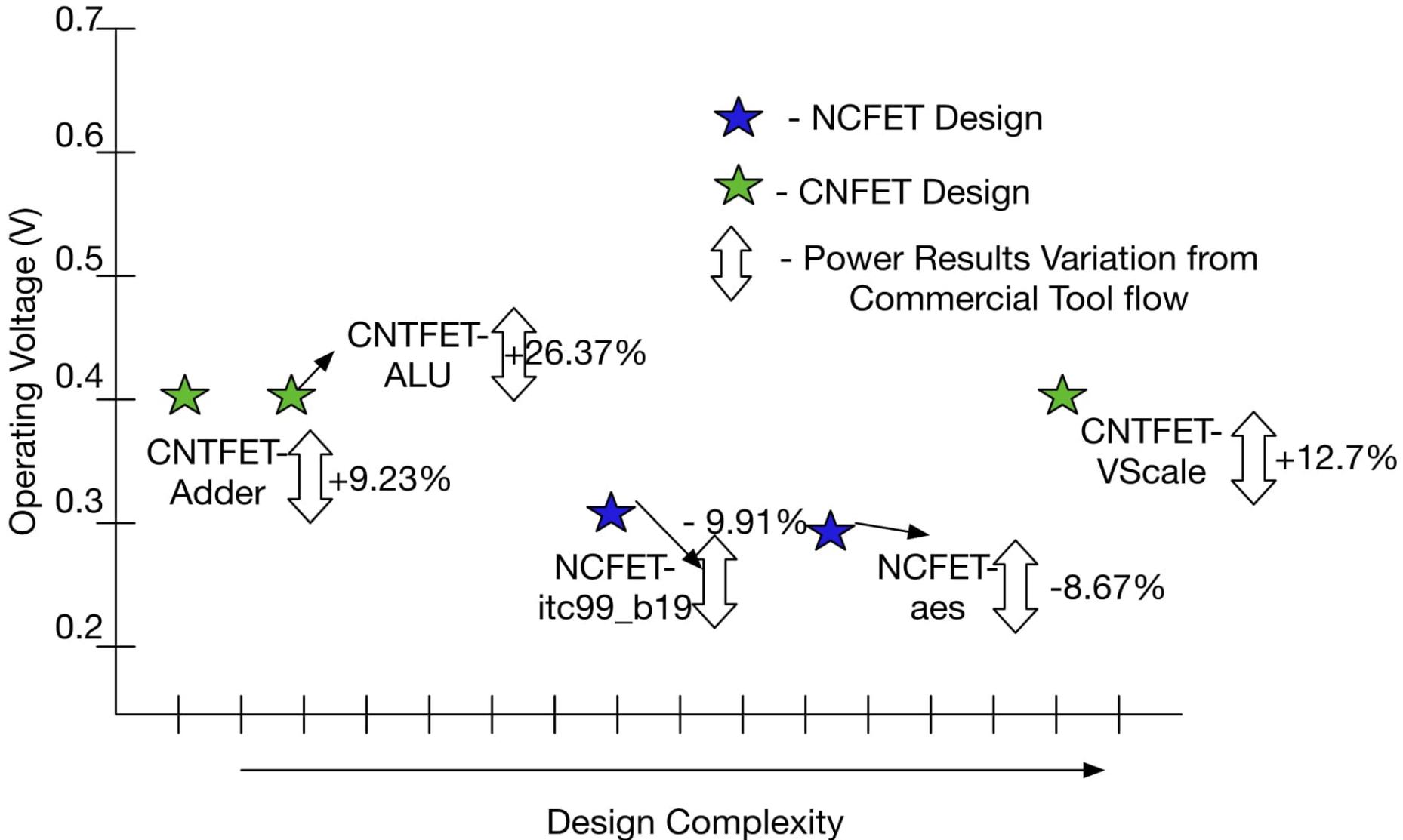
- ◆ Synthesis using Yosys and our own extension for power estimation

PARADISE Level 2, 3 - Logic Synthesis using Standard Cell Library for Adder 8bits - 2048 bits





Design Space Exploration at RTL Level

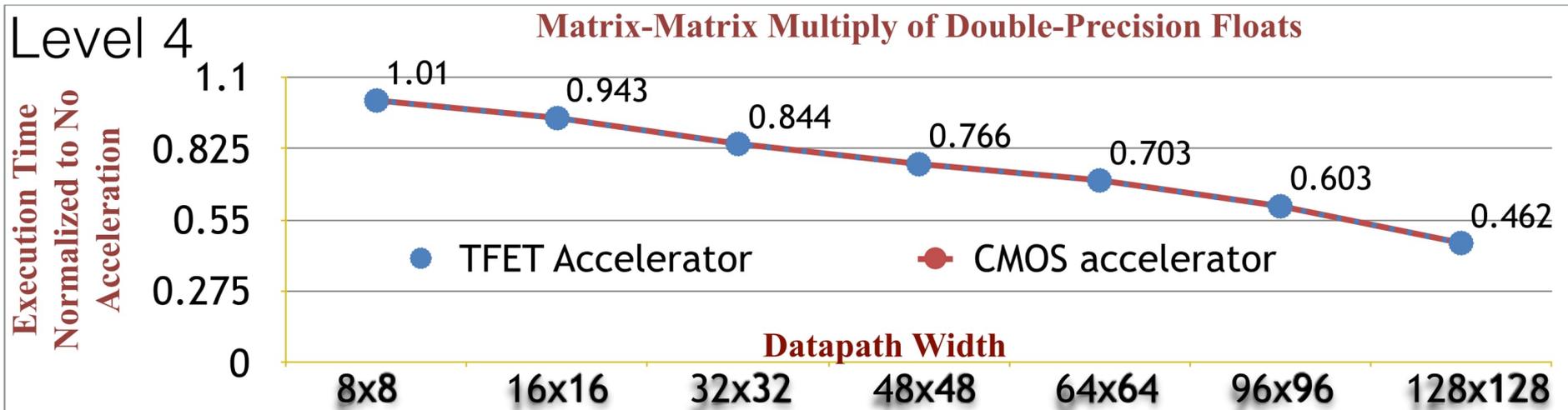




Level 4: Architectural Level



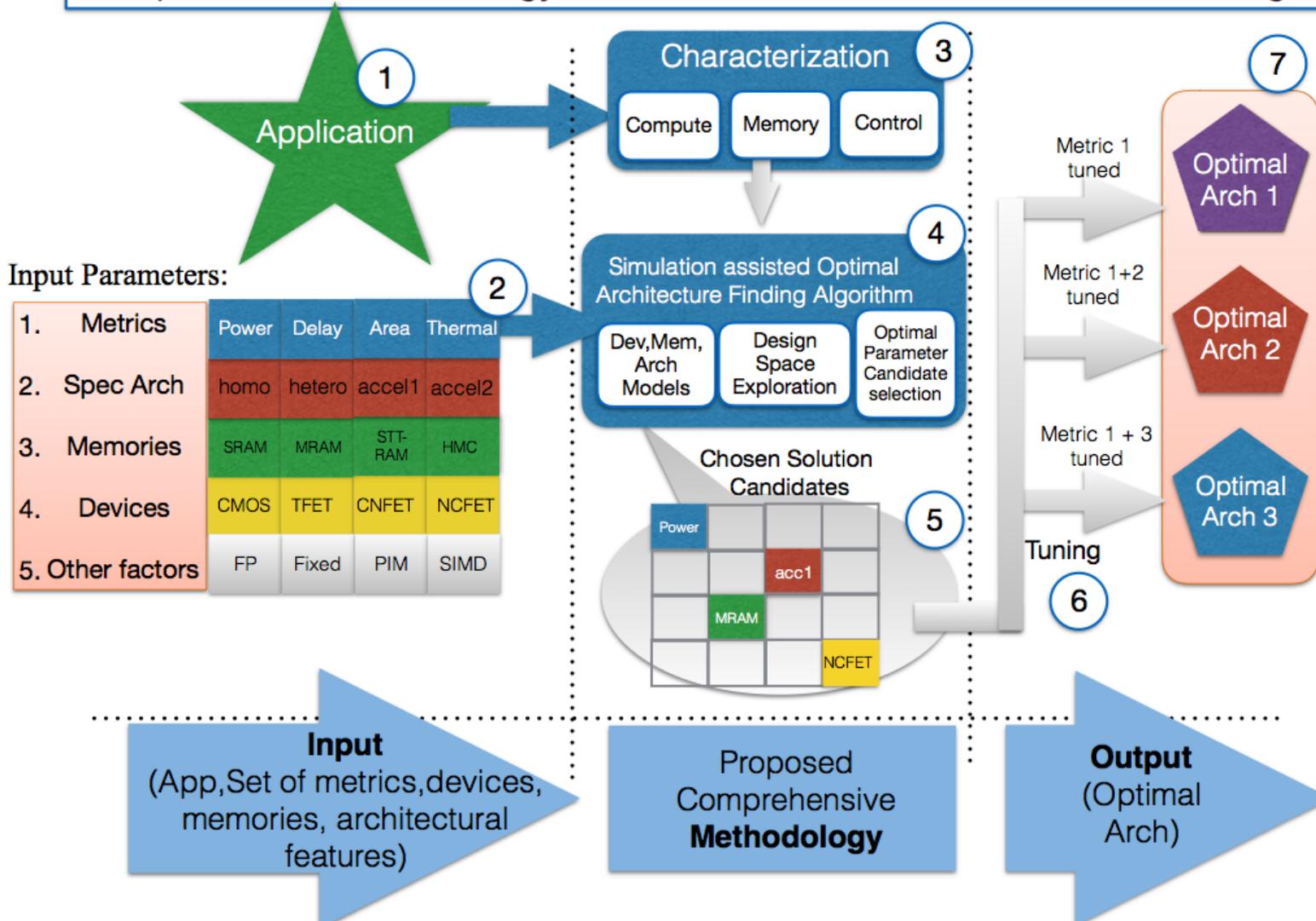
- ◆ Gem5 with Aladdin
- ◆ With small accelerators small delay differences do not have a significant application impact



How To Use These Tools?

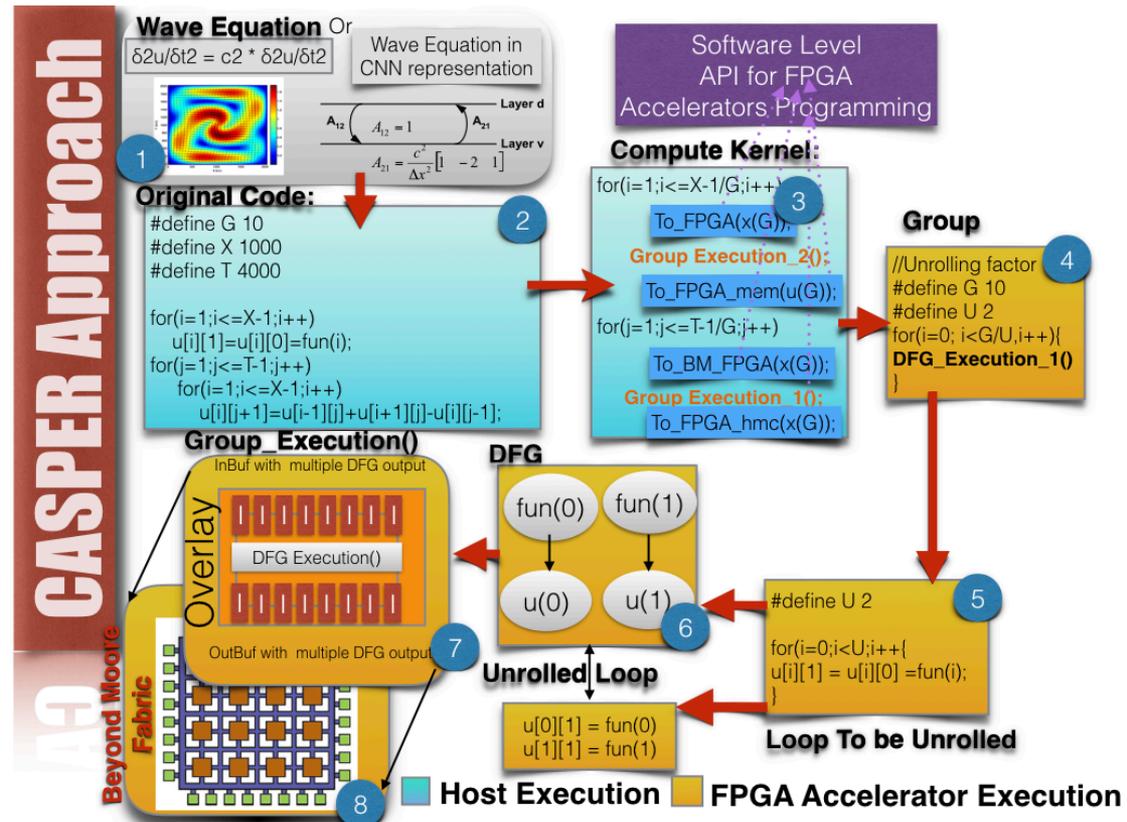


Comprehensive Methodology for Future Devices Based Architecture Design



- ◆ AFM, NCFET, MSET, MRA M based fabric models
- ◆ FPGAs can be heterogeneous too
- ◆ Overlay step understands available FPGA hardware and maps IPs accordingly
- ◆ 50x – 500 performance/energy benefit compared to CMOS FPGAs

End-to-End Open Source Reconfigurable DSE Methodology/Tool Flow for Beyond Moore FPGAs

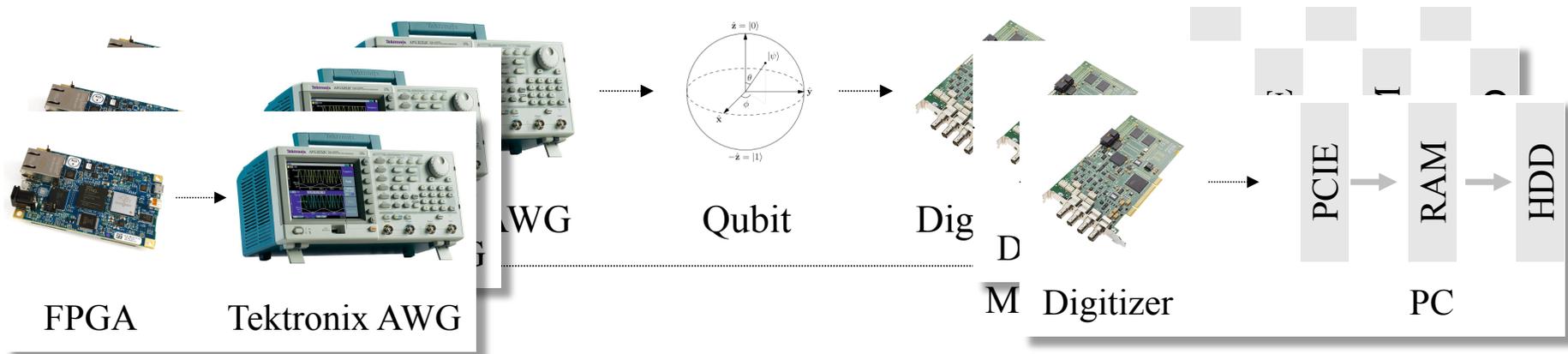


D. Vasudevan et al, "CASPER — Configurable design space exploration of programmable architectures for machine learning using beyond moore devices," 2017

◆ *Quantum Computer = Quantum PU + Control Hardware*

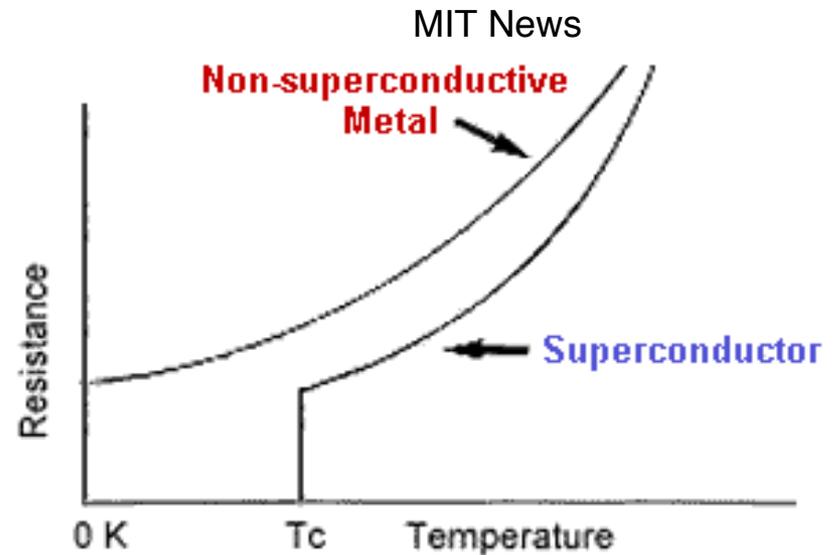
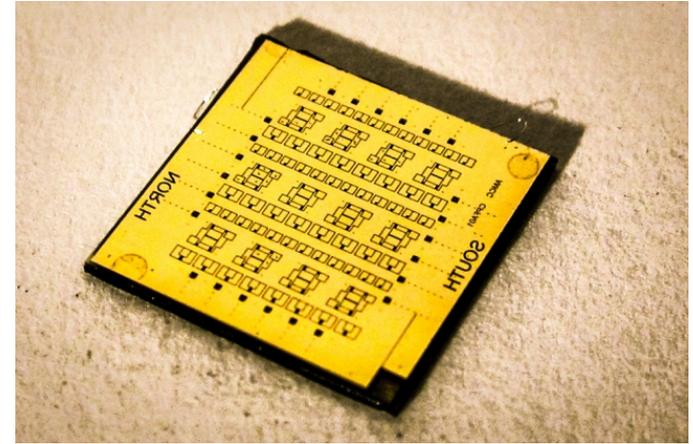
Off the shelf and high cost

Large amount of data and slow speed



*1000 qubits,
gate time 10ns,
3 ops/qubit
300 billion ops per second*

- ◆ Resistance drops to zero
- ◆ 100's of Gigahertz
 - Deep pipelines
- ◆ Memory is a grand challenge
- ◆ Can measure architecture impact and synergy with memory technologies



Gallardo et al "Superconductivity observation in a $(\text{CuInTe}_2)_{1-x}(\text{NbTe})_x$ alloy with $x=0.5$ "



Looking for a PhD Thesis Topic? More Questions to Answer



- ◆ Which device technology will dominate?
 - For what domains, and with what side effects
- ◆ How does architecture change with device technology?
- ◆ How can we best take advantage of deep 3D?
 - With alternating logic and memory layers
- ◆ How large or distant do we make accelerators?
- ◆ How does the memory hierarchy change?
- ◆ How heterogeneous do architectures need to be?



Forewarn Programmers



- ◆ Build an architectural simulation tool that can be used by software developers
- ◆ What is the impact of challenging the far and expensive memory assumption?
 - Also non-volatile
- ◆ What about a heterogeneous memory hierarchy?
- ◆ Can we use reconfigurable accelerators?
- ◆ How to deal with reduced reliability?
 - Approximate computing may see a boost



Conclusion



- ◆ It's an exciting time to be an architect
- ◆ It's hard to predict how digital computing will look like in 20 years
- ◆ Likely more diversified by application domain and even specific algorithm
- ◆ We should focus on a **grand strategy** to best make use of our available options



Questions

