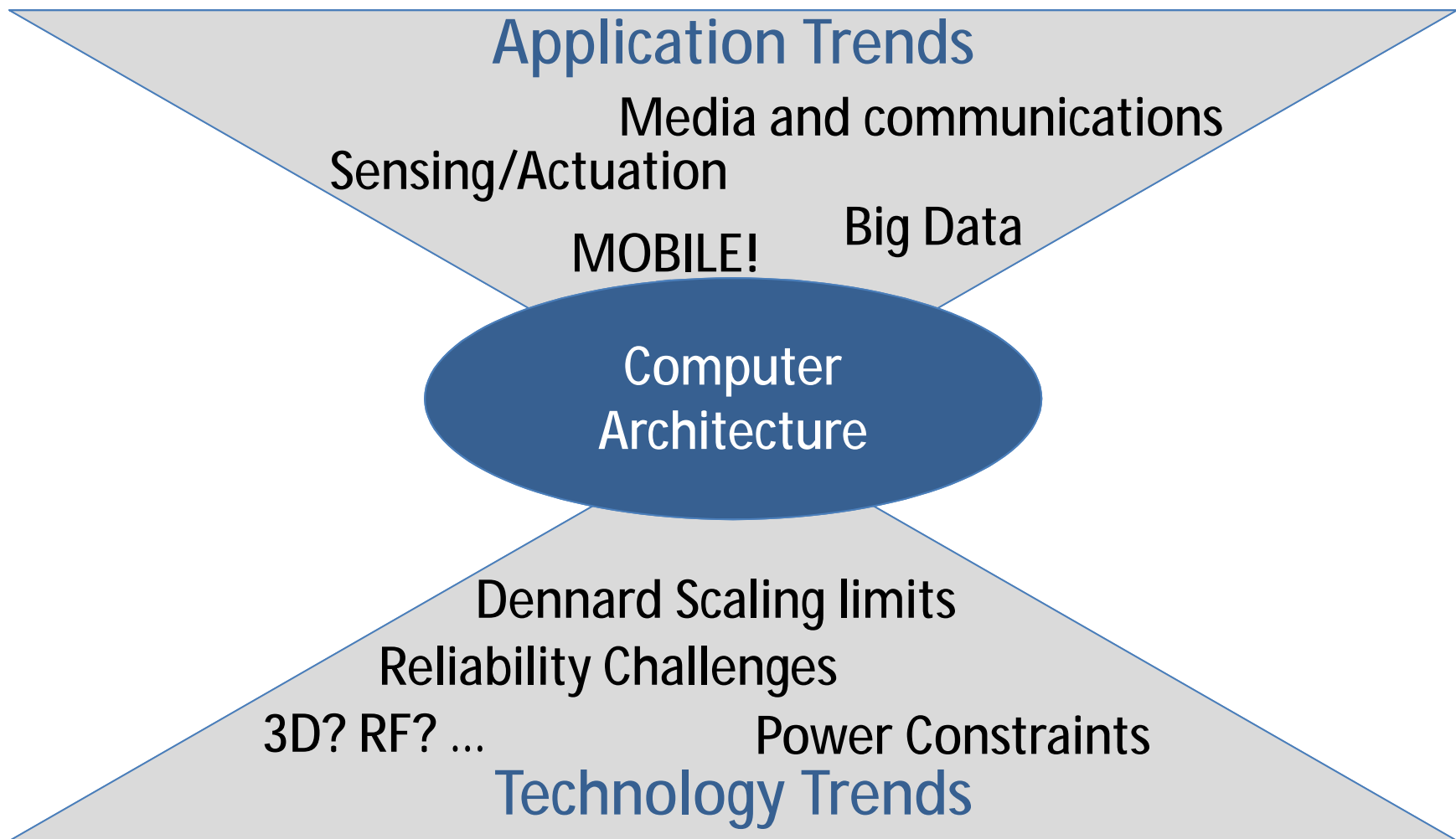


Computer Architecture and the Path of Parallelism and Power Research

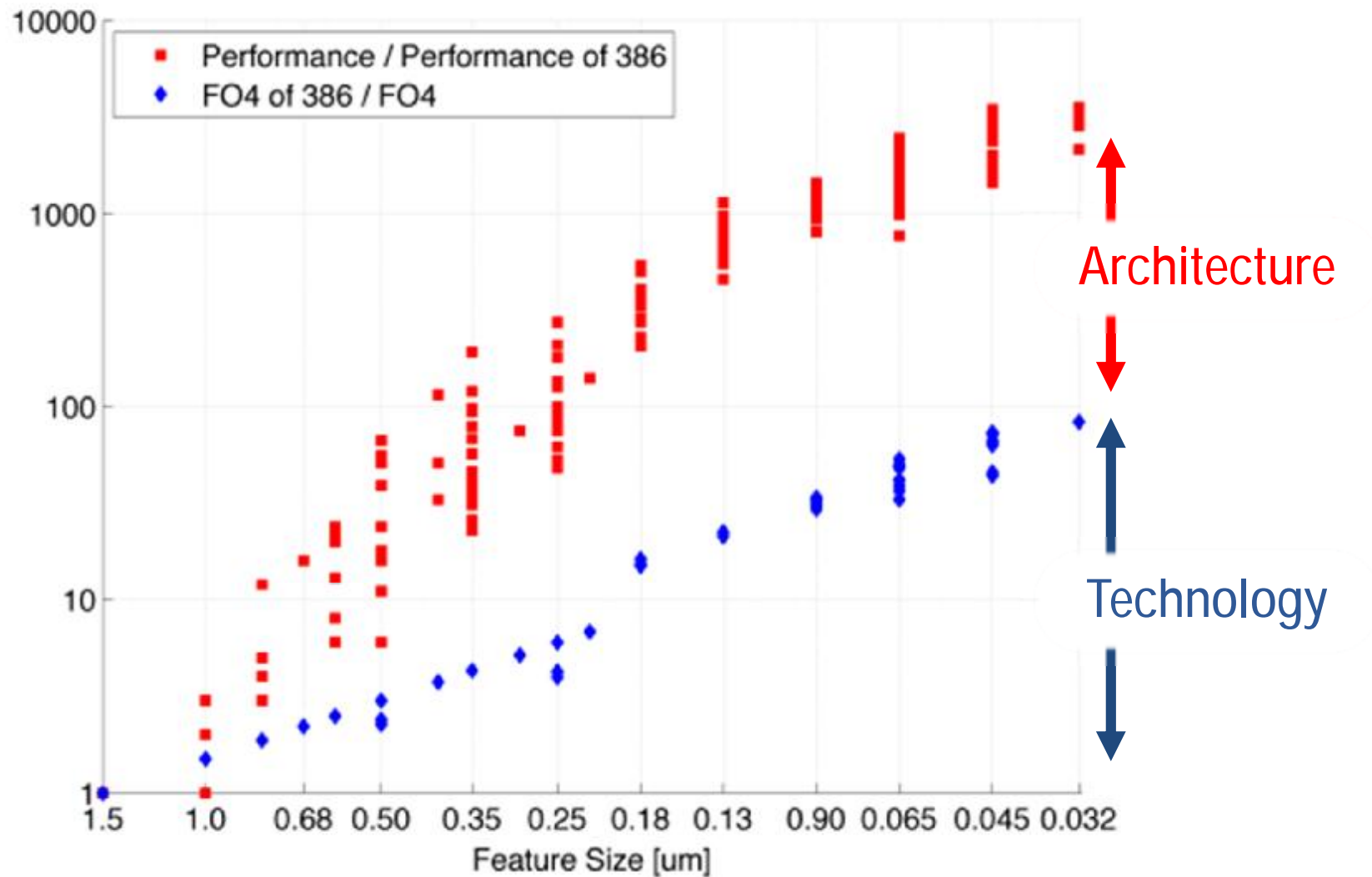
Margaret Martonosi
H. T. Adams '35 Professor of Computer Science
Princeton University



Computer Architecture = Mediator between Technology & Applications

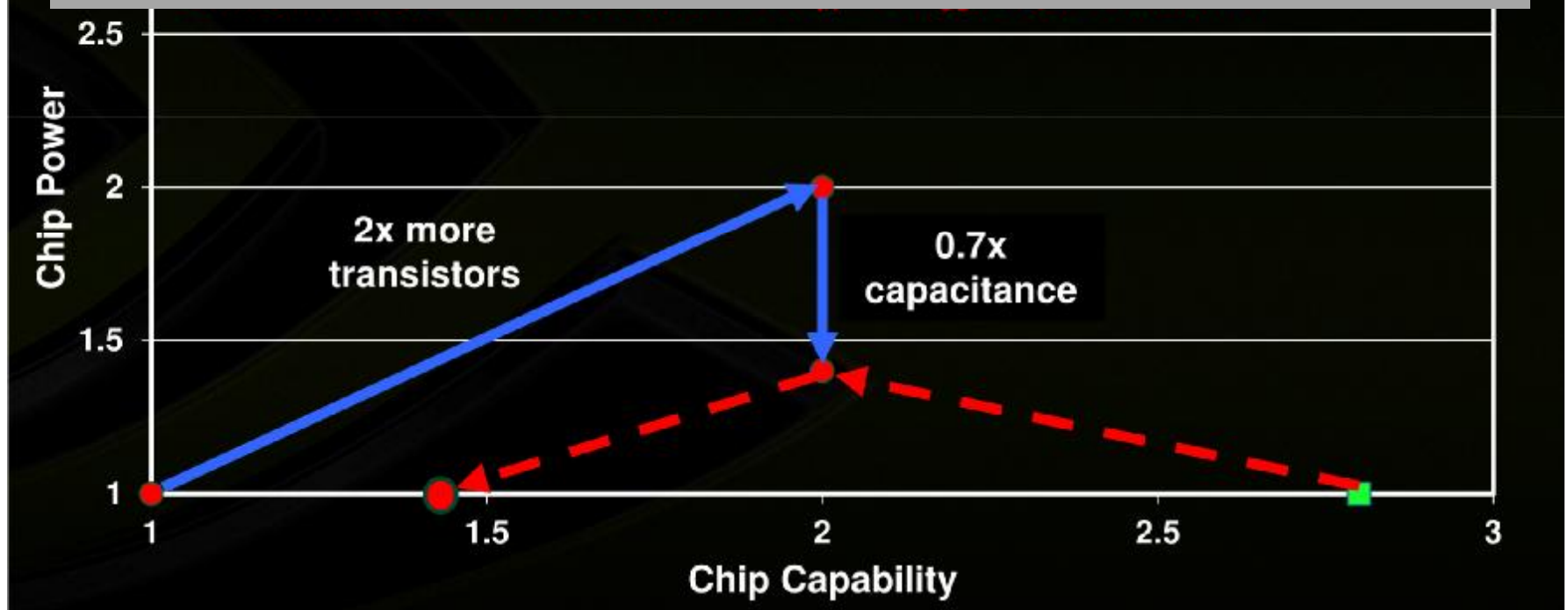


Technology Trends 1



Technology Trends 2 (The bad news...)

For many years, technology was enabler.
But now... Moore's Law & Dennard Scaling slowing =>
Power constraints limit single-core performance scaling.



Upshot

- Application challenges & opportunities:
 - Demand for ever-increasing compute, storage, and communication capabilities...
- Power and thermal challenges:
 - Greatly constrain hardware design alternatives...
 - Direct us to exploit on-chip parallelism.
- *Result: Many parallelism techniques researched decades ago for HPC or “niche” applications are becoming widely-commercialized and mainstream.*



A Timeline of Power & Parallelism Research



Power is not a new problem, but...

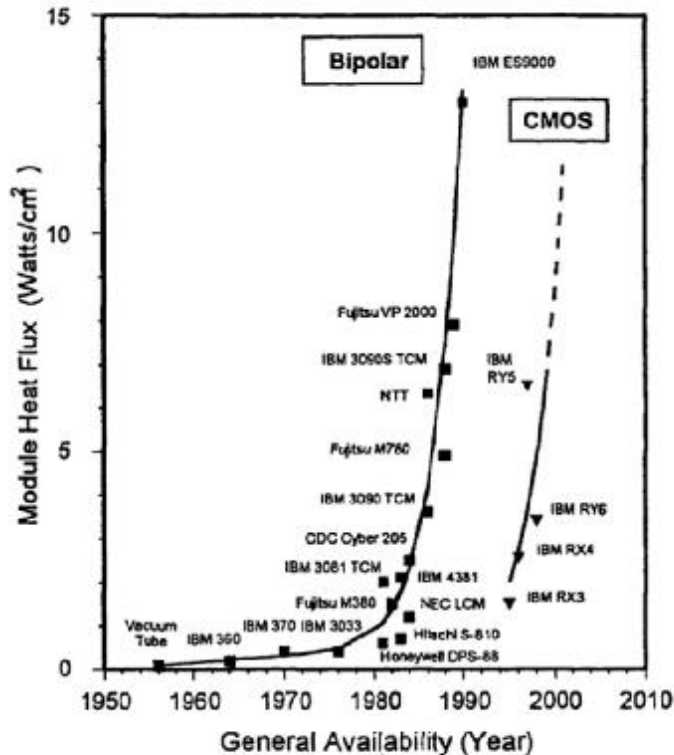


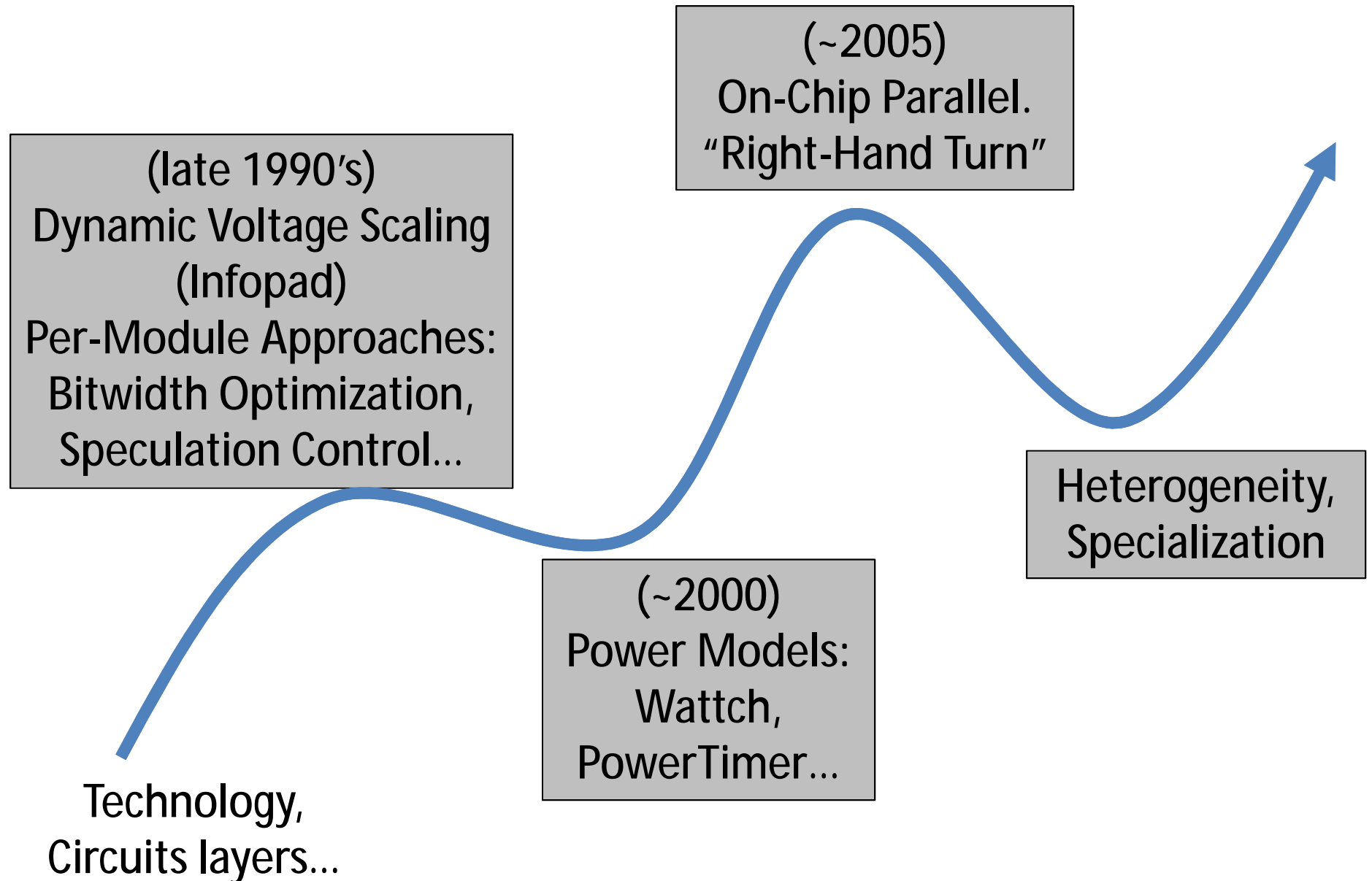
Figure 1. The chronological evolution of module level heat flux in mainframe computers.

Chu et al. SEMITHERM 1999

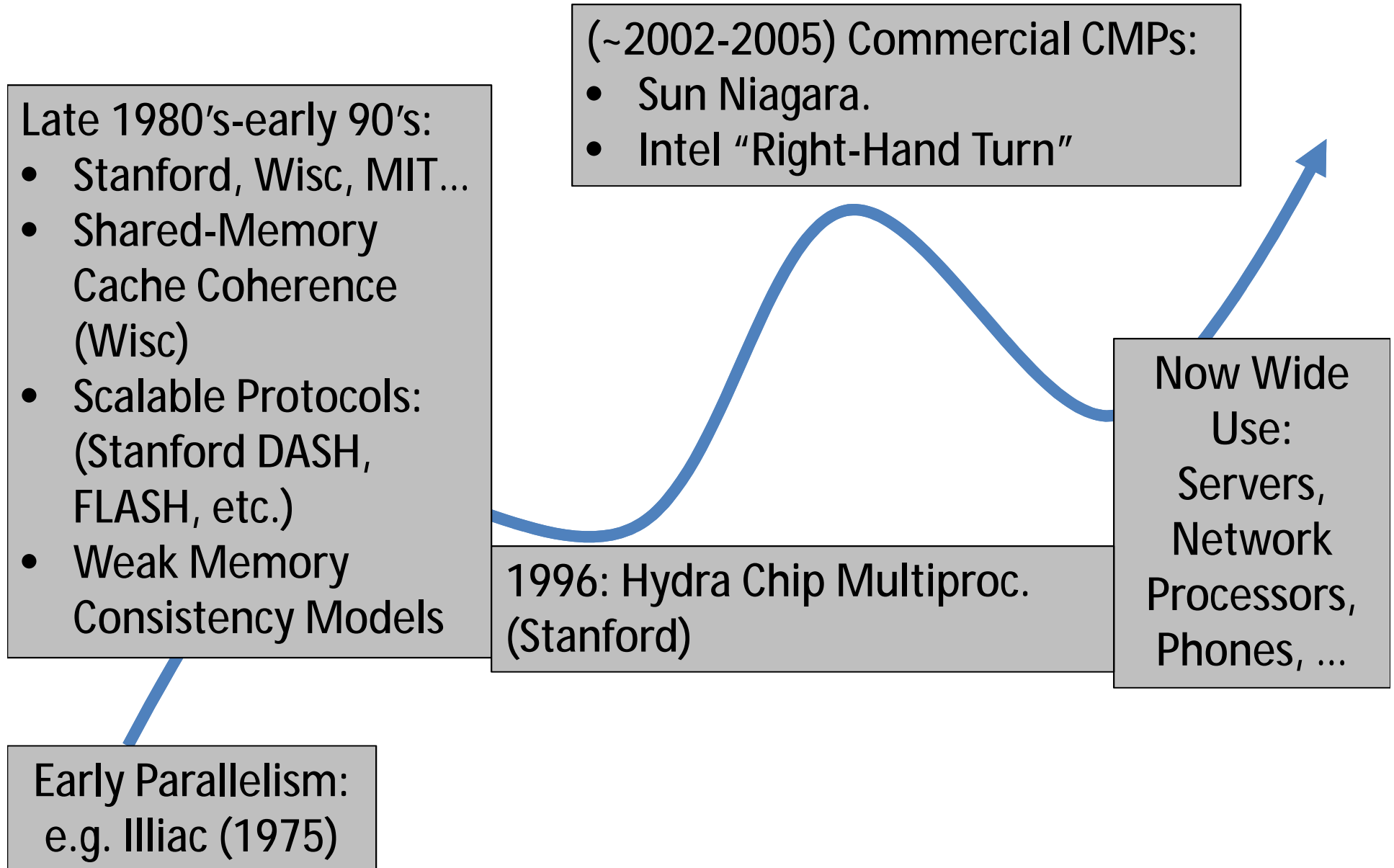
- Computers built from different building blocks over the decades:
 - Vacuum tube... relays... Bipolar transistors...
- Older technologies all reached points where their power was excessive.
- Previous response: Find a new technology and switch to it.
- But now, we don't have something new to switch to!



Architecture's Power Response



Parallelism Timeline Snippets



Parallel Architecture Success Stories

- Shared Memory Parallelism with Cache Coherence
- Weak Memory Consistency Models
- Stanford (DASH, FLASH, ...) scalable shared memory -> SGI Origin -> Many Today...
- On-Chip Parallelism Hydra -> Sun Niagara -> many.
- Distributed Shared Memory -> Many
 - MIT Raw -> Tiler, ...
 - Cavium network processors
- MIT/Stanford Multithreading -> NVIDIA GPUs
- Simultaneous Multithreading -> Intel Hyperthreading.
- Speculative Lock Elision (Wisconsin) -> Intel Transactional Memory

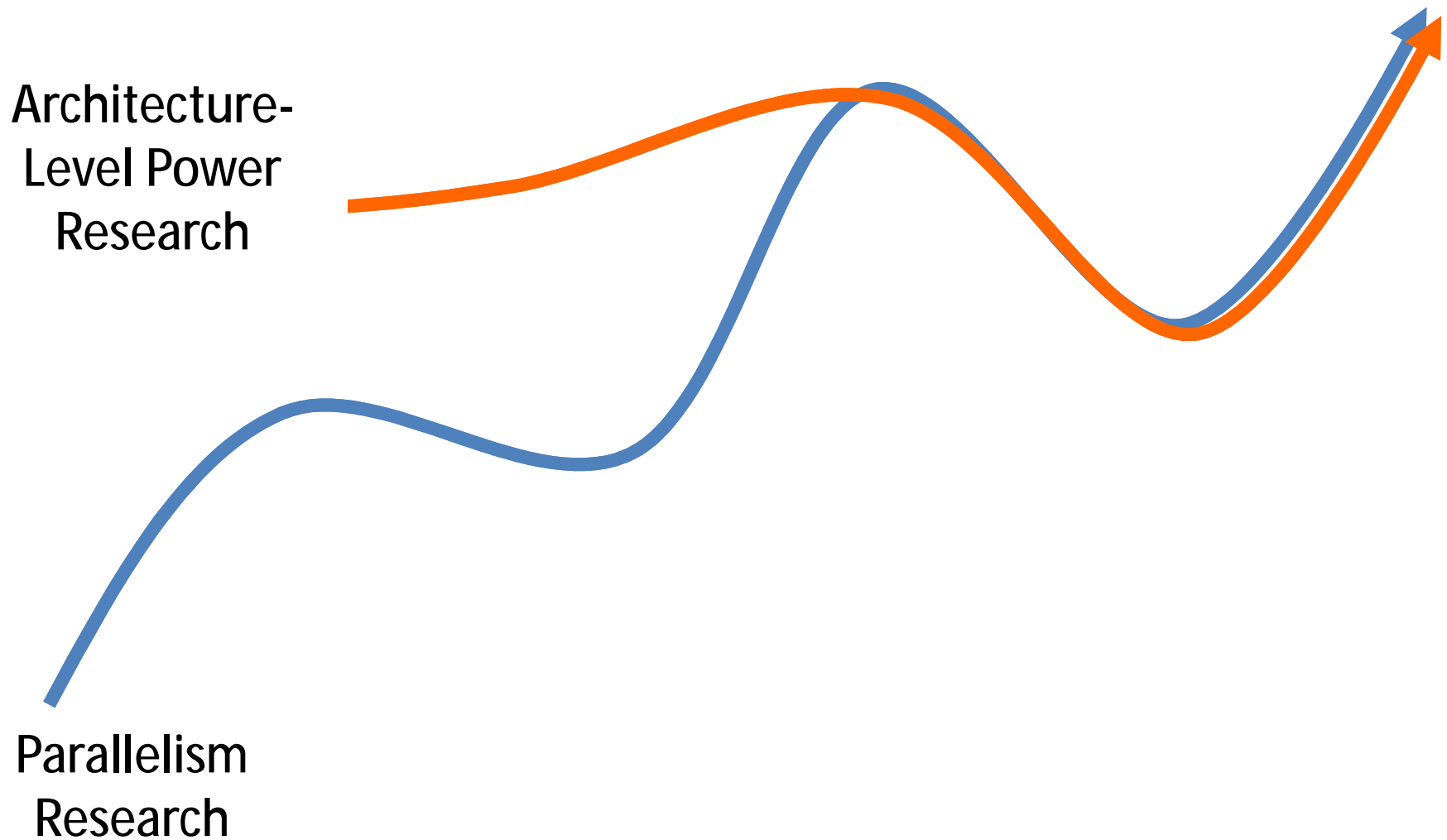


Power-Aware Architecture Research: Success Stories

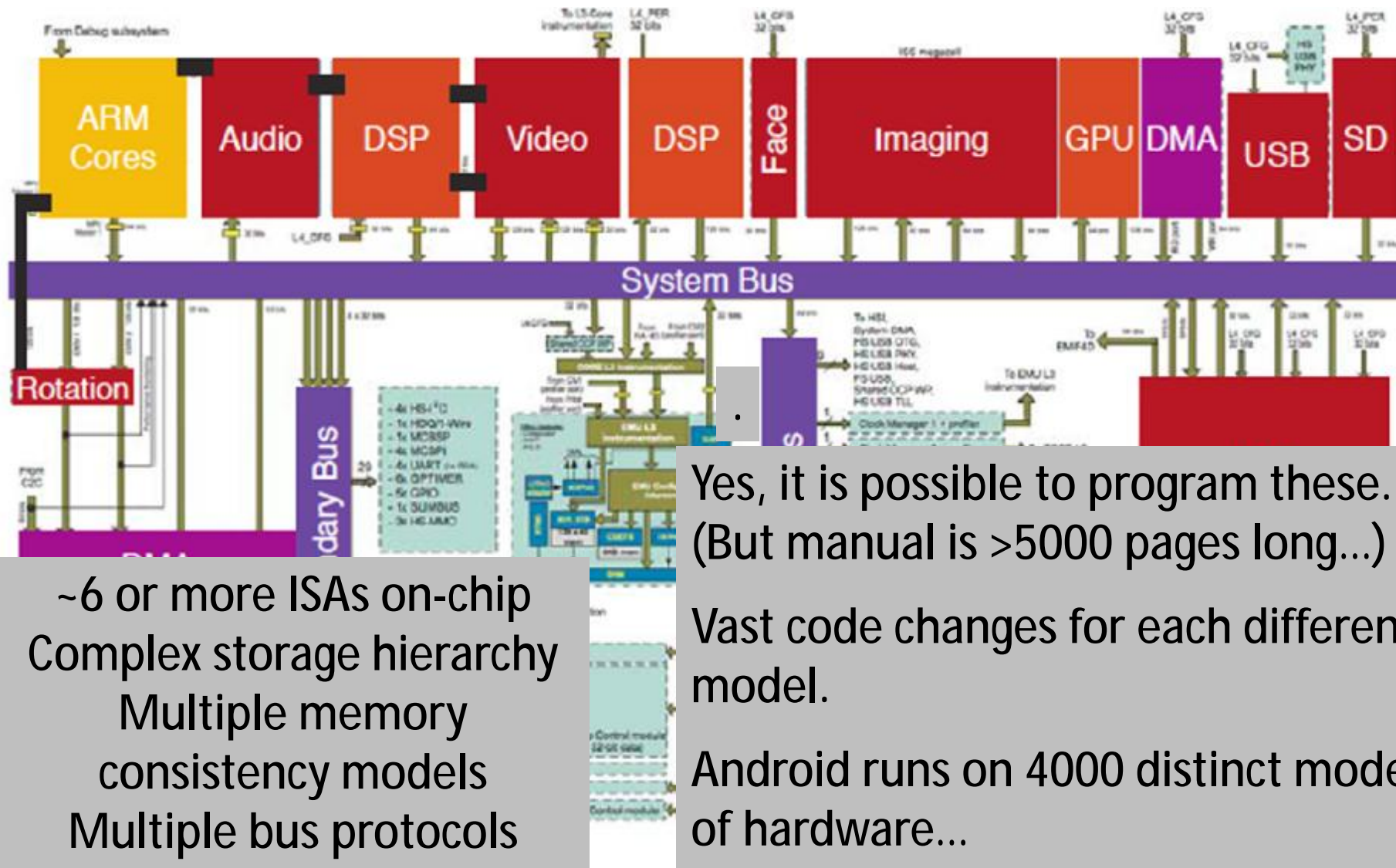
- Power/Thermal Modeling (Wattch, HotSpot, ...):
 - Early-stage design insights: Princeton Wattch-> IBM PowerTimer -> Shaped choice of IBM pipeline depth
 - Enabled Power/Performance Pareto frontier analysis -> drove single-ISA heterogeneity (e.g. UCSD) -> ARM big.Little approach.
- Dynamic Voltage Scaling: Berkeley Infopad (DARPA-funded) -> thousands of other techniques -> Industry approaches ACPI and others.
- Thermal-aware research: Academic research on Dynamic Thermal Management -> mechanisms like Intel TurboBoost, allowing thermally constrained systems to extract more performance for a given thermal budget
Module-Level Optimizations:
 - Narrow bitwidth optimization (and other value-based power optimizations) -> Intel ALUs
 - Cache Memory Leakage Energy Management -> Intel, ARM
- Specialization: Reconfigurable computing research in 90's (DARPA ACS)->
 - Accelerator designs in today's mobile SoCs
 - Microsoft Catapult specialization in datacenters.



Converging Trends + Research Payoffs



Which brings us here...



Computer Architecture = Hardware + HW/SW Abstraction

G. M. Amdahl
G. A. Blaauw
F. P. Brooks, Jr.,

Architecture of the IBM System/360

Abstract: The architecture* of the newly announced IBM System/360 features four innovations:

1. An approach to storage which permits and exploits very large capacities, hierarchies of speeds, read-only storage for microprogram control, flexible storage protection, and simple program relocation.
2. An input/output system offering new degrees of concurrent operation, compatible channel operation, data rates approaching 5,000,000 characters/second, integrated design of hardware and software, a new low-cost, multiple-channel package sharing main-frame hardware, new provisions for device status infor-

*The term *architecture* is used here to describe the attributes of a system as seen by the programmer, i.e., the conceptual structure and functional behavior, as distinct from the organization of the data flow and controls, the logical design, and the physical implementation.

Introduction

The design philosophies of the new general-purpose machine organization for the IBM System/360 are discussed in this paper.† In addition to showing the architecture* of the new family of data processing systems, we point out the various engineering problems encountered in attempts to make the system design compatible, at the program bit level, for large and small models. The compatibility was to extend not only to models of any size but also to their various applications—scientific, commercial, real-time, and so on.

*The term *architecture* is used here to describe the attributes of a system as seen by the programmer, i.e., the conceptual structure and functional behavior, as distinct from the organization of the data flow and controls, the logical design, and the physical implementation.
†Additional details concerning the architecture, engineering design, programming, and application of the IBM System/360 will appear in a series of articles in the *IBM Systems Journal*.

The section that follows describes the objectives of the new system design, i.e., that it serve as a base for new technologies and applications, that it be general-purpose, efficient, and strictly program compatible in all models. The remainder of the paper is devoted to the design problems faced, the alternatives considered, and the decisions made for data format, data and instruction codes, storage assignments, and input/output controls.

Design objectives

The new architecture builds upon but differs from the designs that have gradually evolved since 1950. The evolution of the computer had included, besides major technological improvements, several important systems concepts and developments:



Entering A Post-ISA World

- ISAs still useful operationally, but have little/no relevance as an abstraction layer.
 - Due to Power/performance constraints, chips like the Apple A8 devote ~half their area to accelerators that have no ISA.
 - Vendors increasingly hide their ISAs under other abstraction layers: NVIDIA PTX vs. SASS.
- Disruptive moment: HW/SW interface on which our whole computer systems infrastructure is based, is undergoing a seismic change.
- Why won't industry solve this?
 - Modest abstractions... But messy and short-term: TI OMAP4 software manual >5000 pages long. SW changes required to map to OMAP5.
⇒ Software dev costs € , and software reliability/security €
 - Software and hardware vendors are often separate companies.
 - Even within a single processor, IP from several companies.
 - Programmers and compilers face this Tower of Babel completely unshielded.
- Consider DoD:
 - Military systems compose many heterogeneous parts from many vendors.
 - Without portable abstractions, DoD is beholden to single vendors...
 - Extreme difficulties of creating a system or requalifying it for a different platform.
 - Likelihood of correct and high-performance code?

Going Forward

1. Managing heterogeneity for a Post-ISA World
 - New Abstractions for portable and nimble, reliable software systems.
 - Multi-ISA, multi-chip, multi-location!
2. Communication as First-Class Partner with Computation
3. Maintain and expand architect's role as mediator between application and technology trends.



Computer Architecture and the Path of Parallelism and Power Research

Margaret Martonosi
Princeton University

I happily acknowledge the contributions of my grad students, co-authors, and funding agencies to much of my work. Thanks also to David Brooks, Mark Hill, Trevor Mudge, Shubu Mukherjee, Michael Pellauer and others for providing feedback and information on these slides.



For more information...

- [1] S. Fuller and L. Millett, "The Future of Computing Performance: Game Over or Next Level?," The National Academy Press, 2011
(http://books.nap.edu/openbook.php?record_id=12980&page=R1).
- [2] J. Torrellas and M. Oskin, "Failure is not an Option: Popular Parallel Programming," Workshop on Advancing Computer Architecture Re-search, August 2010 (http://www.cra.org/ccc/docs/ACAR_Report_Popular-Parallel-Programming.pdf).
- [3] M. Oskin and J. Torrellas, "Laying a New Foundation for IT: Computer Architecture for 2025 and Beyond," Workshop on Advancing Computer Architecture Research, September 2010
(<http://www.cra.org/ccc/docs/ACAR2-Report.pdf>).
- [4] M. Hill and C. Kozyrakis, "Advancing Computer Systems without Technology Progress," ISAT Outbrief, April 17-18, 2012, of DARPA/ISAT Workshop, March 26-27, 2012
(http://www.cs.wisc.edu/~markhill/papers/isat2012_ACSWTP.pdf).
- [5] A community white paper, "21st Century Computer Architecture," Computing Community Consortium, May 25, 2012
(<http://cra.org/ccc/docs/init/21stcenturyarchitecturewhitepaper.pdf>).





Stanford Dash Multiprocessor

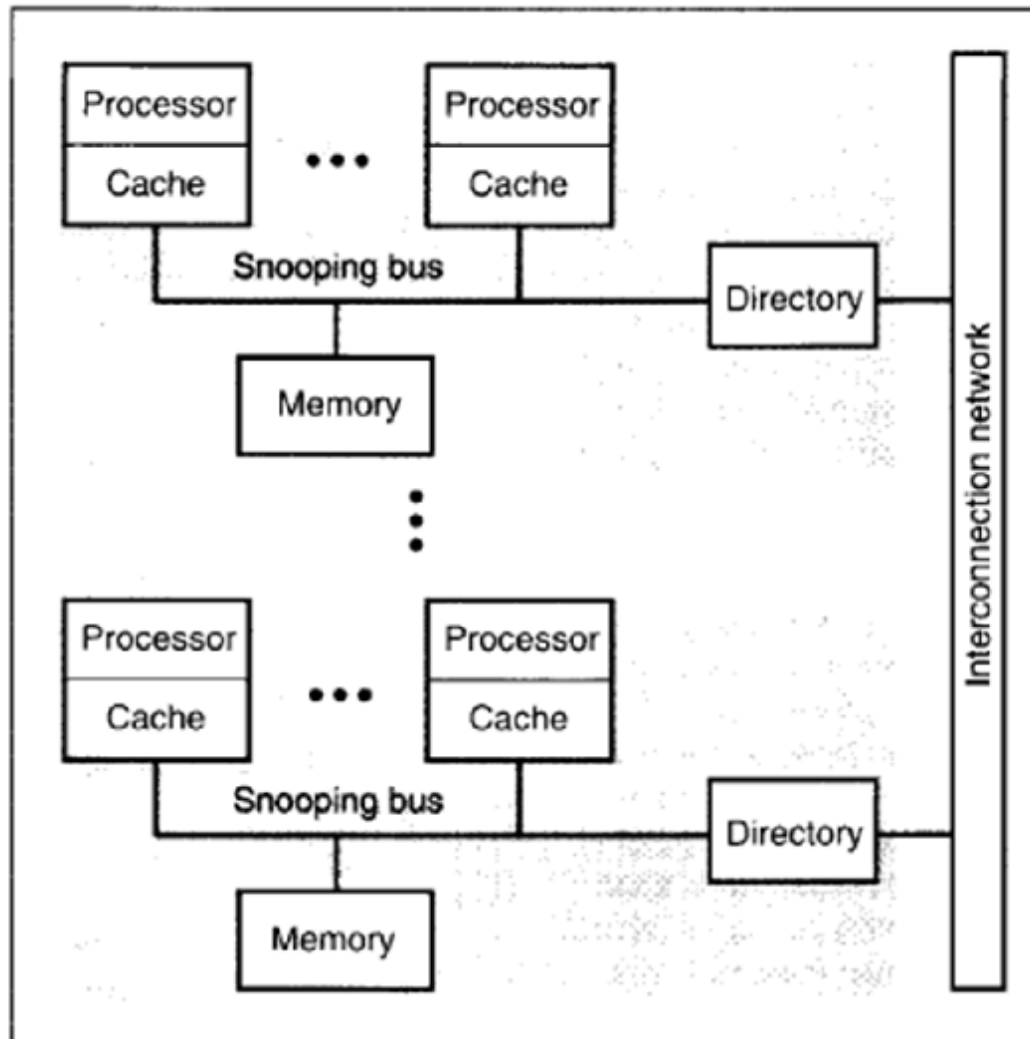
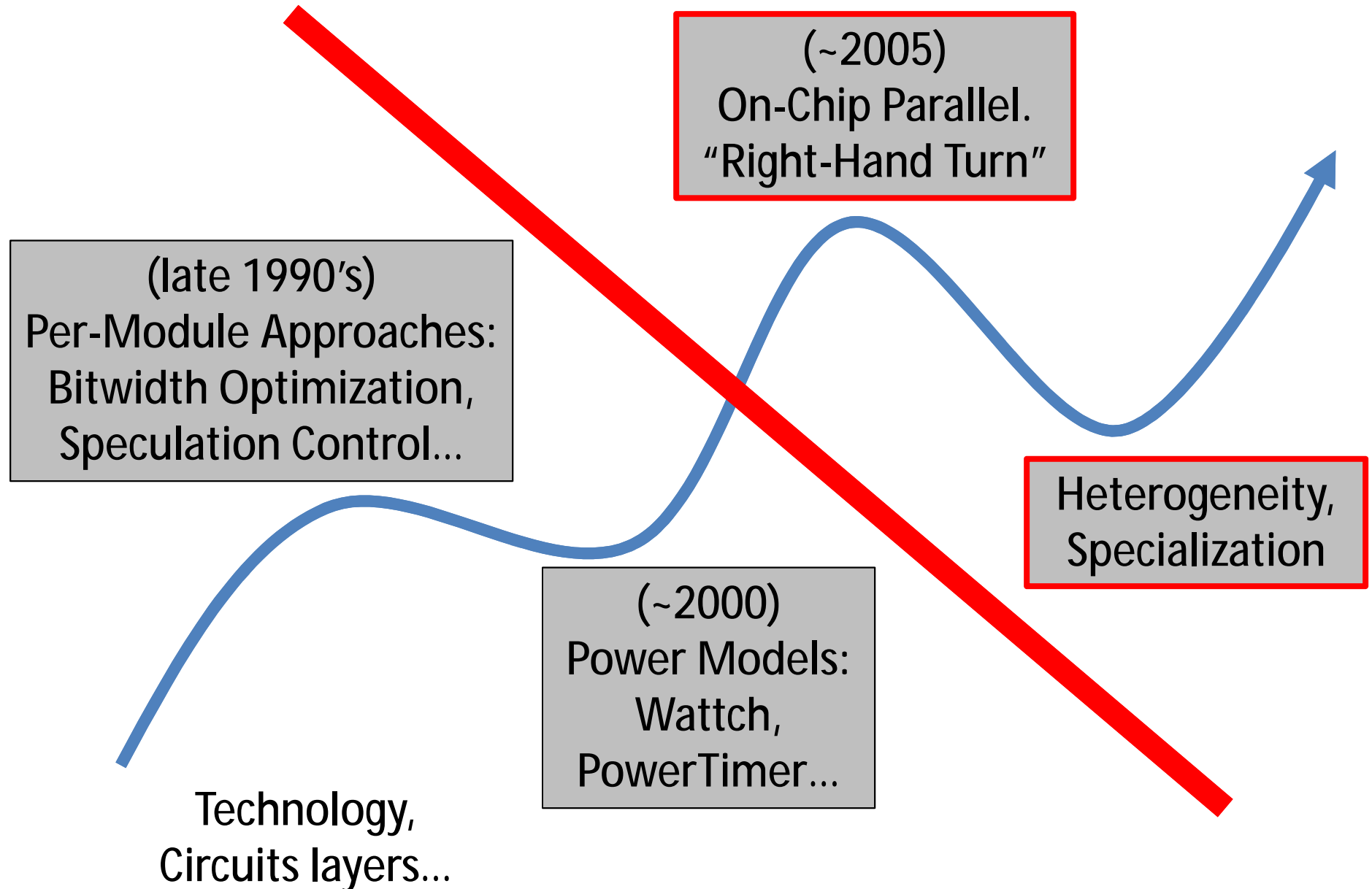


Figure 1. The Dash architecture consists of a set of clusters connected by a general interconnection network. Directory memory contains pointers to the clusters currently caching each memory line.



Architecture's Power Response: **SW View**



How old is the power problem?

FOR RELEASE SATURDAY A.M., FEBRUARY 16, 1946

For Radio Broadcast after
7:00 P.M., EST, February 15, 1946

PHYSICAL ASPECTS, OPERATION OF ENIAC ARE DESCRIBED

The ENIAC (Electronic Numerical Integrator and Computer) is a large scale electronic general purpose computing machine. It occupies a room 30 by 50 feet in size. It weighs 30 tons and has 100 feet of front panels.

This machine is the most intricate and complex electronic device in the world, requiring for its operation 18,000 electronic tubes. Some idea of the machine's complexity can be gained when it is compared with an average radio, which has ten tubes, the largest radar set having 400 tubes and the B-29 bomber with less than 800 tubes. Included in its equipment are 500,000 resistors, 70,000 capacitors, and 10,000 capacitors.

While there are no multitude of vacuum tubes create some noise. This is its specially designed air conditioned proving grounds. Extra are two comparatively small punched cards and receive

The forty main panel panels on each leg and 8 panels arranged from 1

1. Control and Init
2. Cycling Unit
- 3, and 4. Master
- 5, and 6. First Fun
- 7, and 8. Accumulators 1 and 2
9. Divider and Square Rooter
- 10-17. Accumulators 3 - 10
- 18-20. Multiplier
- 21-28. Accumulators 11 - 18
- 29, and 30. Second Function Table
- 31, and 32. Third Function Table
- 33, and 34. Accumulators 19 and 20
- 35 - 37. Constant Transmitters
- 38-40. Printer

The ENIAC consumes 150 kilowatts. This power is supplied by a three-phase regulated, 240-volt, 60-cycle power line. The power consumption may be broken up as follows; 80 kilowatts for heating the tubes 45 kilowatts for generating d.c. voltages, 20 kilowatts for driving the ventilator blower and 5 kilowatts for the auxiliary card machines.

The ENIAC consumes 150 kilowatts... The power consumption may be broken up as follows; 80 kilowatts for heating the tubes 45 kilowatts for generating d.c. voltages, 20 kilowatts for driving the ventilator blower and 5 kilowatts for the auxiliary card machines.



Source: Original ENIAC press release. Feb, 1946