

The Future of Computing Performance: *Game Over or Next Level?*

Kathy Yelick, Study committee member

based on slides from

Samuel H. Fuller, Chair

Lynette I. Millett, Study Director

See <http://www.cstb.org> for the agenda of the 3/22 symposium on this report.

Slides will be posted.

Computer Science and Telecommunications Board (CSTB)

National Research Council (NRC)

THE NATIONAL ACADEMIES



Committee On Sustaining Growth In Computing Performance

THE NATIONAL ACADEMIES

Working Group of Experts Addressed the Problem

- **SAMUEL H. FULLER**, Analog Devices Inc., Chair
- **LUIZ ANDRÉ BARROSO**, Google, Inc.
- **ROBERT P. COLWELL**, Independent Consultant
- **WILLIAM J. DALLY**, NVIDIA Corporation and Stanford University
- **DAN DOBBERPUHL**, P.A. Semi
- **PRADEEP DUBEY**, Intel Corporation
- **MARK D. HILL**, University of Wisconsin–Madison
- **MARK HOROWITZ**, Stanford University
- **DAVID KIRK**, NVIDIA Corporation
- **MONICA LAM**, Stanford University
- **KATHRYN S. McKINLEY**, University of Texas at Austin
- **CHARLES MOORE**, Advanced Micro Devices
- **KATHERINE YELICK**, University of California, Berkeley

Staff

- **LYNETTE I. MILLETT**, Study Director
- **SHENAE BRADLEY**, Senior Program Assistant

The Future of Computing Performance: Game Over or Next Level? 3/22 Symposium

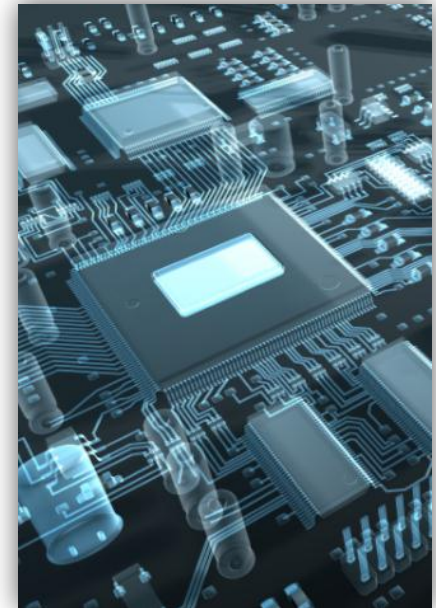


- **Welcome and Report Overview**, Samuel H. Fuller, Analog Devices, Inc. Committee Chair
- **Session 1 – Parallelism & Innovative Programming Models, Algorithms, & Languages**
 - Moderator: Kathryn McKinley, University of Texas, Austin
 - David Grove, IBM; Keshav Pingali, University of Texas, Austin; Guy Steele, Oracle; Katherine Yelick, University of California, Berkeley
- **Perspective on Investment and Resources to Support Continuing Innovation in Computing Performance**
 - David Liddle, U.S. Venture Partners
- **Session 2 – Computing in a Power Constrained WorldPanel**
 - Moderator: Mark A. Horowitz, Stanford University
 - Bob Dennard, IBM; Dan Dobberpuhl, Consultant; Kevin Nowka, IBM; Partha Ranganathan, HP
- **Session 3 – Reaching the Next Level in Computer Science & Engineering Education**
 - Panel Moderator: Mark D. Hill, University of Wisconsin, Madison
 - Guy Blelloch, Carnegie Mellon University; Dan Ernst, University of Wisconsin Eau Claire; David Kirk, NVIDIA; Marcia Linn, University of California, Berkeley
- **Session 4 – Exploring the Terrain: Research Directions, Priorities, and Strategies**
 - Panel Moderator: Samuel H. Fuller
 - Susanne Hambrusch, NSF; Norm Jouppi, HP; Keith Marzullo, NSF; Bill Harrod, DOE

Sustaining Growth in Computing Performance



- Is there a need for continued growth in computing performance?
- What is computer performance?
- What is limiting growth in computing now?
- Can new programming methods that address these challenges be developed and broadly deployed ?
- Recommendations in research, practice and education



What do we mean by Computing Performance?



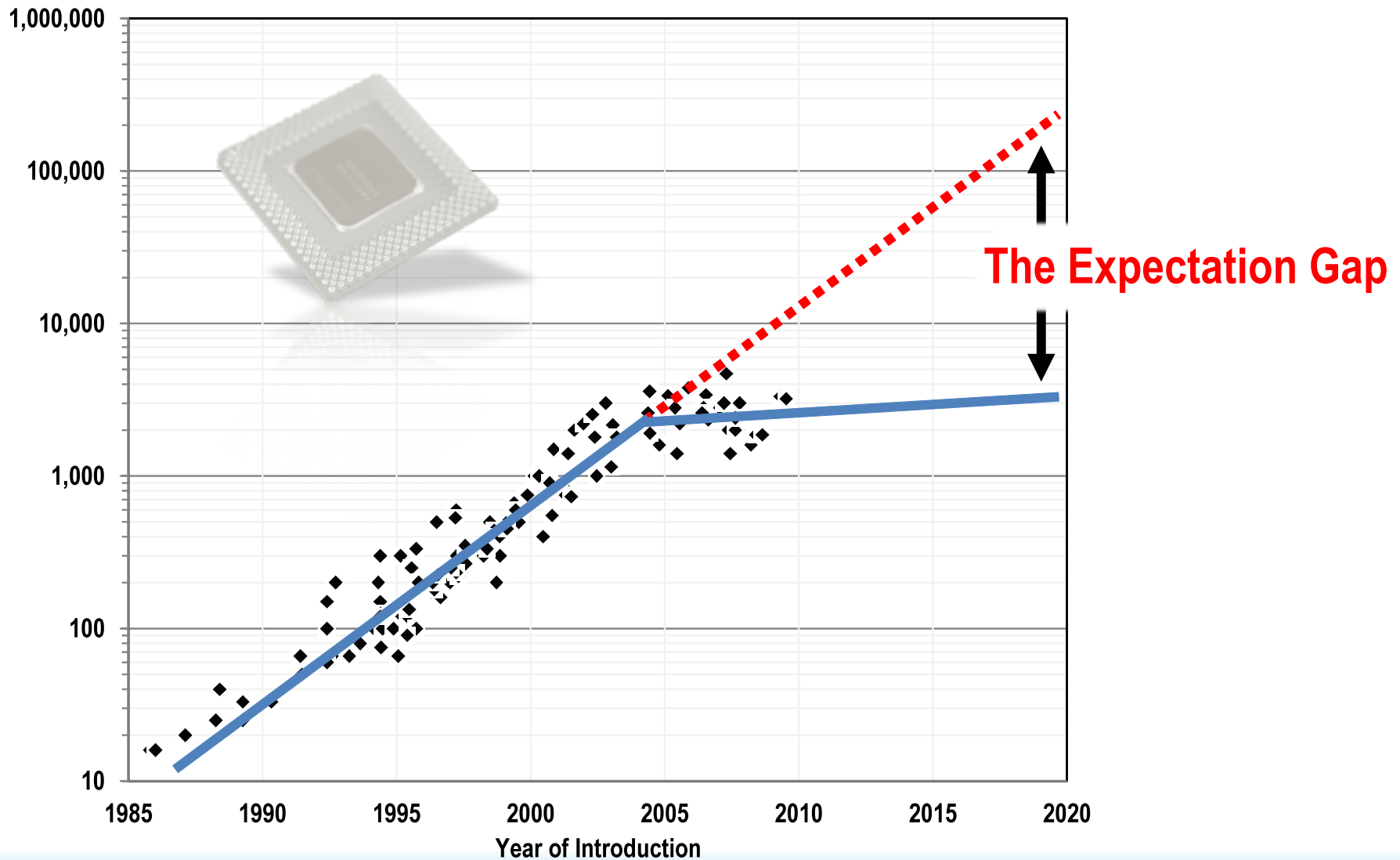
- One measure of single-processor performance is the product of clock rate times instructions per cycle: i.e. raw speed is instructions/sec.
- Computing 'speed' is fungible and can be traded for almost any feature one might want
 - Higher reliability, e.g., error detection/correction
 - Background operations e.g., indexing, compression, decompression
 - Redundancy
 - Near real-time translation
 - Image resolution
 - Signal fidelity
 - I/O bandwidth
- Delivered performance requires balance of processing performance, storage capacity and interconnect bandwidth.



Processor Performance Plateaued Around the Year 2004



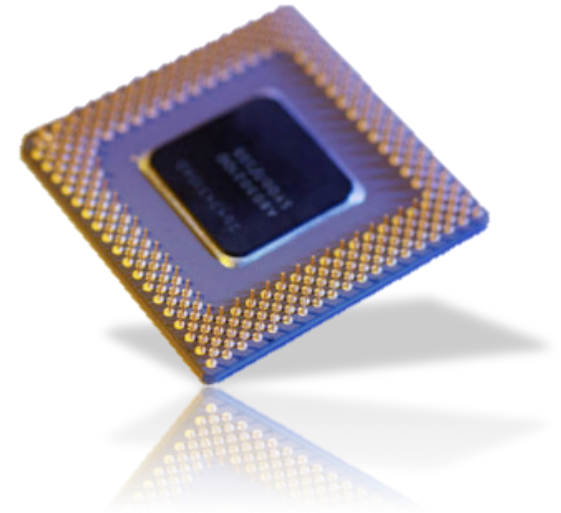
Microprocessor Performance “Expectation Gap” over Time (1985-2020 projected)



“Yes, we know” but Exponential Assumptions Persist



- Even among experts, hard to dislodge an implicit assumption of continuing exponential performance improvements
- “Moore’s law, which the computer industry now takes for granted, says that the processing power and storage capacity of computer chips double or their prices halve roughly every 18 months.” – The Economist, February 2010
- “the software and other custom features become extremely important in constructing a computing system that can take advantage of the intrinsically higher speed provided by Moore’s law of increasing power per chip.” – Defense Science Board, “Advanced Computing”, March 2009 [arguing for parallelism, but still assuming “intrinsically higher speed”]



Classic CMOS Dennard Scaling: the Science behind Moore's Law



Scaling:

Voltage: V/α

Oxide: t_{ox}/α

Wire width: W/α

Gate width: L/α

Diffusion: x_d/α

Substrate: αN_A

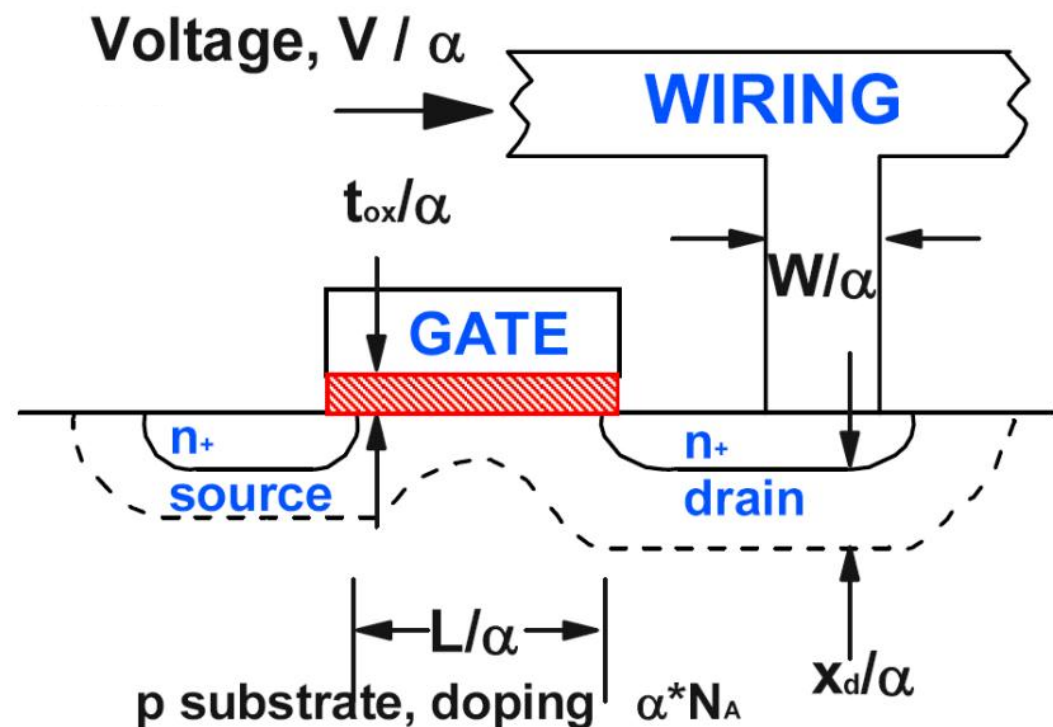
Results:

Higher Density: $\sim\alpha^2$

Higher Speed: $\sim\alpha$

Power/ckt: $1/\alpha^2$

Power Density: $\sim\text{Constant}$



R. H. Dennard et al.,
IEEE J. Solid State Circuits, (1974).

Root Cause is the Fundamental Physical Limitation of Heat-Density per Unit Area of CMOS Semiconductors



Why Has Power/Chip Skyrocketed?

- CMOS threshold voltage (V_t) of at least 200 to 300 millivolts is needed to make it a good switch
 - On current (drive current) must be high for fast switching
 - Off current (leakage current) must be low to minimize power
- Supply voltage (V_{dd}) needs to be 3+ times V_t to enable good digital switch performance
- Therefore, V_{dd} is limited to 0.8 to 0.9 volts, or higher
- Power = $\mathbf{C f V_{dd}^2}$



Root Cause is the Fundamental Physical Limitation of Heat-Density per Unit Area of CMOS Semiconductors



Post Dennard CMOS Scaling Rule

Scaling:

Voltage: ~~V/α~~ V

Oxide: t_{ox}/α

Wire width: W/α

Gate width: L/α

Diffusion: x_d/α

Substrate: αN_A

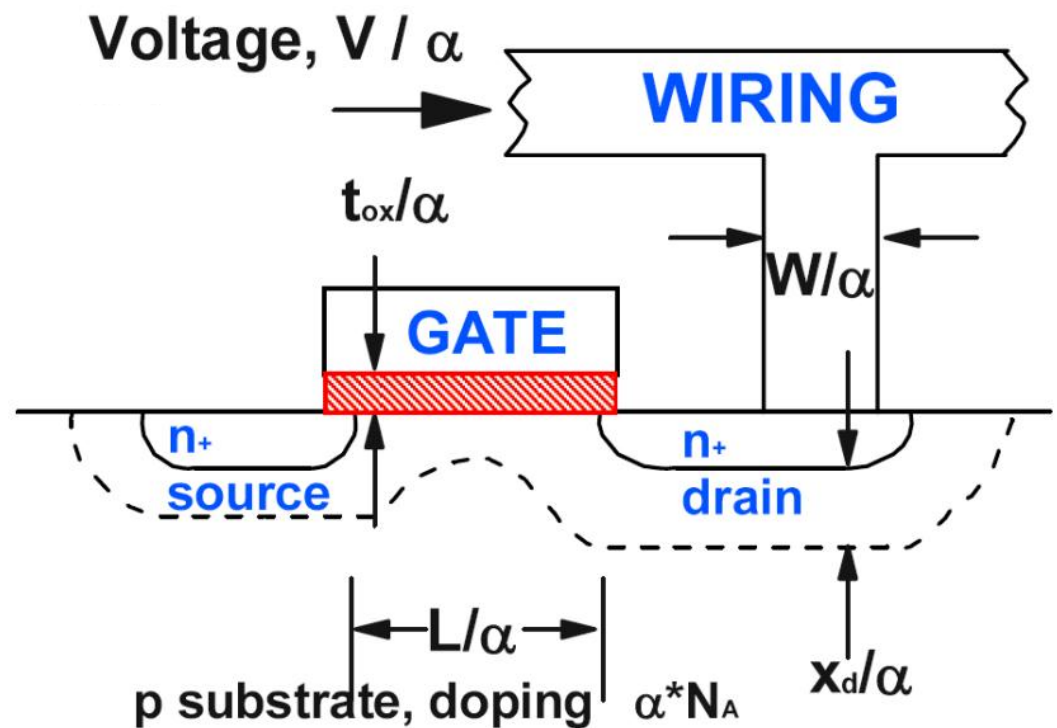
Results:

Higher Density: $\sim \alpha^2$

Higher Speed: $\sim \alpha$

Power/ckt: ~~$1/\alpha^2$~~ 1

Power Density: ~~$\sim \text{Constant}$~~ α^2



R. H. Dennard et al.,
IEEE J. Solid State Circuits, (1974).

Alternatives of CMOS

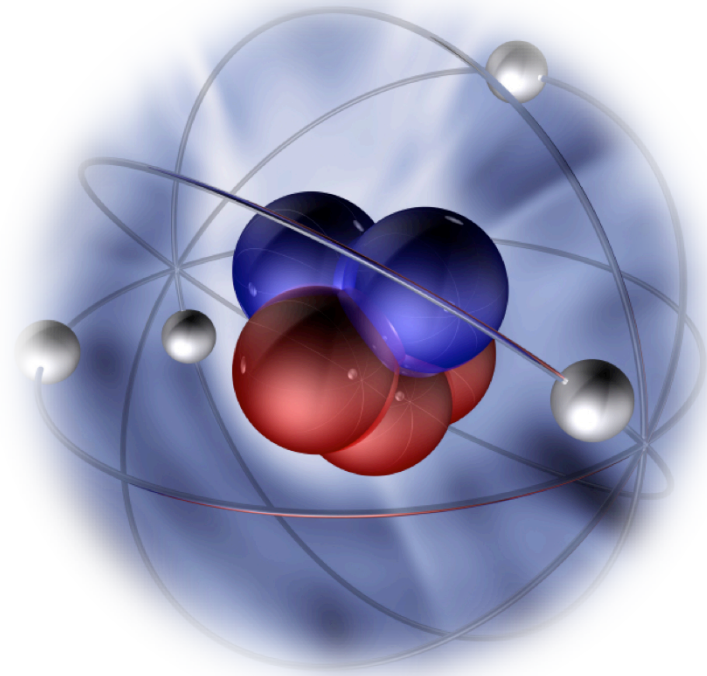


Near Term – But Limited Relief to Power Constraints

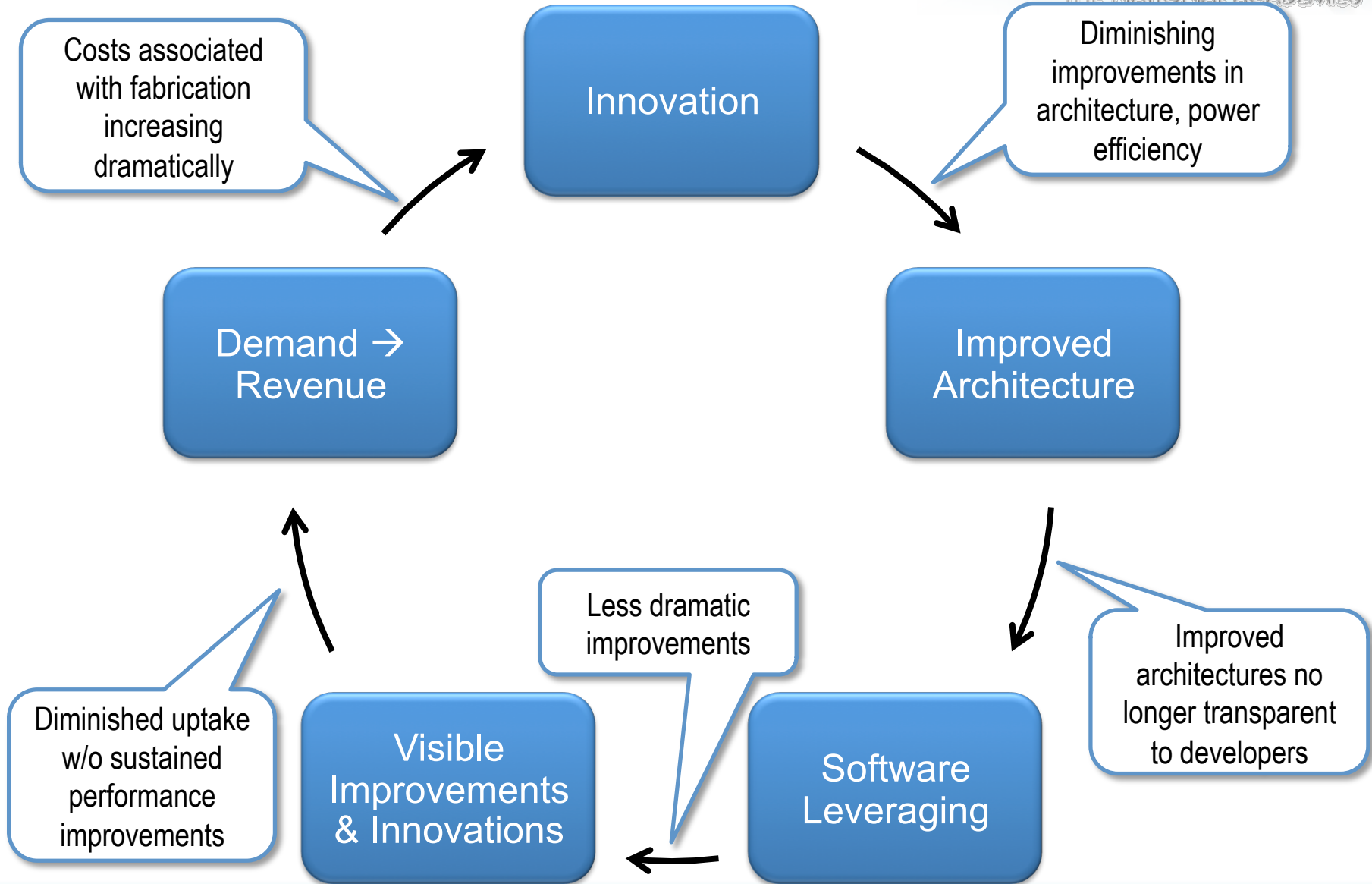
- III – V materials for MOSFETs, e.g. GaAs
- Carbon nanotubes or grapheme based devices

Longer Term – Much Work Required to Bring to Commercial Reality

- Electron spin, versus electron charge., i.e. Spintronics
- Quantum devices



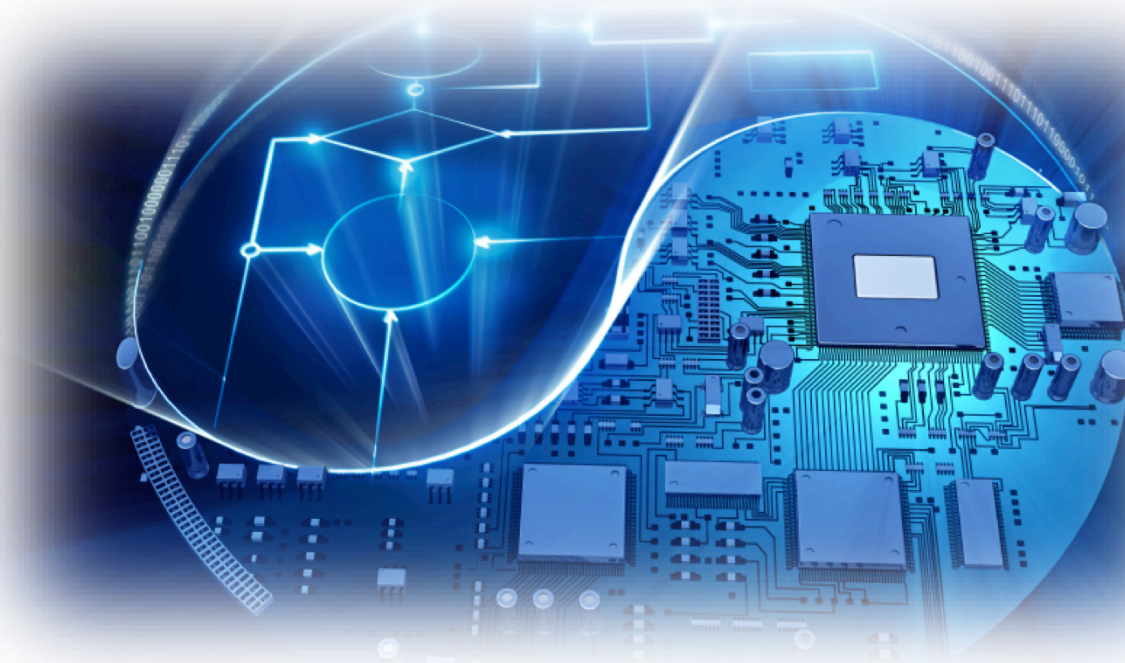
Cracks in the Virtuous Cycle





Recommendations

From the Committee on Sustaining Growth in Computing Performance



Summary of Recommendations



Place a much greater on improvements and innovations in parallel processing

1. Invest in **algorithms** that can exploit parallel processing
2. Invest in **programming methods** to enable efficient use of parallel systems
3. Focus long-term efforts on **rethinking of the canonical computing “stack”** in light of parallelism and resource-management challenges
4. Invest in **parallel architectures driven by applications**, including enhancements of chip multiprocessor systems, data-parallel architectures, application-specific architectures, and radically different approaches
5. Invest to make computer systems more **power efficient**
6. Promote cooperation and innovation of **open interface standards** for parallel programming
7. Invest in **tools and methods** to transform legacy apps to parallel systems
8. Increased emphasis on **parallelism in computer science education**

Highlights of the Symposium



- Bob Denard on Denard Scaling and Dan Doberpuhl on processor design
 - 250W PC, Compute gets most
 - Laptop 50W, Compute gets most
 - SmartPhone 2-3W Compute gets 1W (constrained by heat)
- Kevin Nowka, IBM on cloud computing challenges
 - Deployment model, involves sharing of resources
 - Someone (not Kevin) predicted that power use by Cloud will drop
 - What we need for HPC is a direct analogy to what we need for Clouds
 - 50X in sustained performance/\$
 - 20X improvement in sustained performance/Watt
- Marcia Lin, UC Berkeley
 - Students bring many preconceptions to programming; study what works rather than asserting what will work
- Funding agencies
- David Liddle on VC perspective

VC Perspective: David Liddle, US Venture Partners



- Processor industry was running at "maneuvering speed"
- Traditional funding sources will not work
 - Federal agency grants: ARPA net: time scales and costs were different
 - Mission oriented agency have little ability to provide sustained exploratory funding and require near-term deliverables
 - NSF general needs to provide small grants to a very large constituency with a low "hit rate"
 - No sympathy for the IT industry in Washington, as it is seen as robust compared to others
 - Hard in an industrial lab
 - Next product segments are in the direction of lower cost and larger markets, resulting in step changes in component cost/performance
 - It is extremely hard for industrial hardware labs to make a determined change
 - The changes flow of computer science software talent
 - The best no longer to to companies that sell software they go to companies that sell services online
 - This significantly biases the kinds of problems they work on
- Venture Capital Perspective
 - Faster Better Cheaper (FBC) vs Brave New World (BNW)
 - FBC = 95%; Known market, entrenched competitors, fails in the lab cheaply
 - BNW = 5%; Latent market, no competitors, fails in the marketplace expensively
 - VCs are looking for big impact (big \neq worthy)
 - You can't invest in a programming model or a science project

Algorithms and Software Recommendations



1. **Invest in research in and development of algorithms** that can exploit parallel processing
2. **Invest in research in and development of programming methods** that will enable efficient use of parallel systems not only by parallel-systems experts but also by typical programmers
3. **Focus long-term efforts on rethinking of the canonical computing “stack”** in light of parallelism and resource-management challenges
 - Applications
 - Programming language
 - Compiler
 - Runtime
 - Virtual machine
 - Operating system
 - Hypervisor
 - Architecture



Architecture Recommendations



- 4. Invest in research on and development of parallel architectures driven by applications**, including enhancements of chip multiprocessor systems and conventional data-parallel architectures, cost-effective designs for application-specific architectures, and support for radically different approaches



Power Efficiency Recommendation



5a. Invest in research and development to make computer systems more power efficient at all levels of the system, including software, application-specific structures, and alternative devices.

R&D efforts should address ways in which software and system architectures can improve power efficiency, such as by **exploiting locality** and the use of **domain-specific execution units**.

5b. R&D should also be aimed at making **logic gates more power-efficient**. Such efforts should address alternative physical devices beyond incremental improvements in today's CMOS circuits.



Practice and Education Recommendations



6. To promote **cooperation and innovation** by sharing and encouraging development of **open interface standards** for parallel programming rather than proliferating proprietary programming environments.
7. Invest in the **development of tools and methods** to transform legacy applications to parallel systems.
8. **Incorporate in computer science education** an increased emphasis on parallelism, and use a variety of methods and approaches to prepare students better for the types of computing resources that they will encounter in their careers.

