
Computing's Energy Problem:

(and what we can do about it)

Mark Horowitz

Stanford University

horowitz@ee.stanford.edu

Everything Has A Computer Inside



The Reason is Simple: Moore's Law Made Gates Cheap

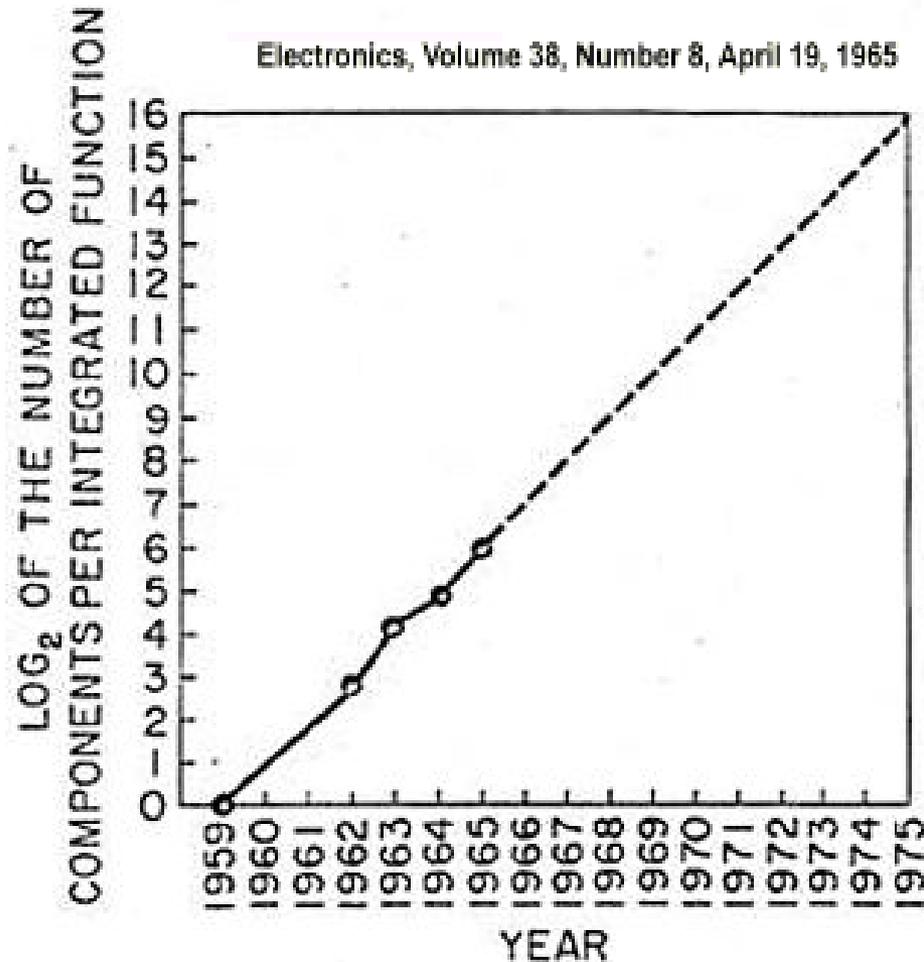
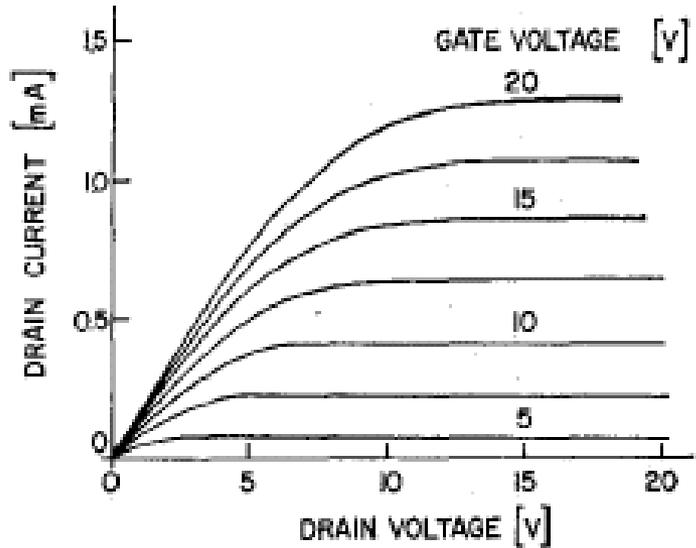


Fig. 2 Number of components per integrated function for minimum cost per component extrapolated vs time.

Dennard's Scaling Made Them Fast & Low Energy

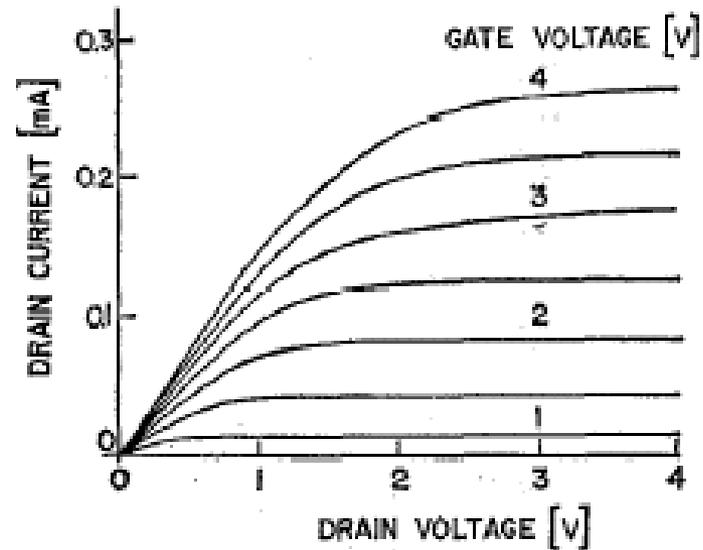


$$t_{ox} = 1000 \text{ \AA}$$

$$L = W = 5 \mu\text{m}$$

$$V_{sub} = -7 \text{ V}$$

$$\psi_s = 0.65 \text{ V}$$



$$t'_{ox} = 200 \text{ \AA}$$

$$L' = W' = 1 \mu\text{m}$$

$$V'_{sub} = -1 \text{ V}$$

$$\psi'_s = 0.73 \text{ V}$$

The triple play:

- Get more gates,
- Gates get faster,
- Energy per switch

$$1/L^2$$

$$1/\alpha^2$$

$$CV/i$$

$$\alpha$$

$$CV^2$$

$$\alpha^3$$

Dennard, JSSC, pp. 256-268, Oct. 1974

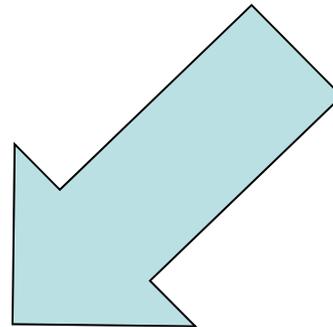
Our Expectation

Cray-1: world's fastest computer 1976-1982

- 64Mb memory (50ns cycle time)
- 40Kb register (6ns cycle time)
- ~1 million gates (4/5 input NAND)
- 80MHz clock
- 115kW

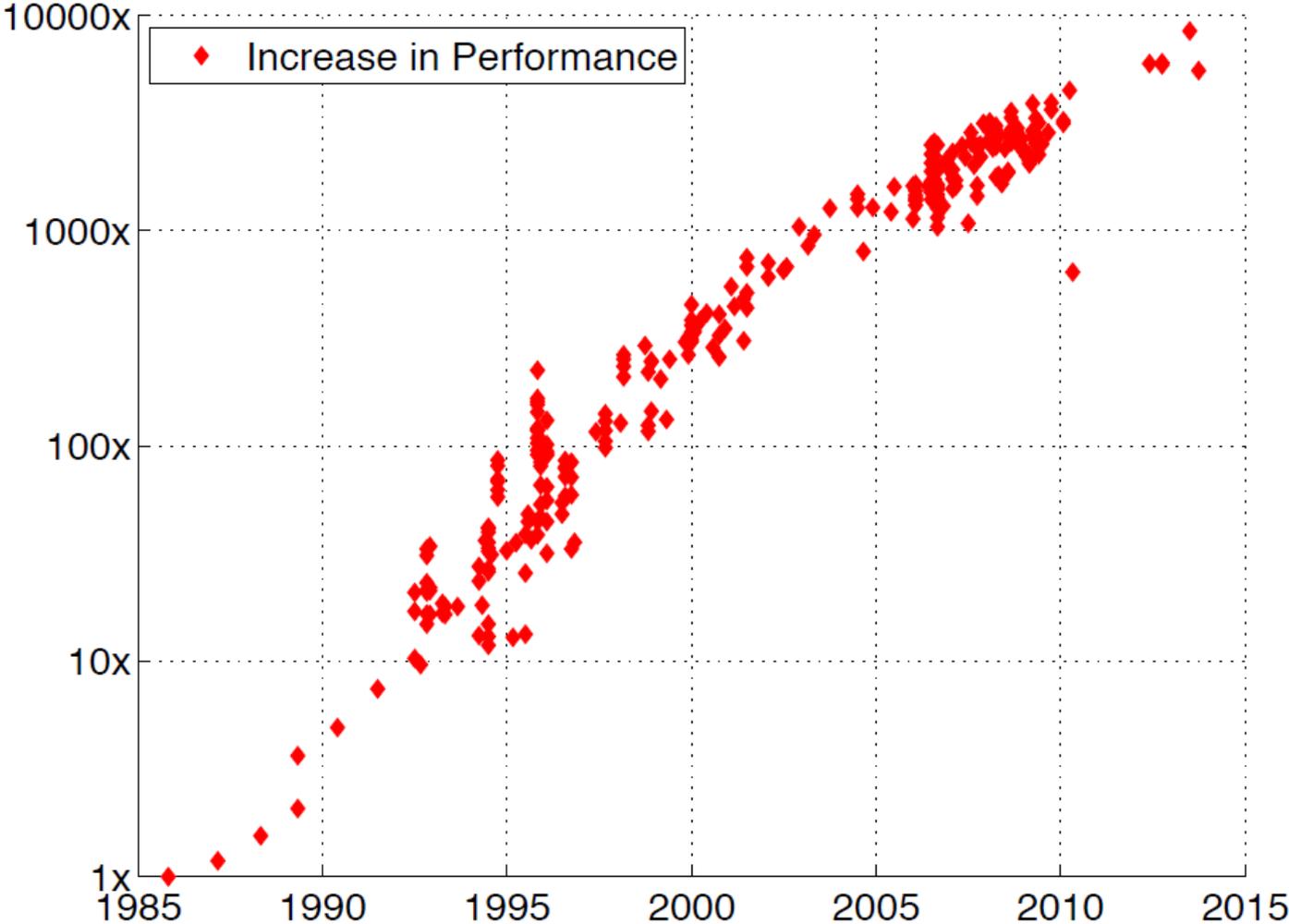
In 45nm (30 years later)

- $< 3 \text{ mm}^2$
- $> 1 \text{ GHz}$
- $\sim 1 \text{ W}$



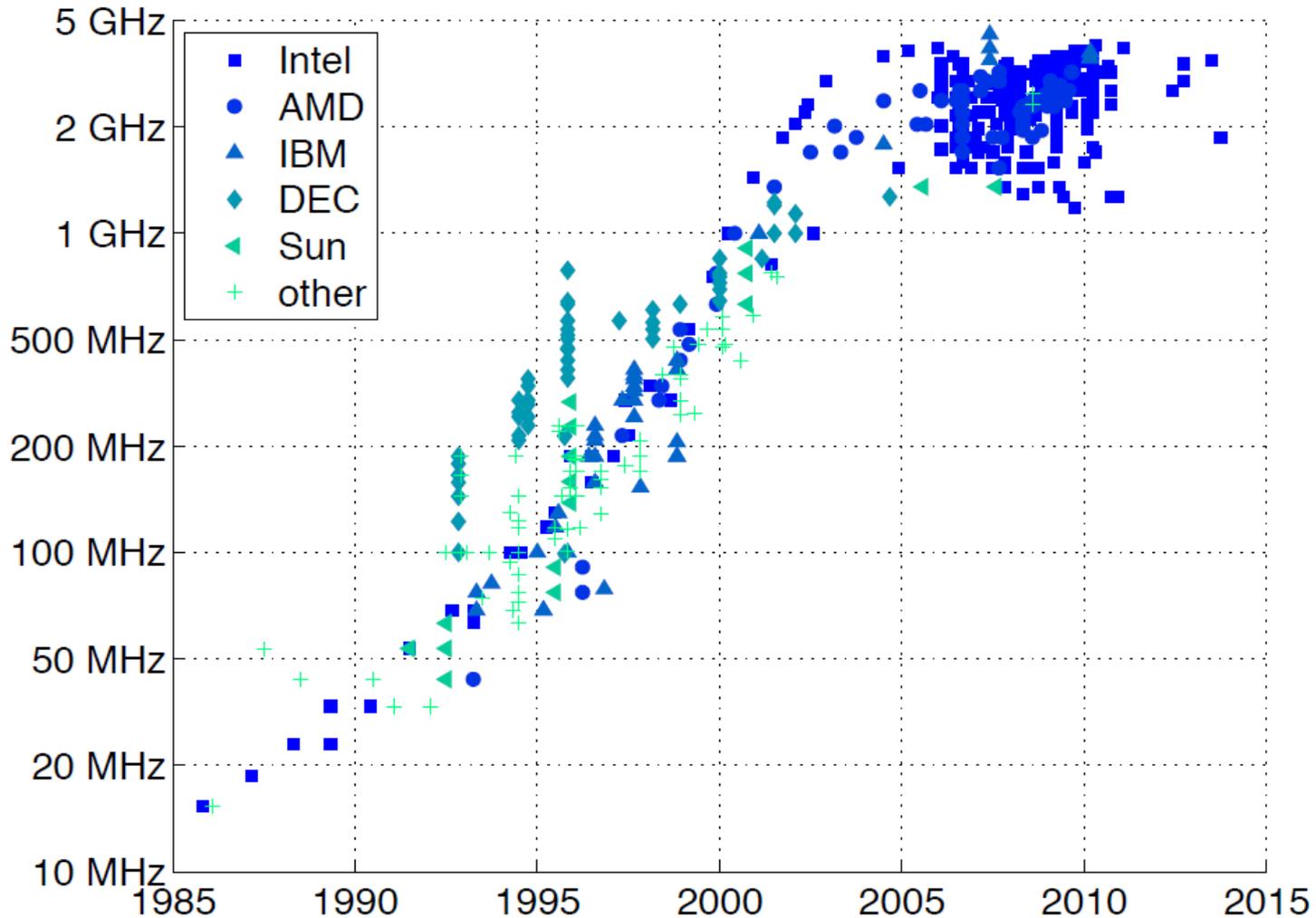
CRAY-1

Supporting Evidence



<http://cpudb.stanford.edu/>

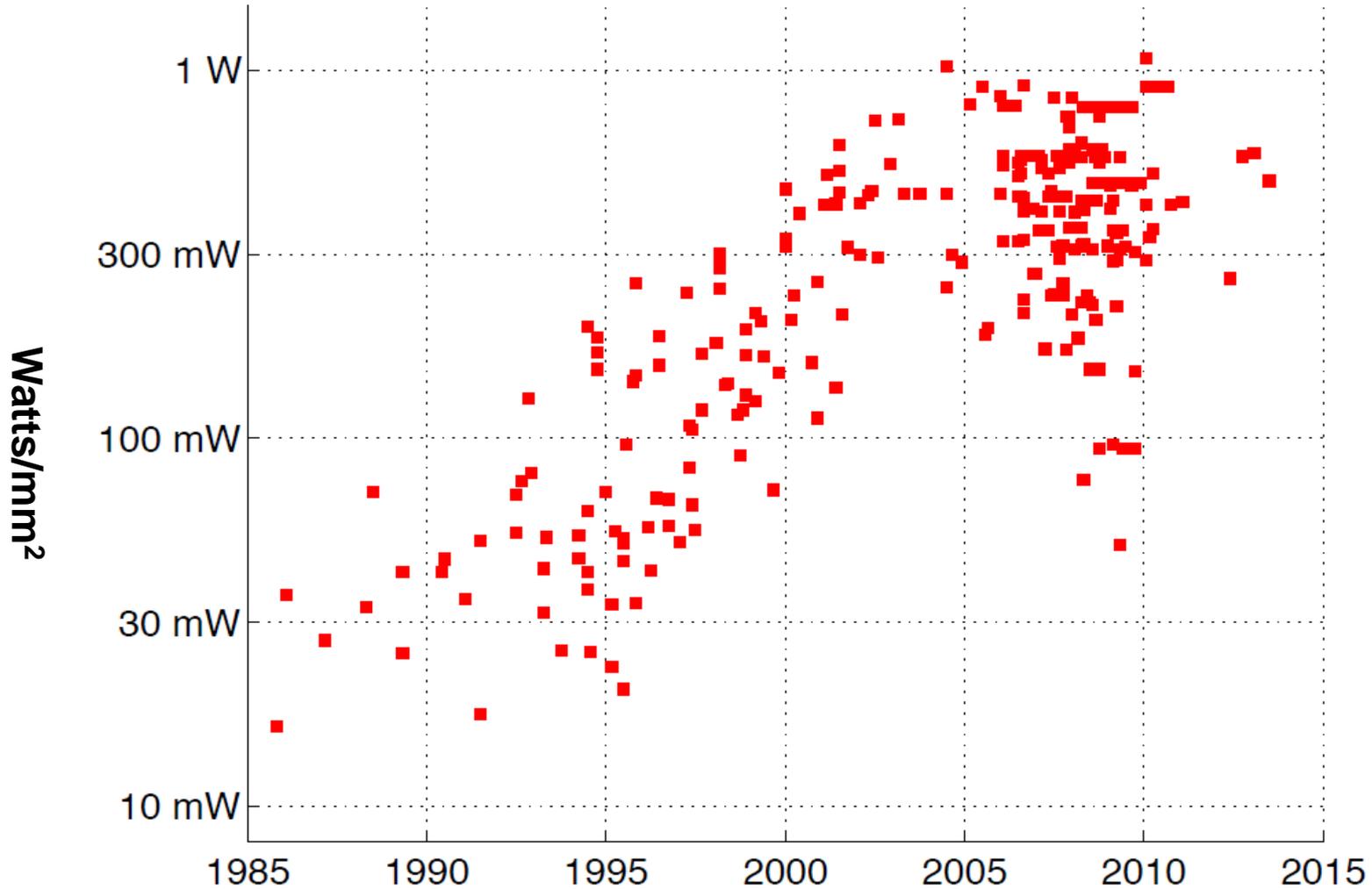
Houston, We Have A Problem



<http://cpudb.stanford.edu/>

1.1: Computing's Energy Problem: (and what we can do about it)

The Power Limit

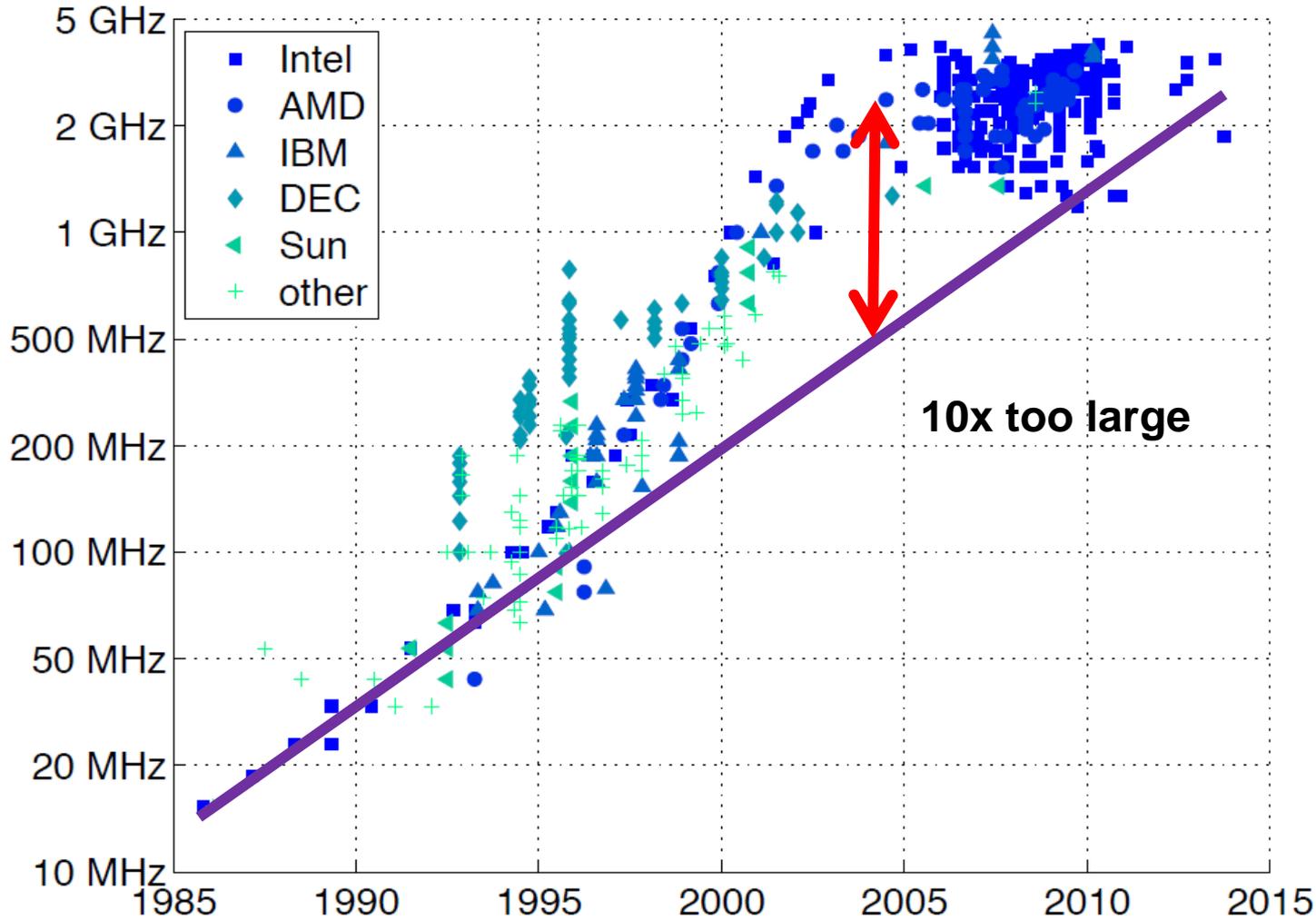


<http://cpudb.stanford.edu/>

1.1: Computing's Energy Problem: (and what we can do about it)

Clever

Power Increased Because We Were Greedy



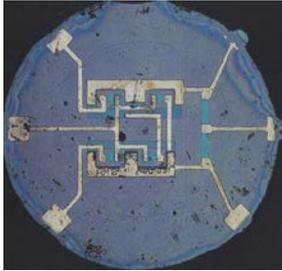
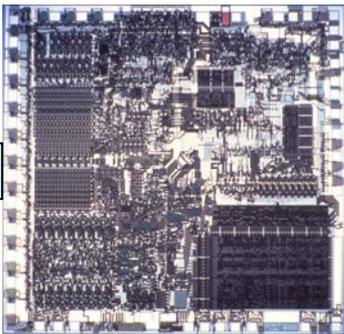
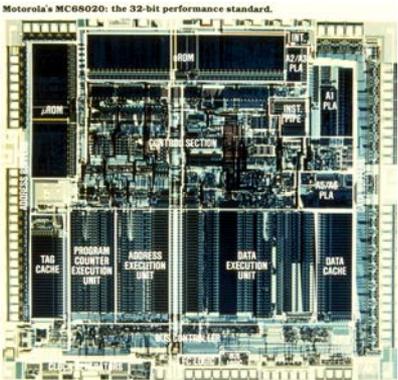
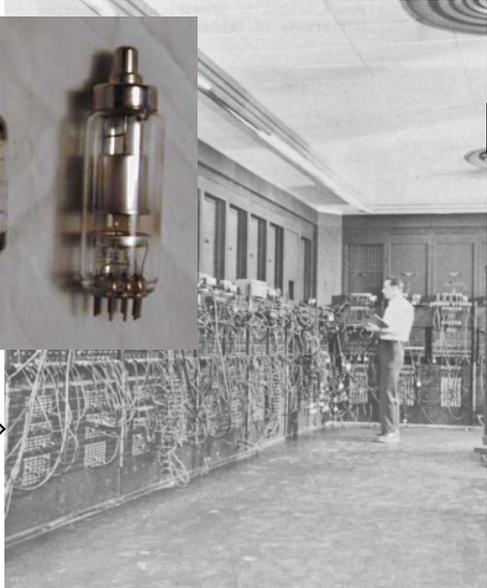
<http://cpudb.stanford.edu/>

1.1: Computing's Energy Problem: (and what we can do about it)

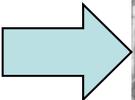
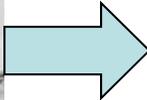
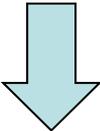
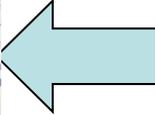
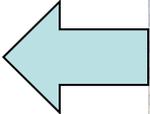
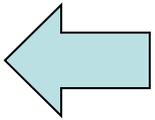
Think About It

$$P = \frac{\text{ENERGY}}{\text{OP}} \frac{\text{OPS}}{\text{S}}$$

Technology to the Rescue?



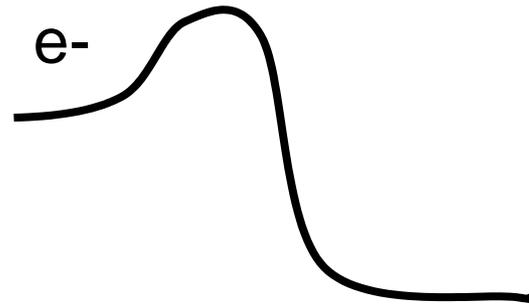
X



Problems w/ Replacing CMOS

Pretty fundamental physics

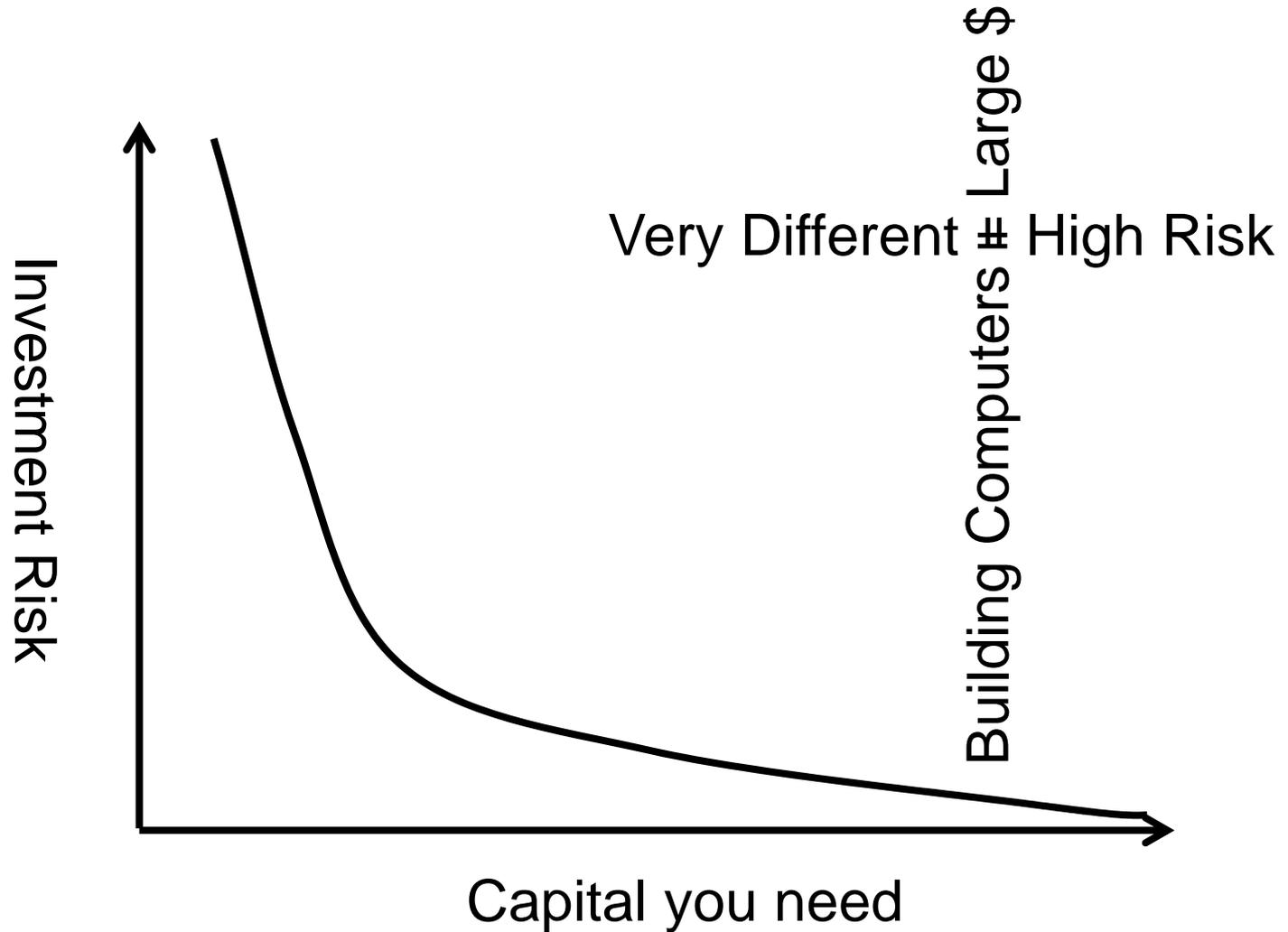
- Avoiding this problem will be hard



Its capability is pretty amazing

- fJ/gate, 10ps delays, 10^9 working devices

Catch - 22



The Truth About Innovation

Google™



ARM®

amazon.com®



ebay®

Start by creating new markets

Our CMOS Future

Will see tremendous innovative uses of computation

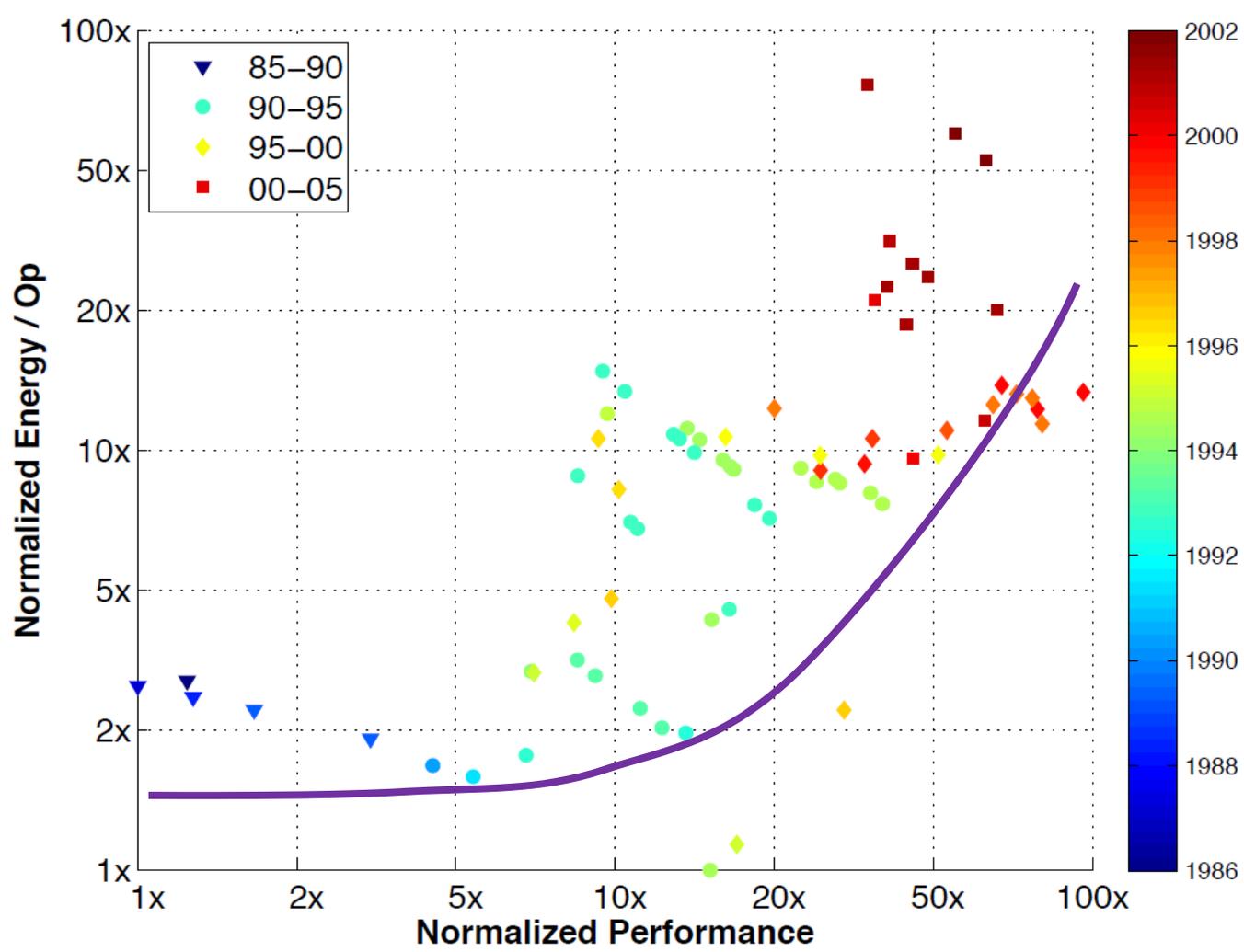
- Capability of today's technology is incredible
- Can add computing and communication for nearly \$0
- Key questions are what problems need to be solved?

Most performance system will be energy limited

- These systems will be optimized for energy efficiency

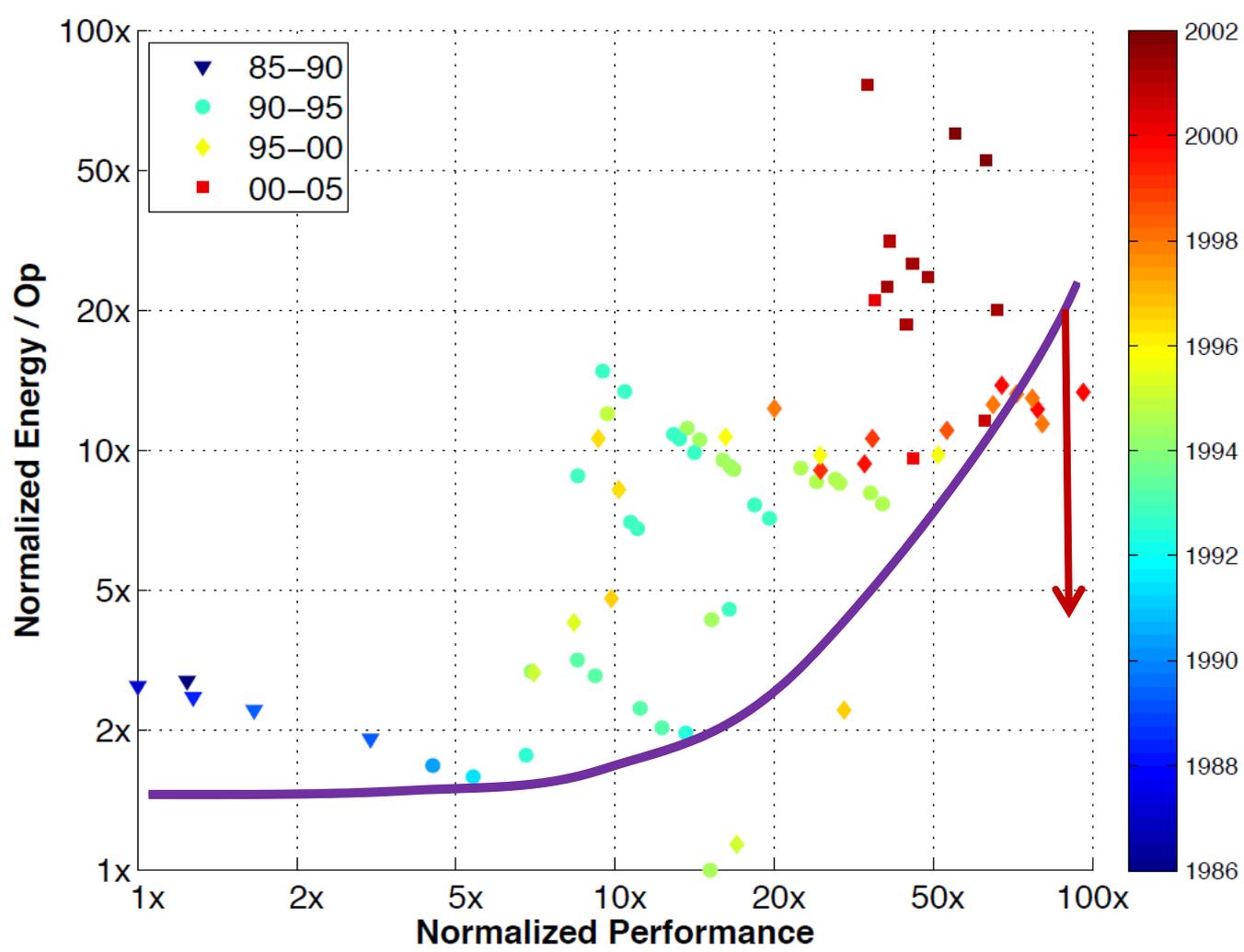
$$\text{Power} = \text{Energy/Op} * \text{Ops/sec}$$

Processor Energy – Delay Trade-off



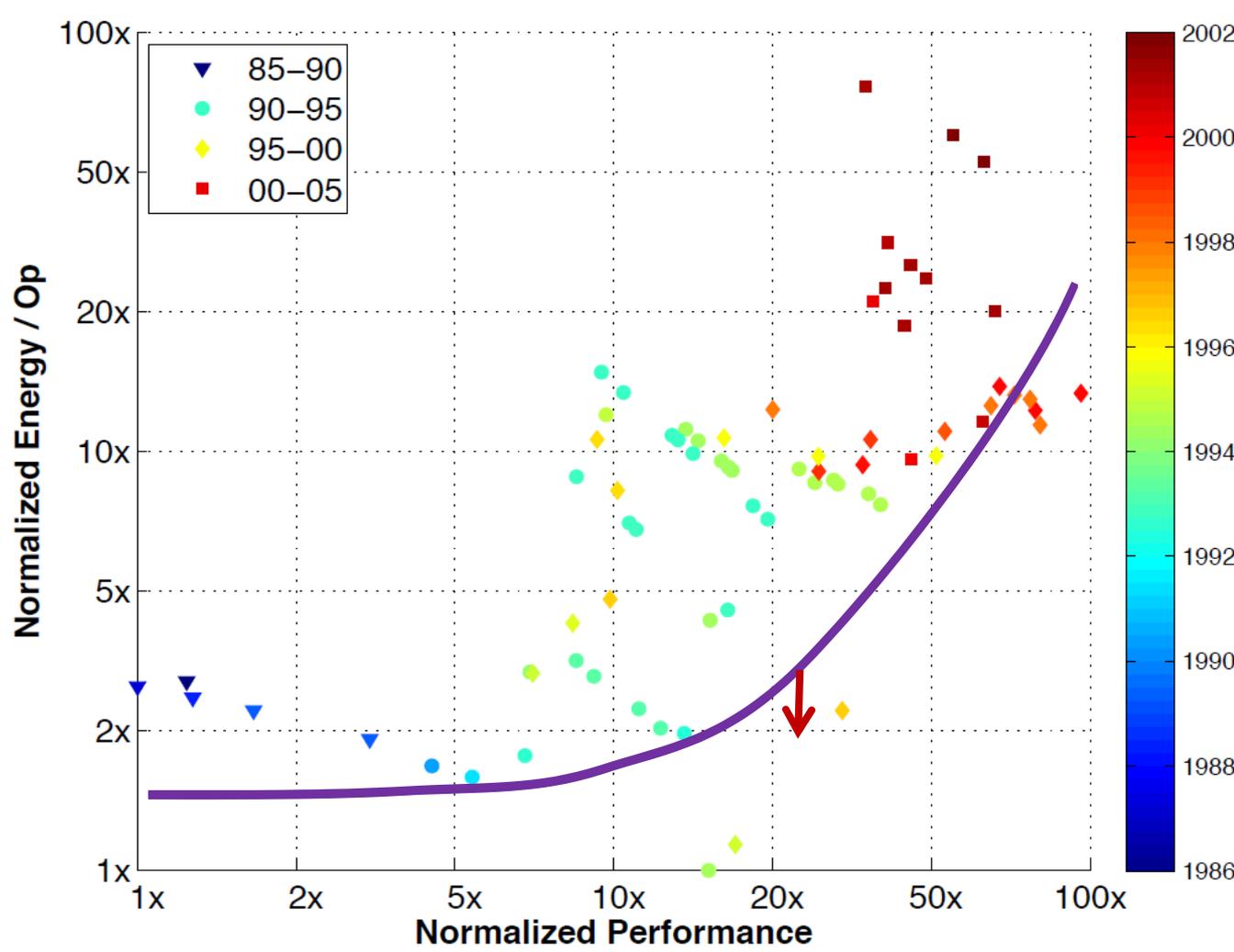
<http://cpudb.stanford.edu/>

The Rise of Multi-Core Processors



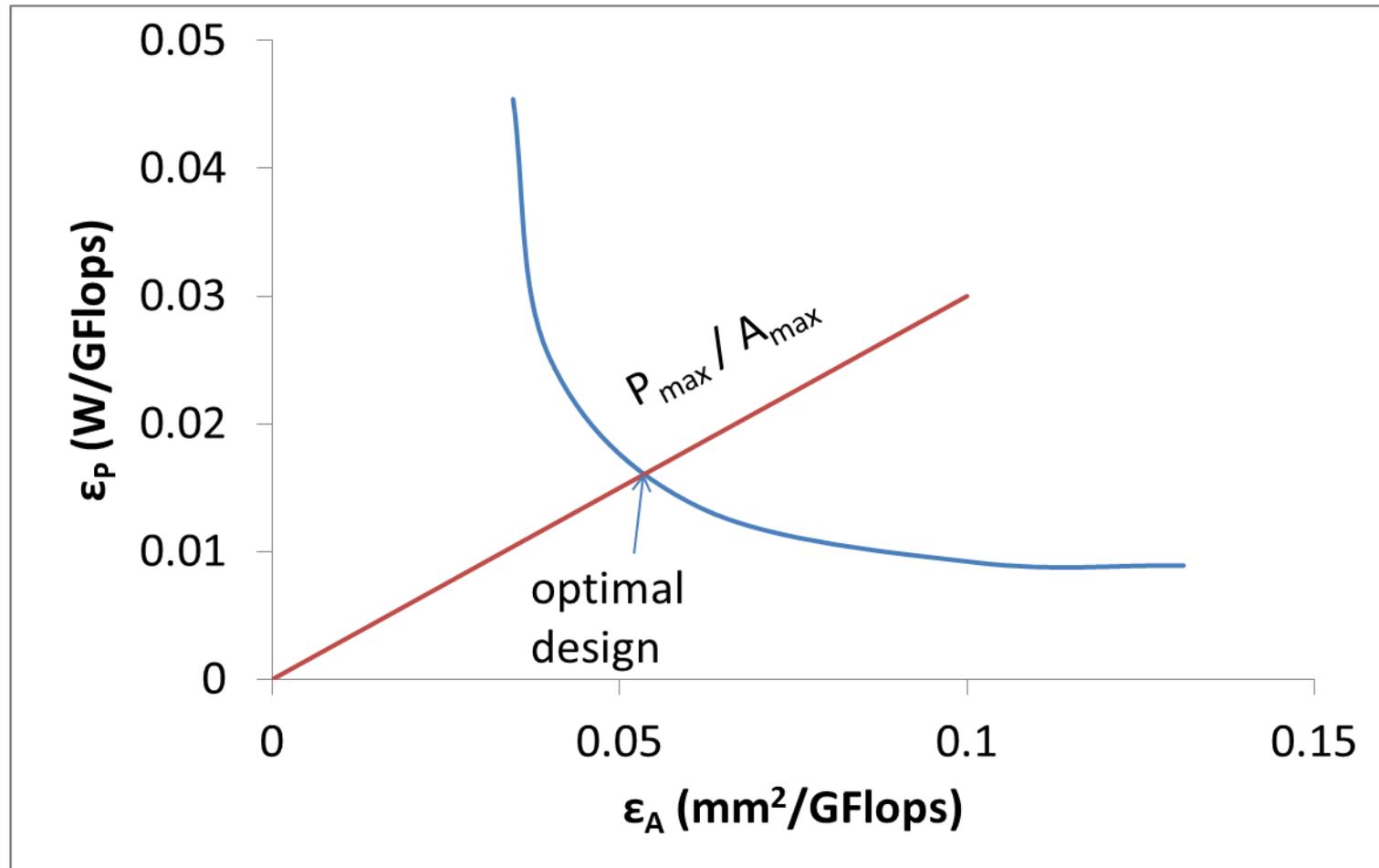
<http://cpudb.stanford.edu/>

The Stagnation of Multi-Core Processors



<http://cpudb.stanford.edu/>

Optimizing Parallel Machines (GPUs)

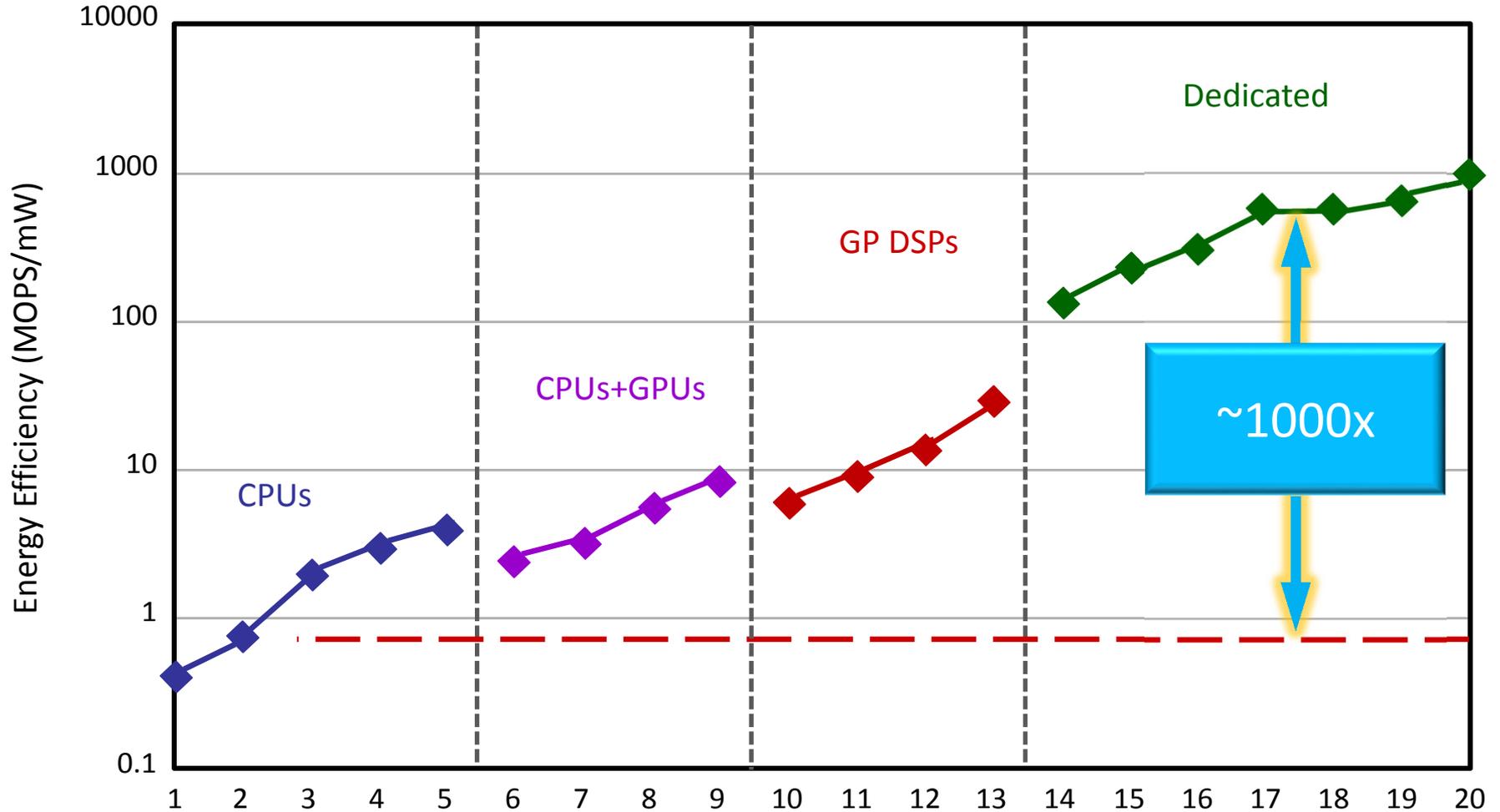


Galal et al, Trans on Computers, July 2011, pp 913,922

Have A Shiny Ball, Now What?



Signal Processing ASICs



Markovic, EE292 Class, Stanford, 2013

1.1: Computing's Energy Problem: (and what we can do about it)

The Push For Specialized Hardware

Dark Silicon and the End of Multicore

Hadi Esmaeilzadeh^{*} Emily Blem[‡] Renée St. Amant[†]
^{*}University of Washington [†]University of Texas at Austin
hadia@cs.washington.edu blem@cs.wisc.edu

ABSTRACT

Since 2005, processor designers have increasingly exploit Moore's Law scaling, rather than focusing on performance. The failure of Dennard scaling, to which multicore parts is partially a response, may soon limit multicore scaling limits by combining device scaling, scaling, and multicore scaling to measure the speedup potential for device scaling for the next five technology generations. For more conservative device scaling parameters, we use a set of parallel workloads to derive Pareto-optimal frontiers and power-to-derive Pareto-optimal frontiers for area/performance and power/performance. Finally, to model multicore scaling and power-bound core power. The multicore designs we study include single-threaded CPU-like and massively threaded GPU-like multicore chip organizations with symmetric, asymmetric, dynamic, and composed topologies. The study shows that regardless of chip organization and topology, multicore scaling is power limited to a degree not widely appreciated by the computing community. Even at 22 nm (just one year from now), 21% of a fixed-size chip must be powered off, and at 8 nm, this number grows to more than 50%. Through parallel workloads, leaving a nearly 24-fold gap from a target of doubled performance per generation.

Categories and Subject Descriptors: C.0 [Computer Systems Organization] General — Modeling of computer architecture; C.0 [Computer Systems Organization] General — System architectures

General Terms: Design, Measurement, Performance

Keywords: Dark Silicon, Modeling, Power, Technology Scaling, Multicore

Reducing the Energy of Mature Computations: Conservation Cores:

Ganesh Venkatesh
Viadyslav Bryksin
Jack Sampson
Jose Lugo-Martinez

Nathan Goulding
Steven Swanson

Saturmino Garcia
Michael Bedford Taylor

Department of Computer Science & Engineering
University of California, San Diego
[gvenkatesh, jsampson, ngoulding, sat, vbryksin, jlugomar, swanson, mbtaylor]@cs.ucsd.edu

Abstract

Growing transistor counts, limited power budgets, and the breakdown of voltage scaling are currently conspiring to create a utilization wall that limits the fraction of a chip that can run at full speed at one time. In this regime, specialized, energy-efficient processors can increase parallelism by reducing the per-computation power requirements and allowing more computations to execute under the same power budget. To pursue this goal, this paper introduces conservation cores, which focus on reducing energy and energy-delay instead of increasing performance. This focus on energy makes c-cores an excellent match for many applications (e.g., irregular integer codes). We present a toolchain for automatically synthesizing c-cores from application source code and demonstrate that they can significantly reduce energy and energy-delay for a wide range of applications. The c-cores are not patching, a form of targeted reconfigurability, that allows them to be used in new versions of the software they target. Our approach to conservation cores can reduce energy consumption by up to 2.1x for whole applications and by up to 2.1x for individual processors.

power. Consequently, the rate at which we can switch transistors is far outpacing our ability to dissipate the heat created by those transistors.

The result is a technology-imposed utilization wall that limits the fraction of the chip we can use at full speed at one time. Our experiments with a 45 nm TSMC process show that we can switch less than 7% of a 300mm² die at full frequency within an 80W power budget. ITRS roadmap projections and CMOS scaling theory suggests that this percentage will decrease to less than 3.5% in 32 nm, and will continue to decrease by almost half with each process generation—and even further with 3-D integration.

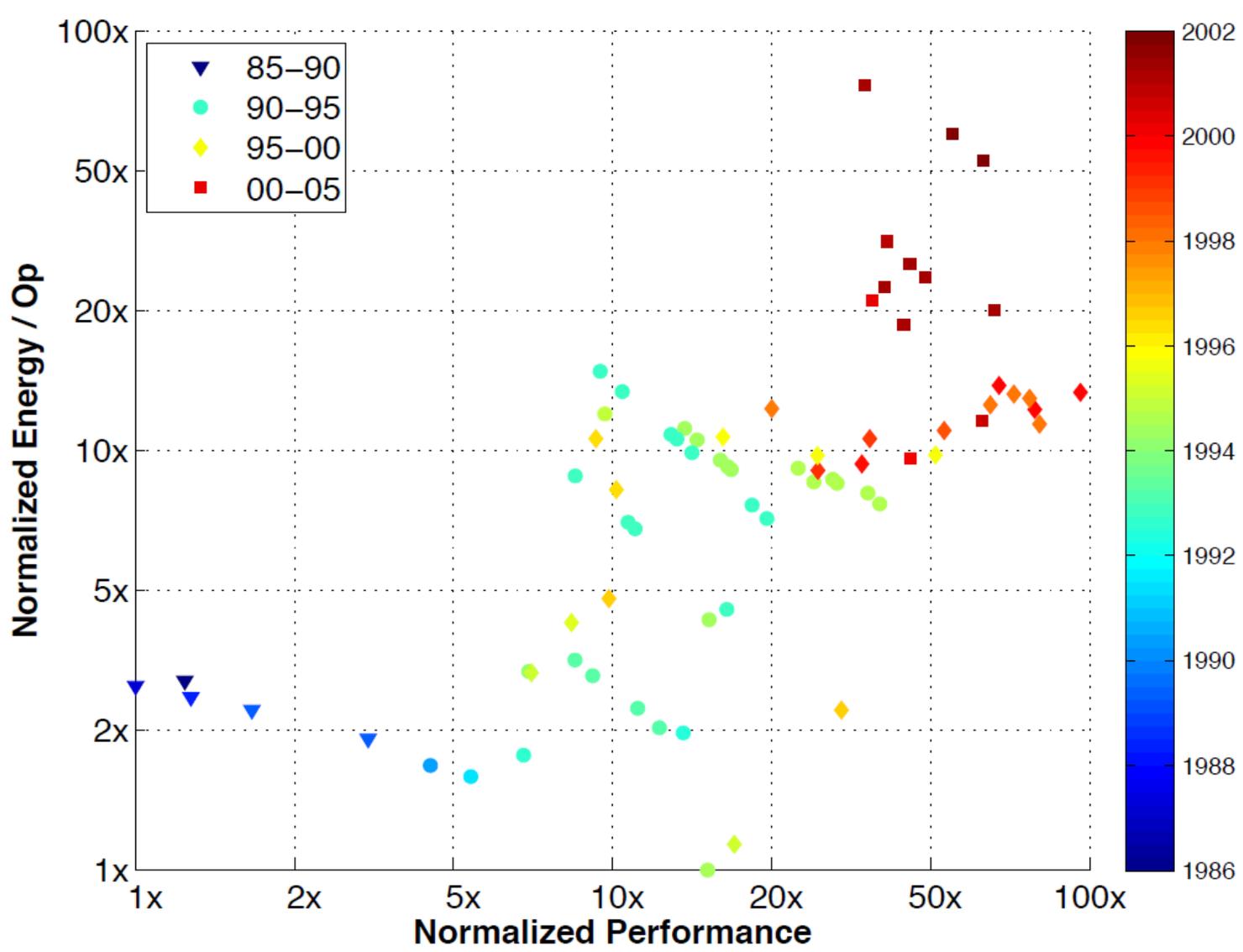
The effects of the utilization wall are already indirectly apparent in modern processors: Intel's Nehalem provides a "turbo mode" that powers off some cores in order to run others at higher speeds. Another strong indication is that even though native transistor switching speeds have continued to double every two process generations, processor frequencies have not increased substantially in this regime, reducing per-operation energy [19] translates directly into increased potential parallelism for the system. If given computation can be made to consume less power at the same level of performance, other computations can be run in parallel without violating the power budget.

This paper attacks the utilization wall with conservation cores, or c-cores, an application-specific hardware circuitry created for the purpose of reducing energy consumption to run the entire chip at full frequency at once, it makes possible for the application at hand. In effect, it allows architects to trade area for energy in order to minimize the portions of the application area for energy. The utilization wall has made this trade-off possible for architects to trade area for energy in order to minimize the portions of the application area for energy. The utilization wall has made this trade-off possible for architects to trade area for energy in order to minimize the portions of the application area for energy.

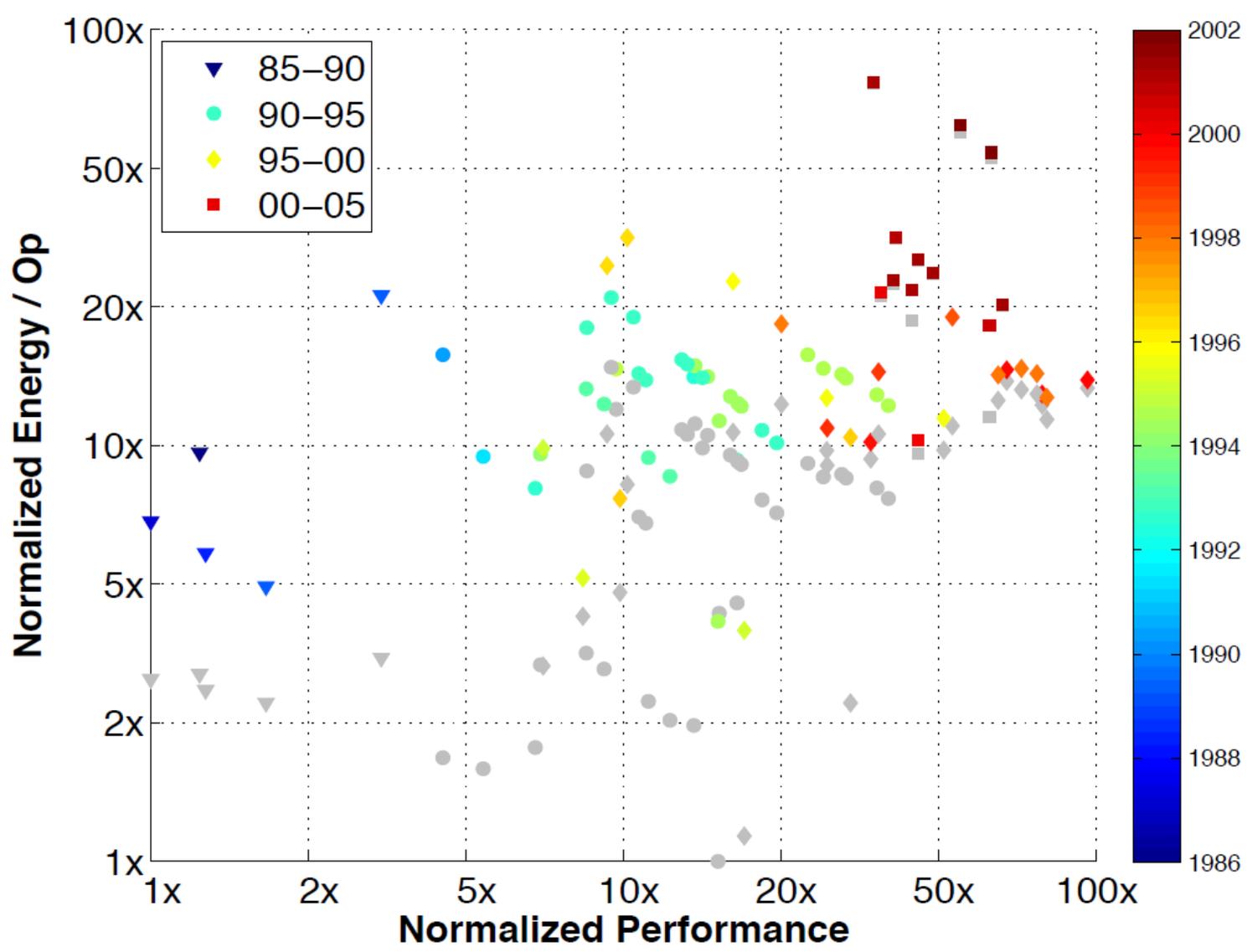
Before Talking About Specialization

**WE SHOULD CHECK
ONE MORE THING FIRST**

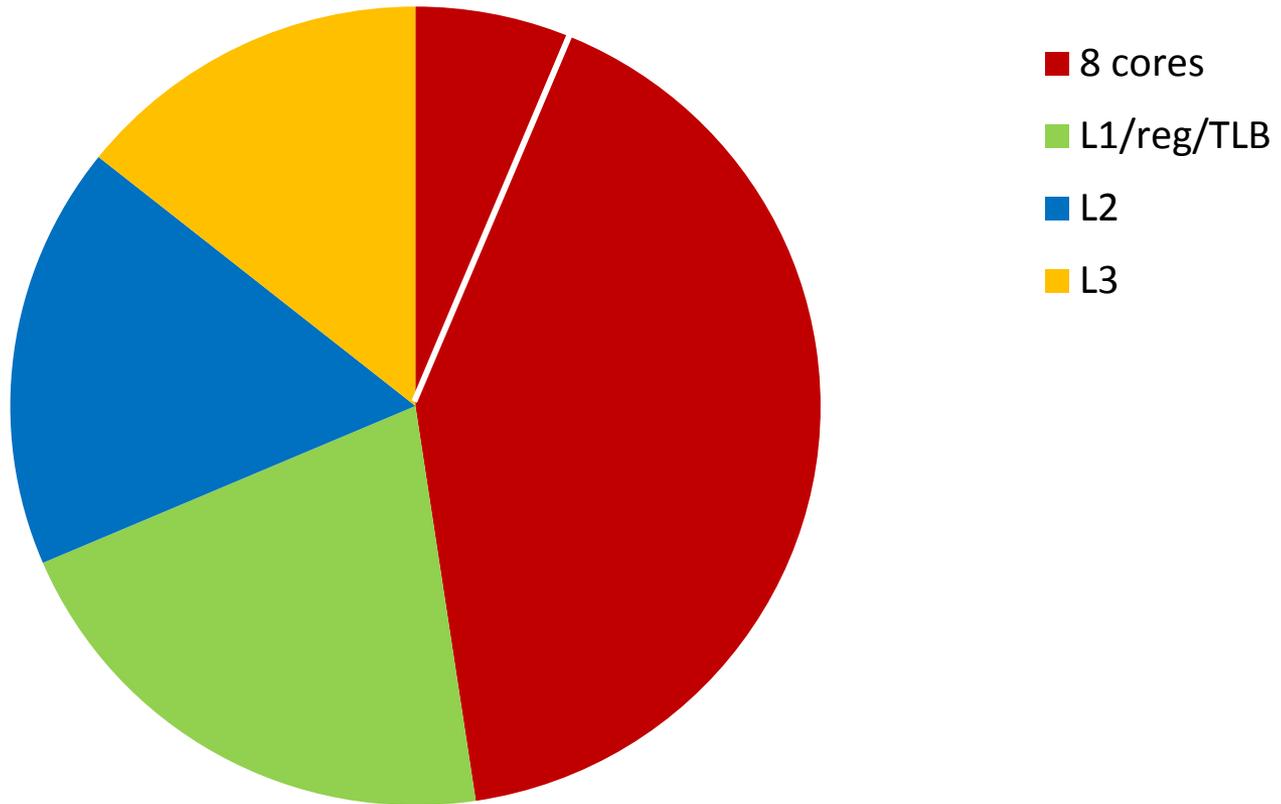
Don't Forget Memory System Energy



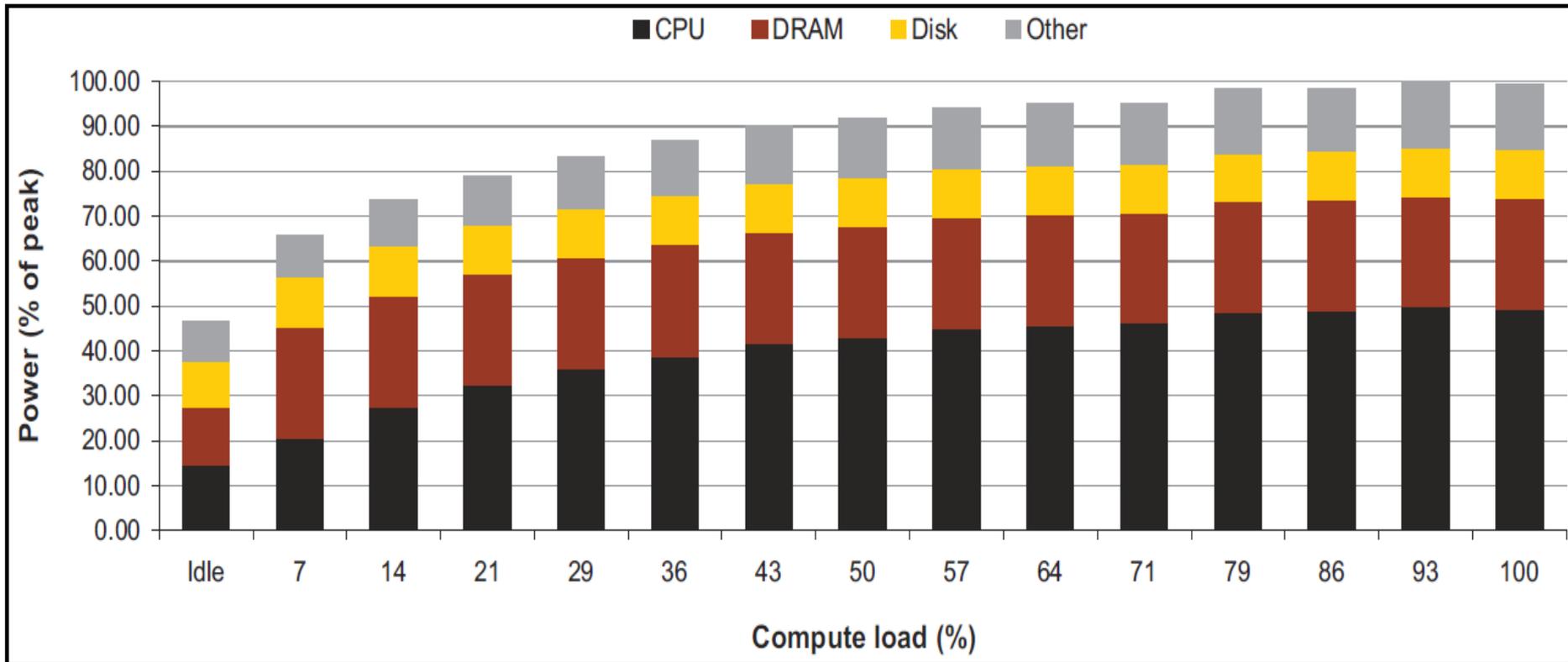
Processor Energy w/ Corrected Cache Sizes



Processor Energy Breakdown



Data Center Energy Specs

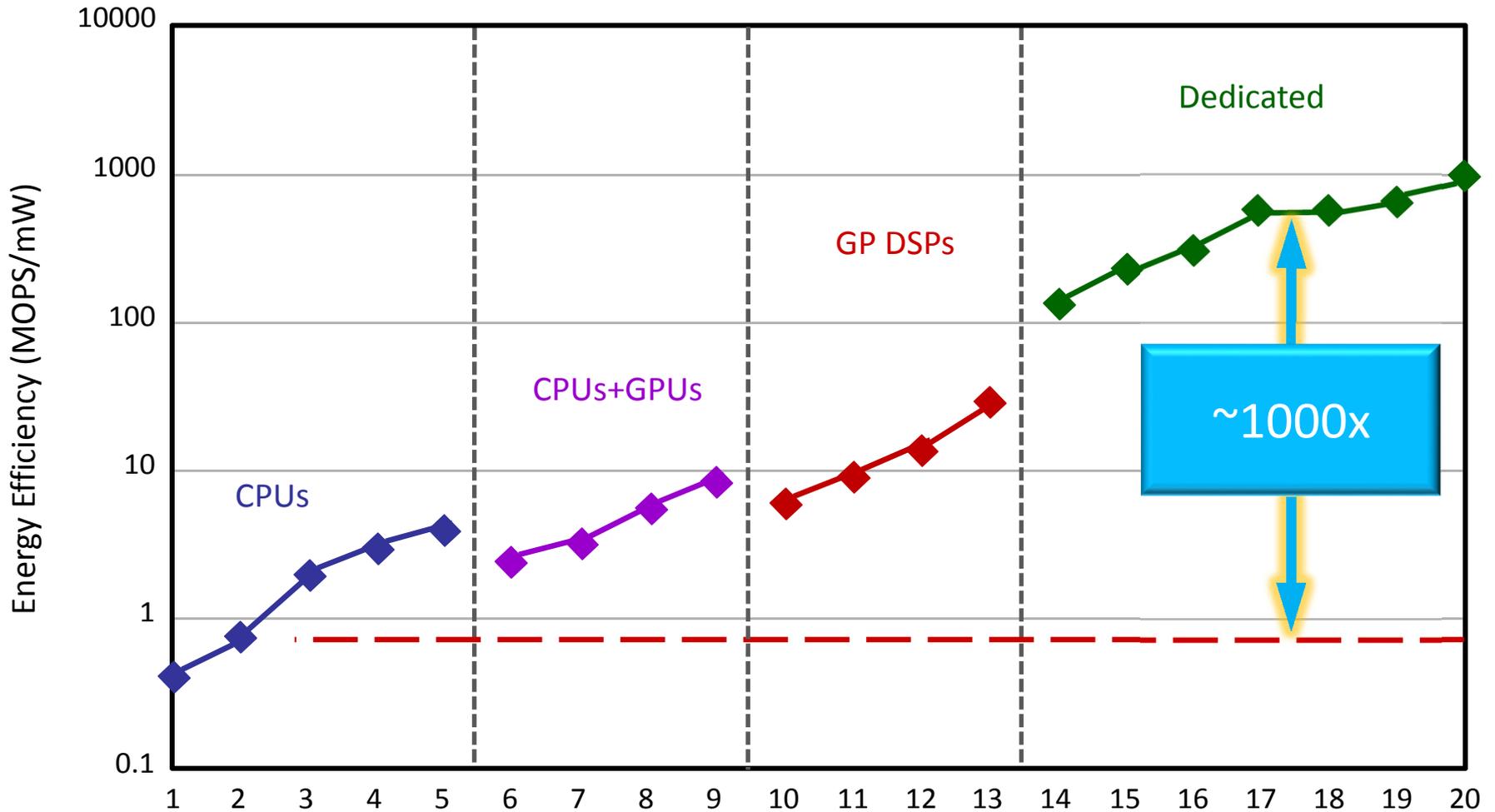


Malladi, ISCA, 2012

1.1: Computing's Energy Problem: (and what we can do about it)

SO HOW WILL ACCELERATORS HELP?

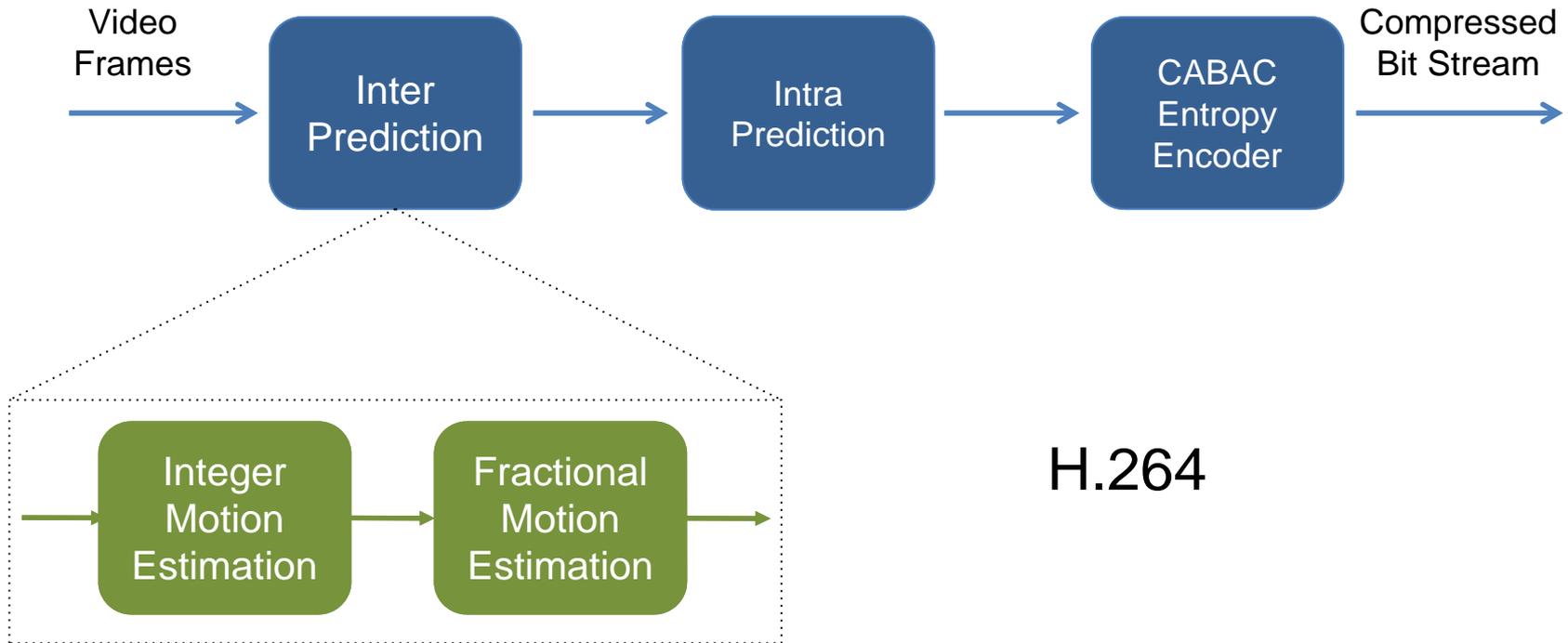
What Is Going On Here?



ASIC's Dirty Little Secret

All the ASIC applications have absurd locality

- And work on short integer data



90% of Execution time is here

Hamid et al, ISCA, 2010

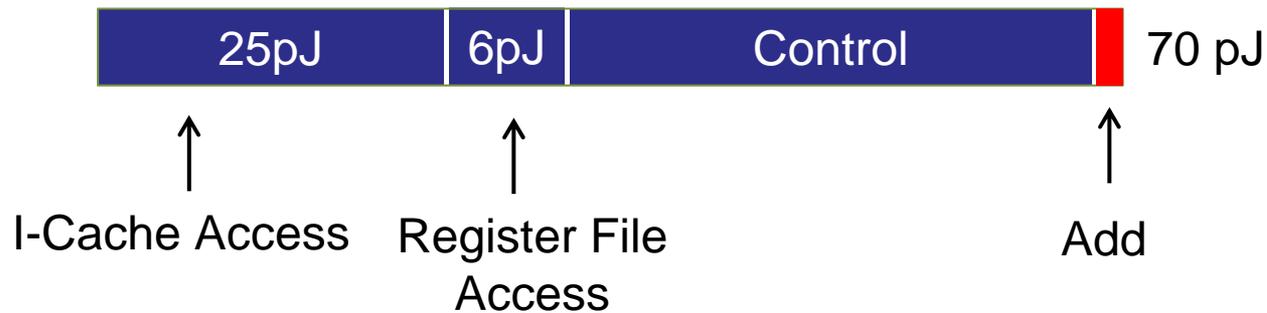
Rough Energy Numbers (45nm)

Integer	
Add	
8 bit	0.03pJ
32 bit	0.1pJ
Mult	
8 bit	0.2pJ
32 bit	3 pJ

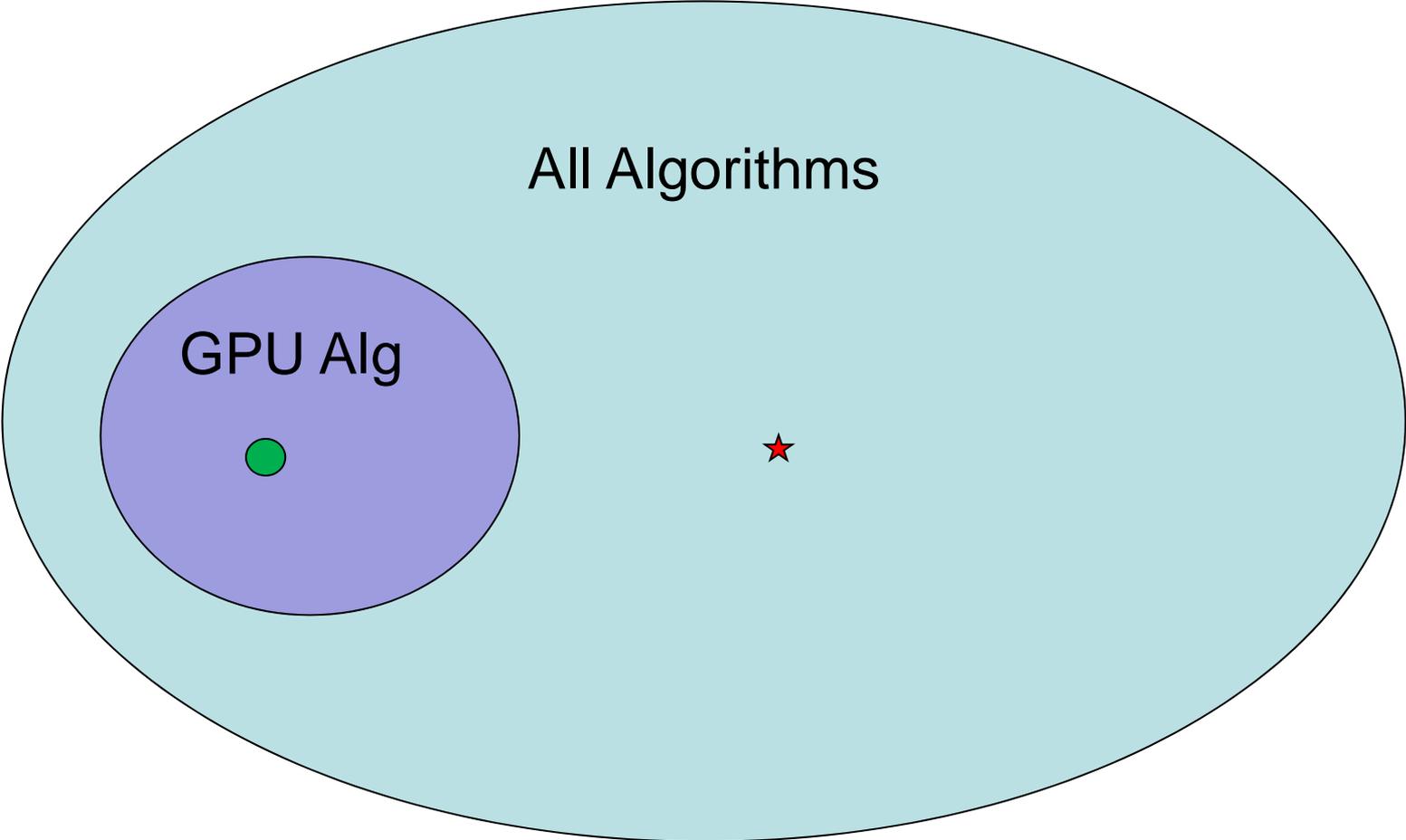
FP	
FAdd	
16 bit	0.4pJ
32 bit	0.9pJ
FMult	
16 bit	1pJ
32 bit	4pJ

Memory	
Cache	(64bit)
8KB	10pJ
32KB	20pJ
1MB	100pJ
DRAM	1.3-2.6nJ

Instruction Energy Breakdown

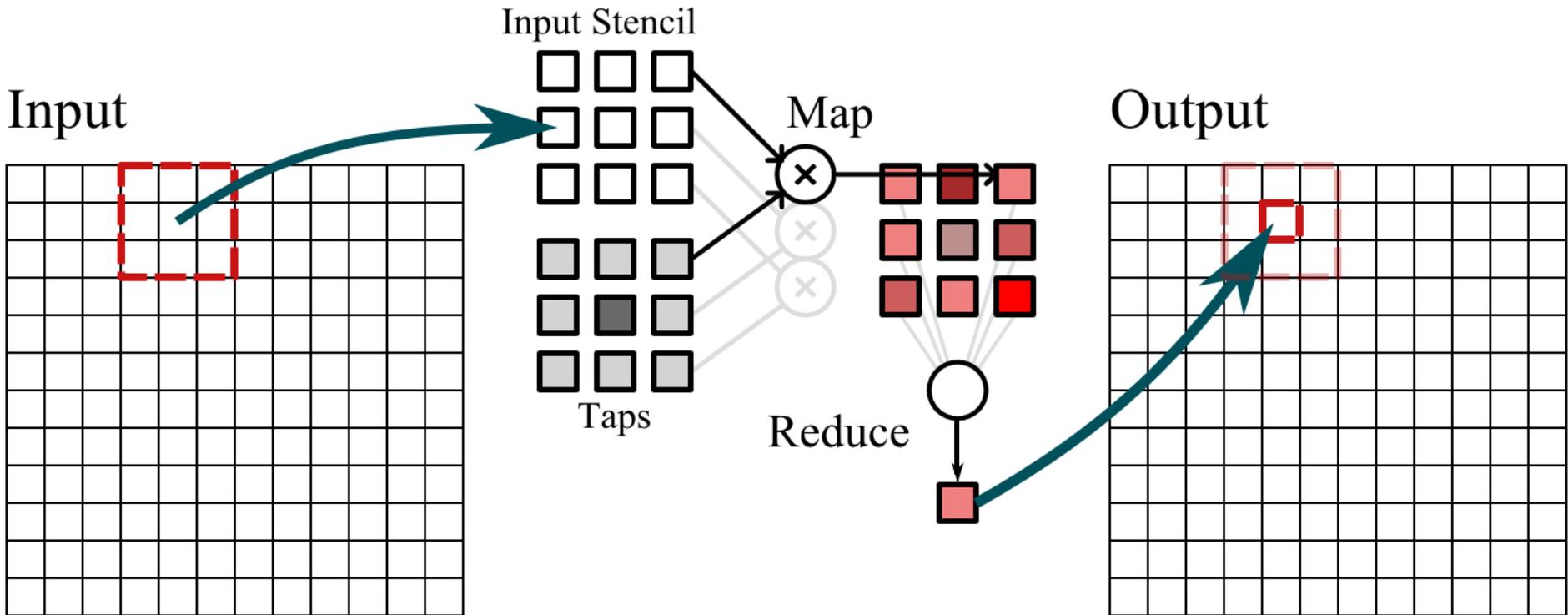


The Truth: It's More About the Algorithm than the Hardware

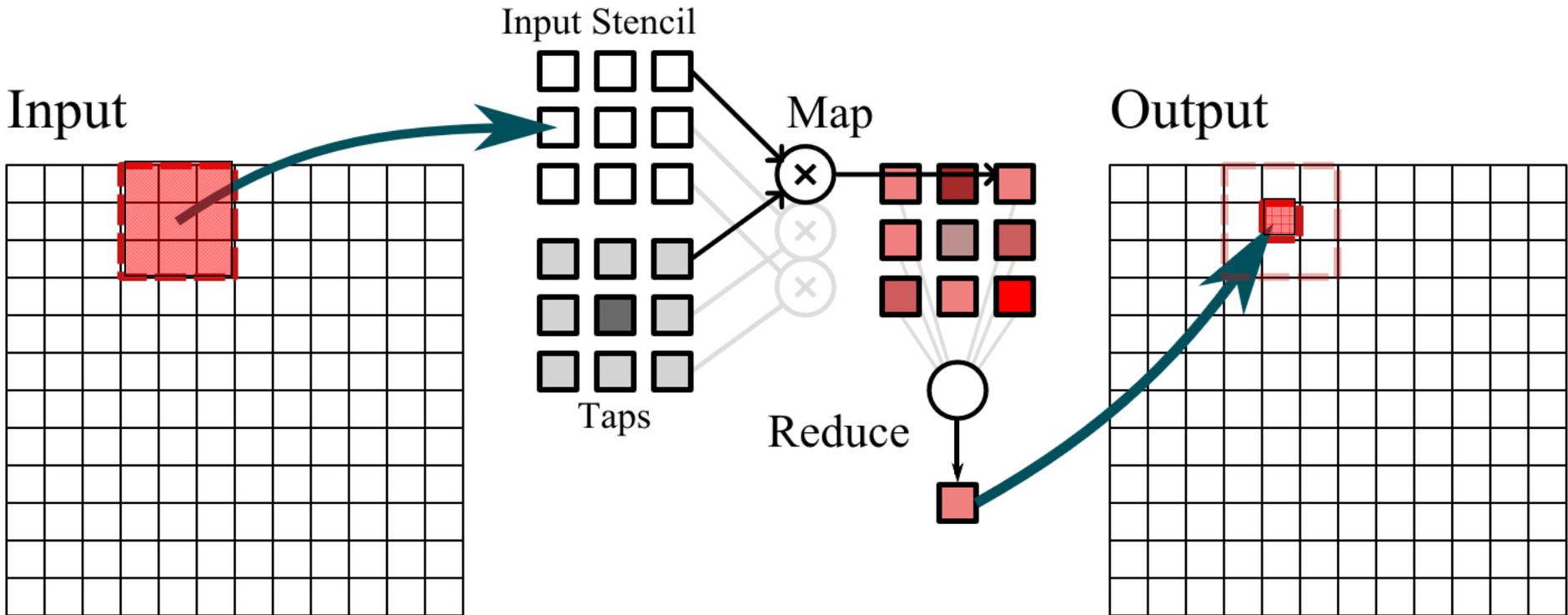




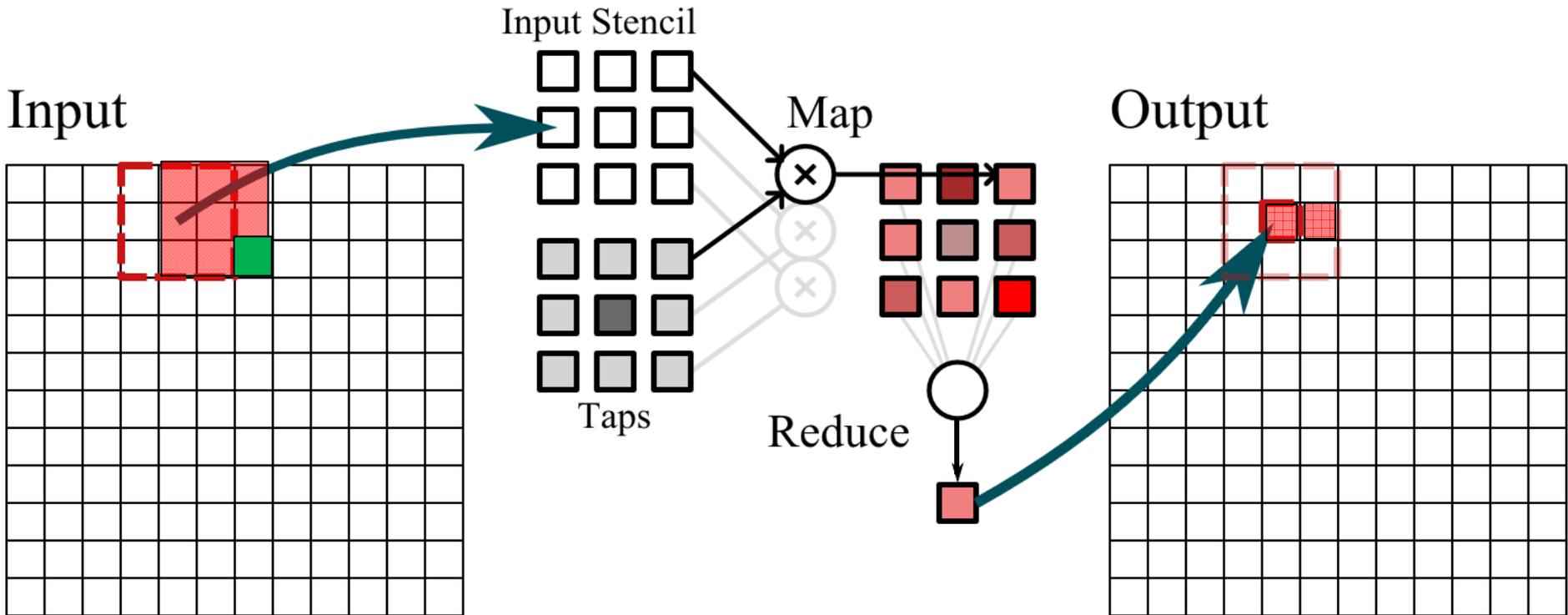
Highly Local Computation Model



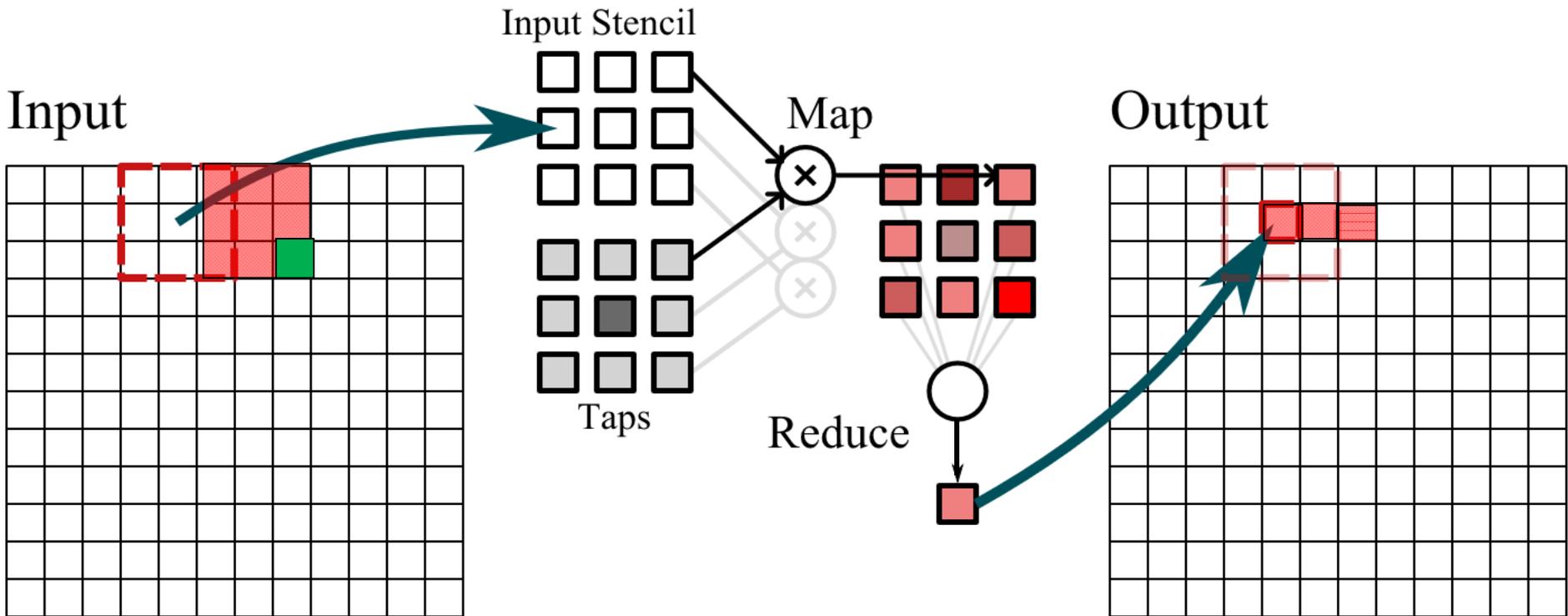
Highly Local Computation Model



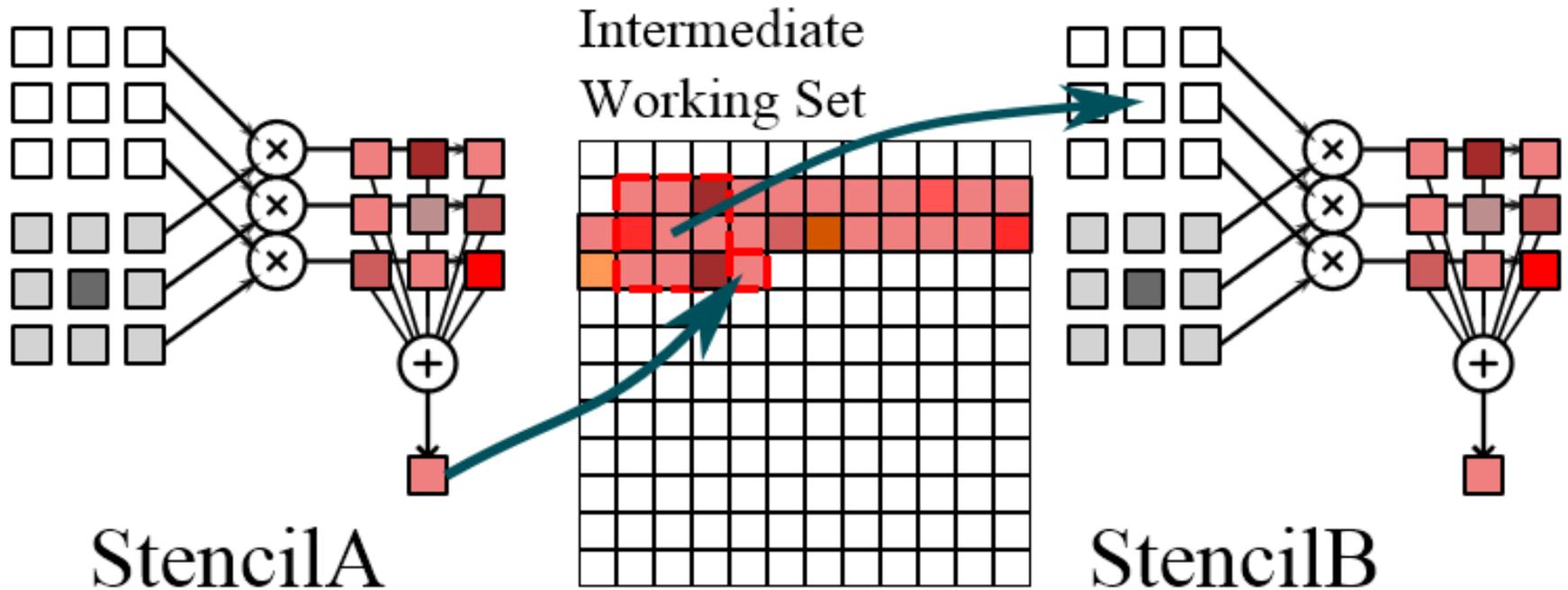
Highly Local Computation Model



Highly Local Computation Model



Compose These Cores into a Pipeline



Program in space, not time

- Makes building programmable hardware more difficult

Working on System to Explore This Space

Takes high-level program

- Graph of stencil kernels

Maps to hardware level assembly

- Compute graph of operations for each kernel

Currently we map the result to:

- FPGA, custom ASIC

Enabling Innovation

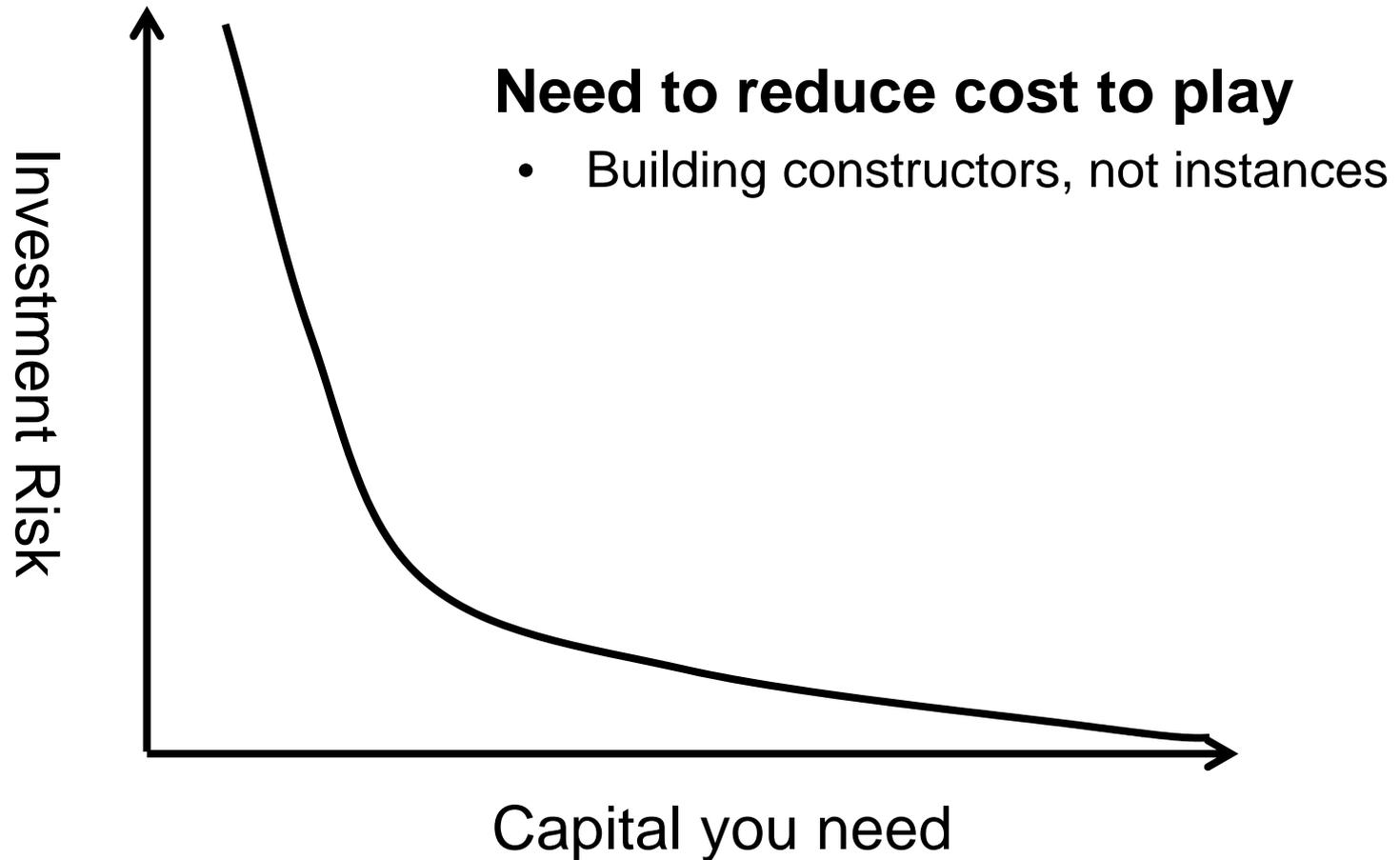
You don't just compile applications to efficiency

- Need to tweak the application to fit constraints

Need to enable application experts to play

- They know how to “cheat” and still get good results

Remember This Trade-off?



Not All Systems Are On The Bleeding Edge



App Store For Hardware



There's almost no limit
to what iPhone can do.

The App Store has the best selection of mobile apps — from Apple and third-party developers. And they're all designed specifically for iPhone. The more apps you download, the more you'll realize your iPhone can do just about anything you can imagine.



1.1: Computing's Energy Problem: (and what we can do about it)

Challenge



What Arduino can do

Arduino can sense the environment by receiving input from a variety of sensors and can affect its surroundings by controlling lights, motors, and other actuators. The microcontroller on the board is programmed using the [Arduino programming language](#) (based on [Wiring](#)) and the Arduino development environment

Community

The community of Arduino enthusiasts is vast, and includes region specific groups and special interest groups. The community is an excellent further source of assistance on all topics such as accessory selection, project assistance, and ideas of all sorts.

A New Hope

If technology is scaling more slowly

- We can incorporate current design knowledge into tools
- To create extensible system constructors

If killer products are going to be application driven

- Application experts need to design them

We can leverage the 1st bullet to enable the 2nd

- To usher in a new wave of innovative computing products