# 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



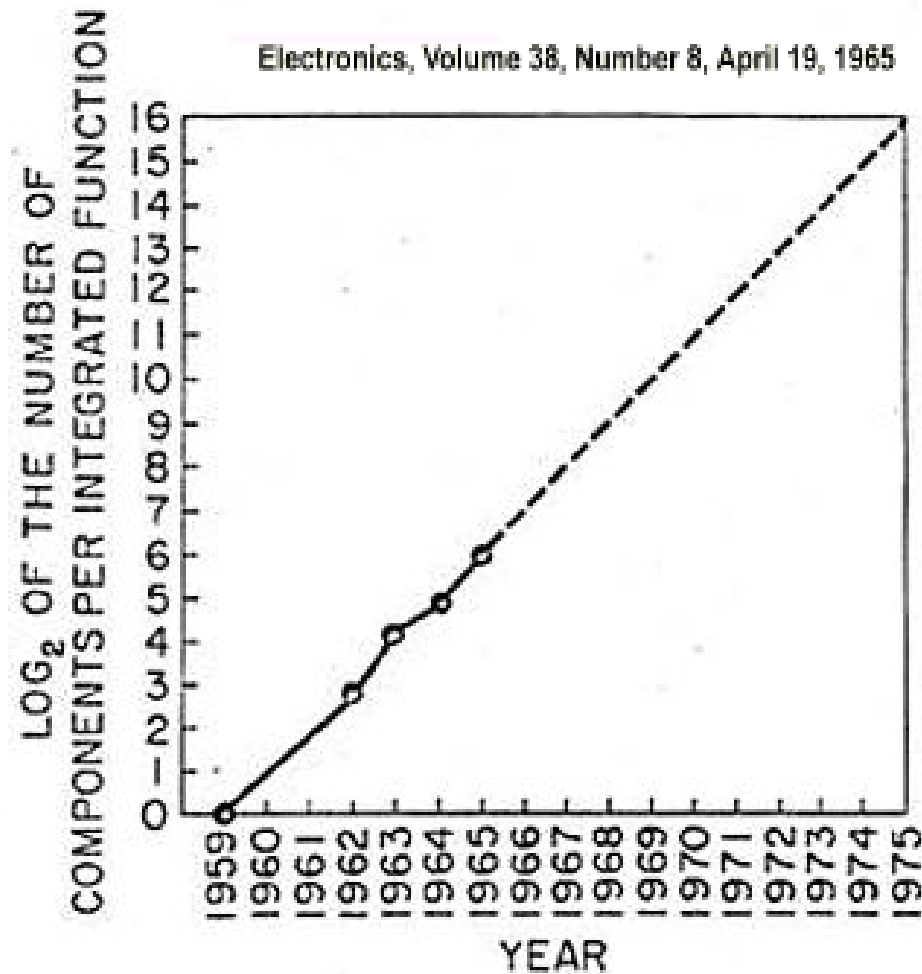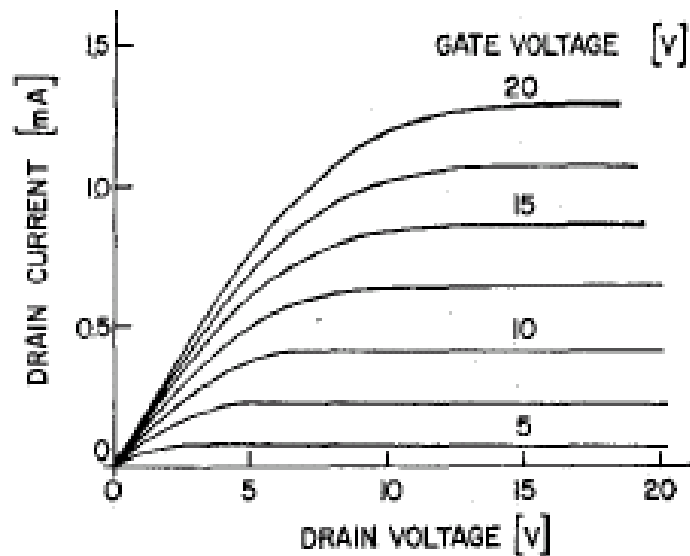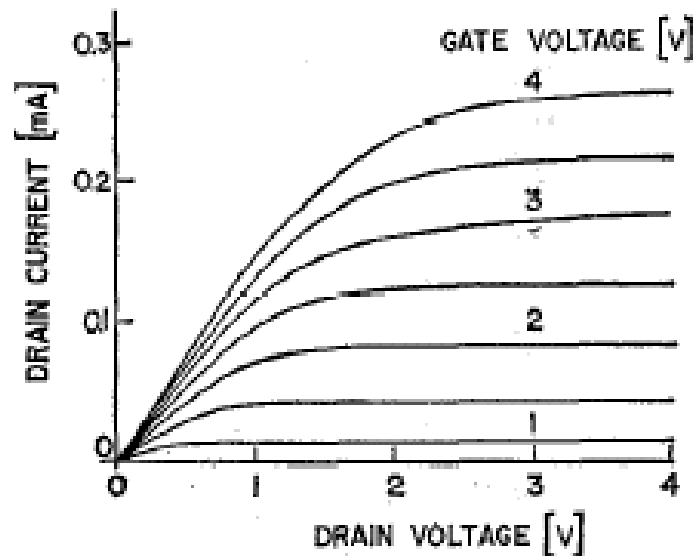Electronics, Volume 38, Number 8, April 19, 1965

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

# Dennard's Scaling
# Made Them Fast & Low Energy



## The triple play:

| | | |
|---|---|---|
| • Get more gates, | $1/L^2$ | $1/\alpha^2$ |
| • Gates get faster, | $CV/i$ | $\alpha$ |
| • Energy per switch | $CV^2$ | $\alpha^3$ |

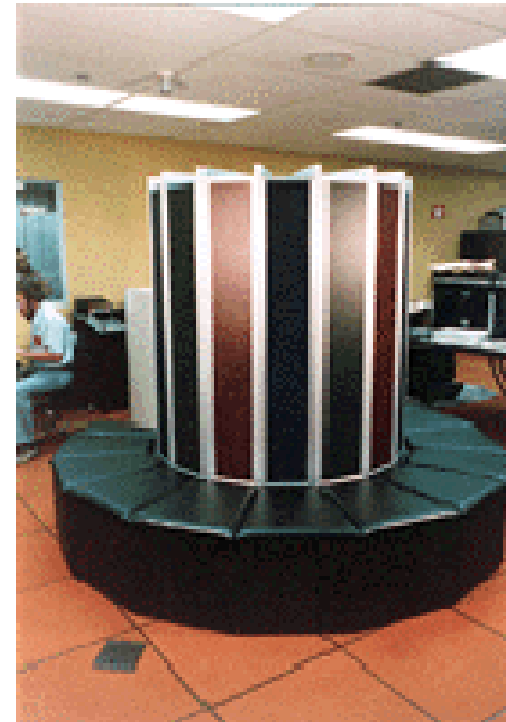Dennard, JSSC, pp. 256-268, Oct. 1974

# Our Expectation

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

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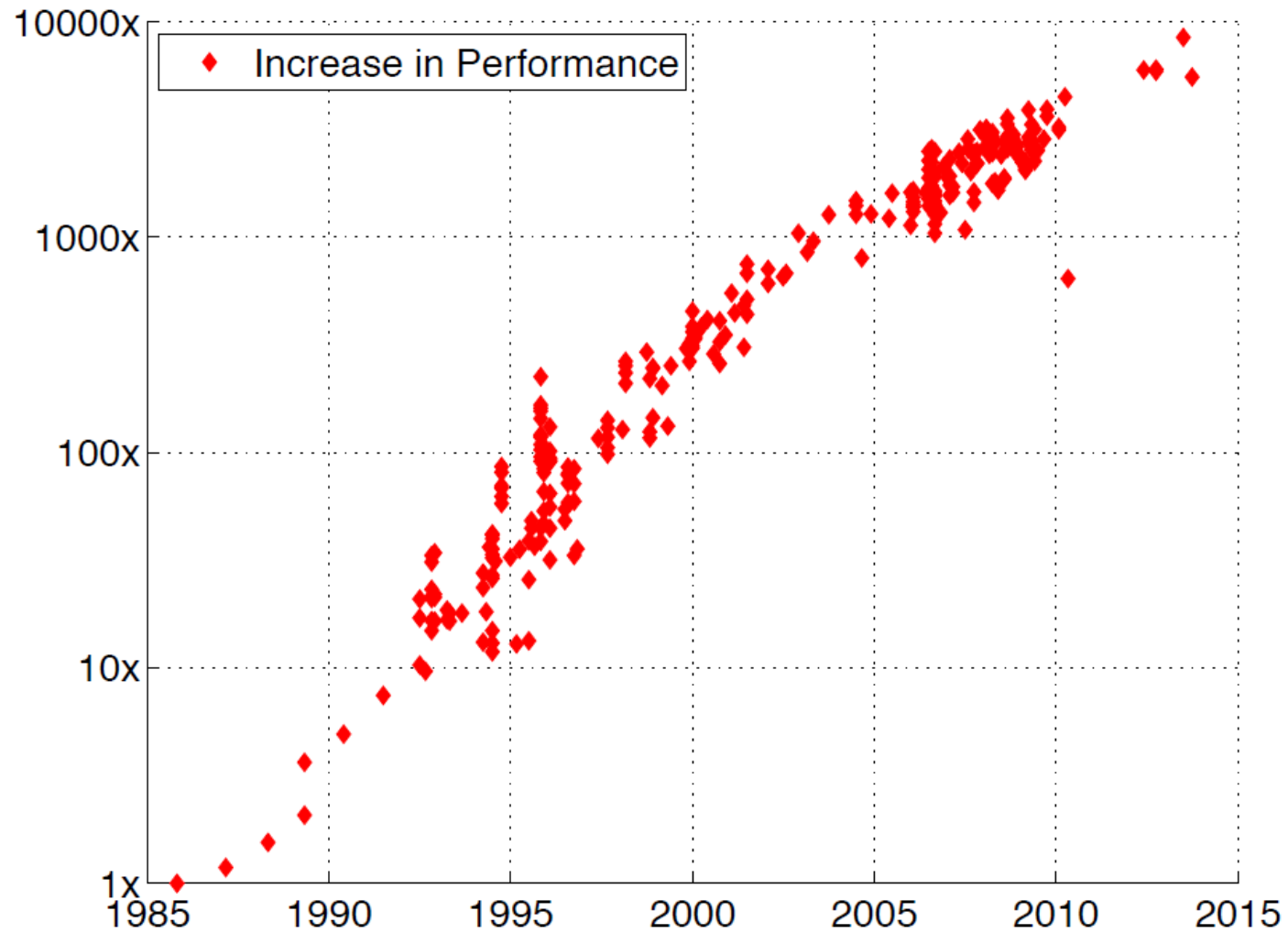- ~1 million gates (4/5 input NAND)
- 80MHz clock
- 115kW

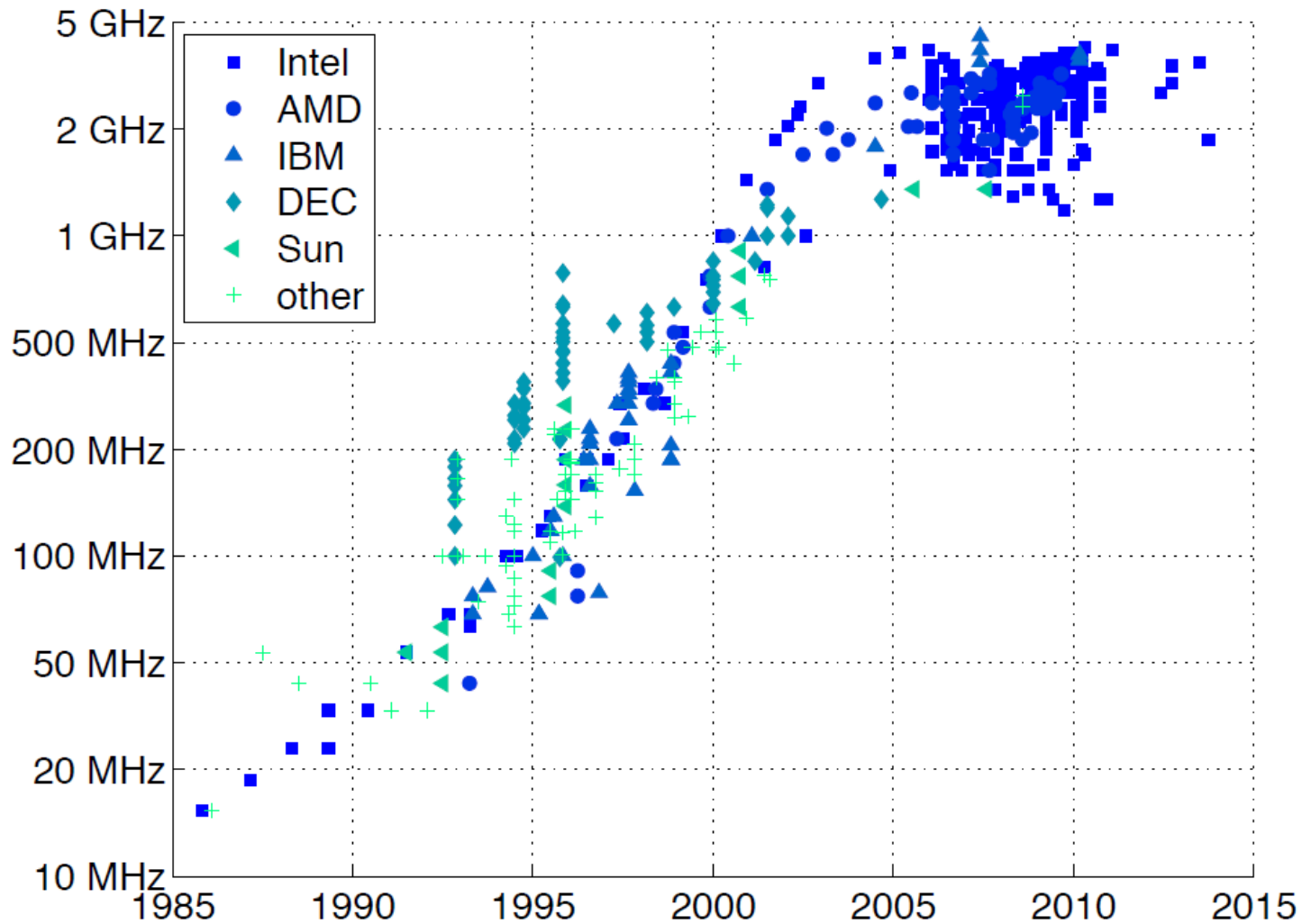## In 45nm (30 years later)

- < 3 mm$^2$
- > 1 GHz
- ~ 1 W



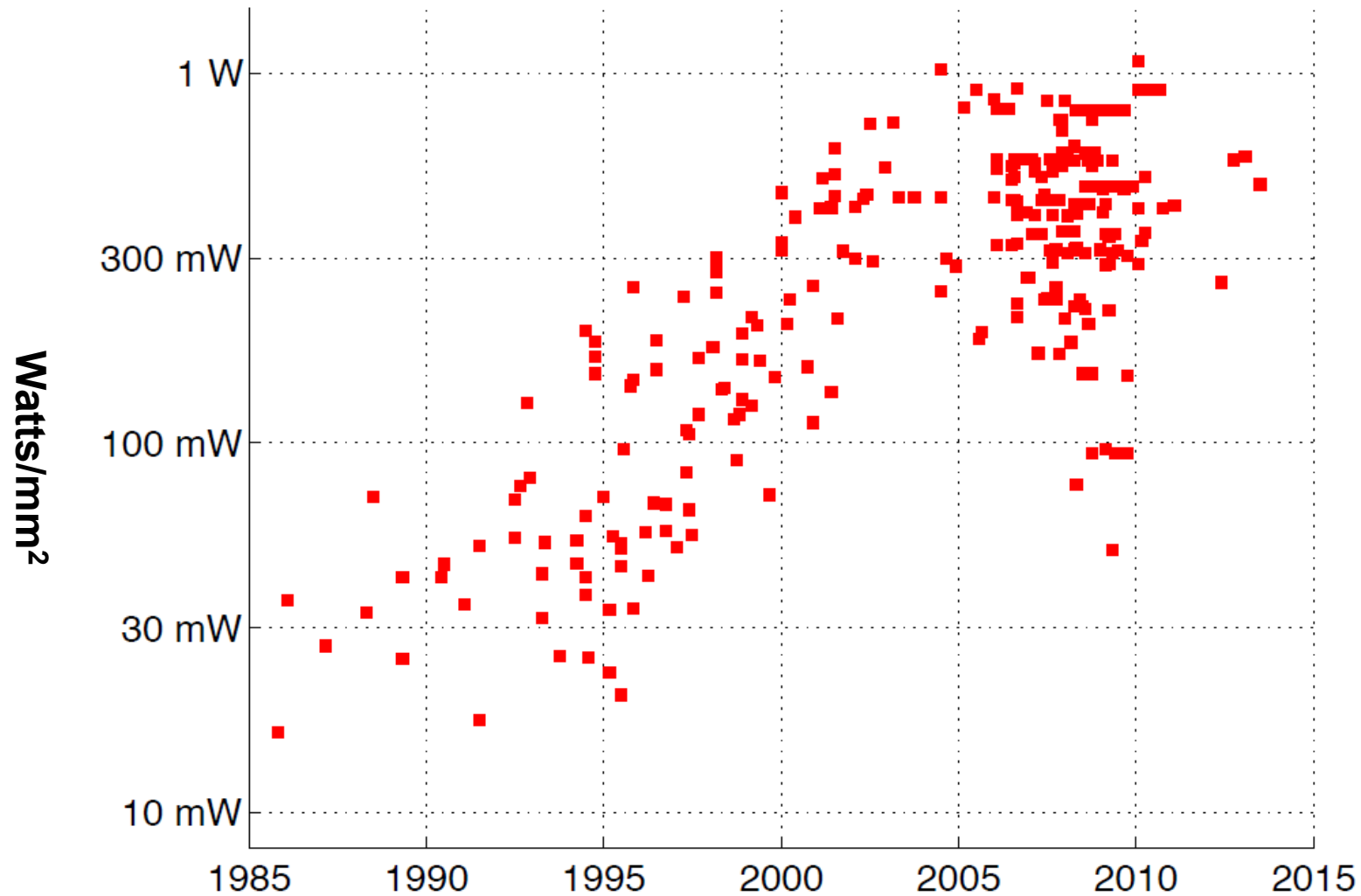CRAY-1

# Supporting Evidence



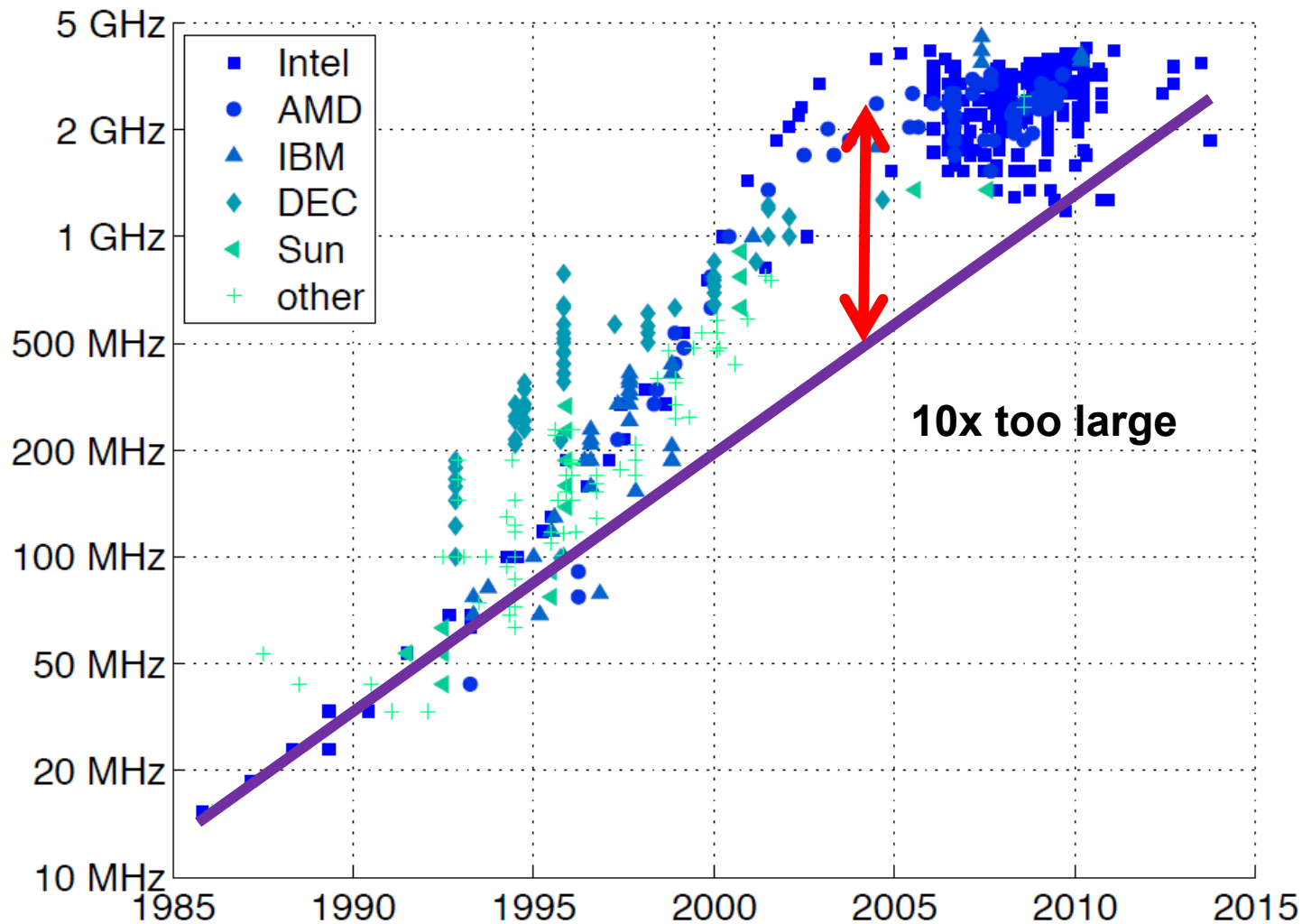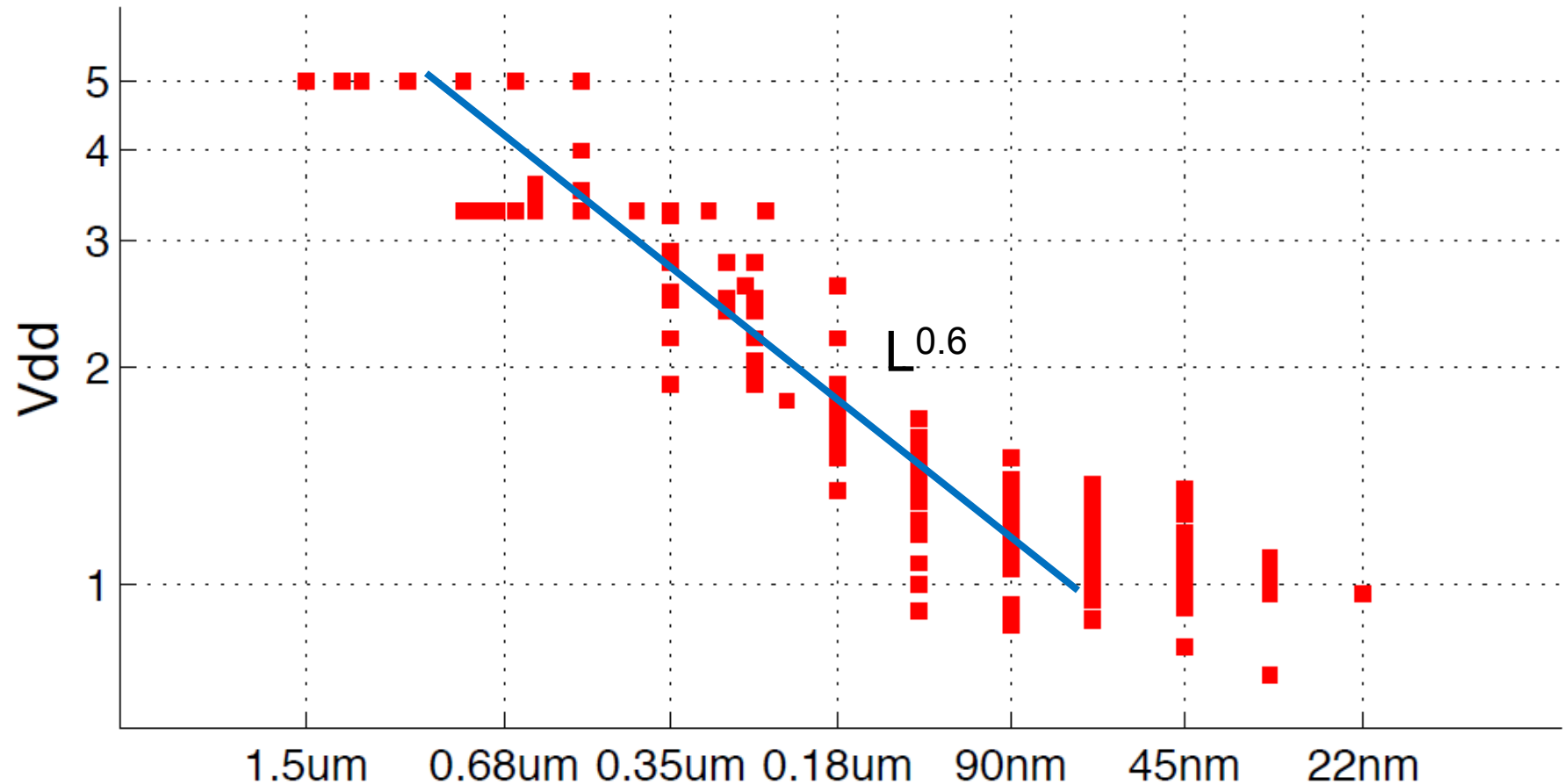http://cpudb.stanford.edu/

# Houston, We Have A Problem

# The Power Limit

# Power Increased Because We Were ~~Clever~~ Greedy



10x too large

# This Power Problem Is Not Going Away:
## $P = \alpha\, C * Vdd^2 * f$



$L^{0.6}$

Vdd

1.5um   0.68um   0.35um   0.18um   90nm   45nm   22nm

http://cpudb.stanford.edu/

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

# Technology to the Rescue?

# Problems w/ Replacing CMOS

## Pretty fundamental physics

- Avoiding this problem will be hard

e-

## Its capability is pretty amazing

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

# Catch - 22

Investment Risk

Capital you need

Very Different ≠ High Risk

Building Computers ≠ Large $

# The Truth About Innovation



## 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

### Power = Energy/Op * 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/

# Aside:
# Throughput Optimized FP

# Aside: Floating Point Optimization
# 180nm – ITRS 10nm

# Have A Shiny Ball, Now What?

# Signal Processing ASICs

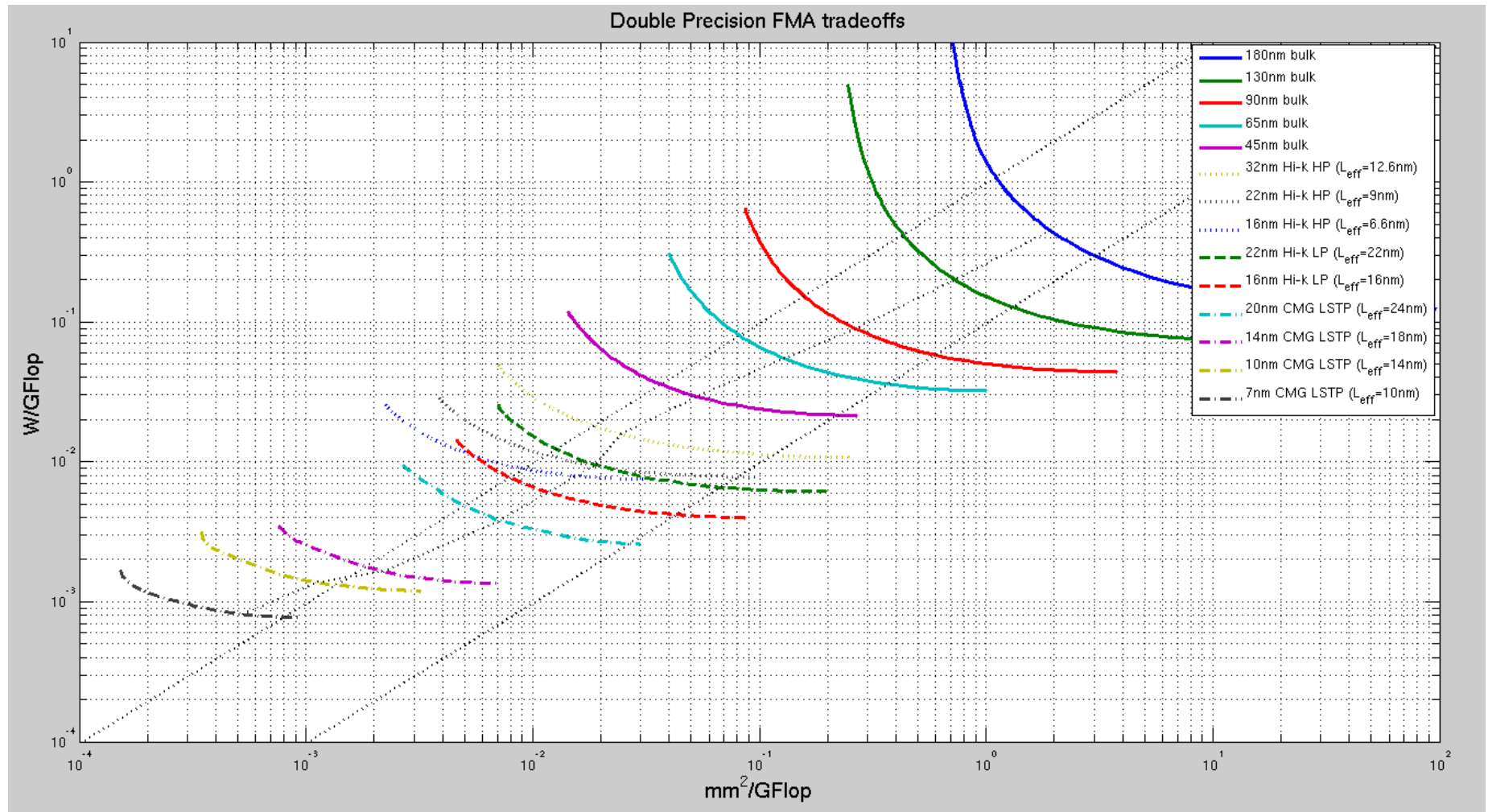# The Push For Specialized Hardware

## Dark Silicon and the End of Multic...

Hadi Esmaeilzadeh[+]  Emily Blem[‡]  Renée St. Amant[§]
[+]University of Washington  [‡]Un...
[§]The University of Texas at ...
hadianeh@cs.washington.edu  blem@cs.wisc.edu ...

**ABSTRACT**

Since 2005, processor designers have increas... ploit Moore's Law scaling, rather than focusing ... formance. The failure of Dennard scaling, to which ... ticore parts is partially a response, may soon limit mu... just as single-core scaling has been curtailed. This p... multicore scaling limits by combining device scaling, ... scaling, and multicore scaling to measure the speedup pot... For device scaling, we use both the ITRS projections and ... of more conservative device scaling parameters. To model sing... core scaling, we combine measurements from over 150 processor... to derive Pareto-optimal frontiers for area/performance and pow... er/performance. Finally, to model multicore scaling, we build a de... tailed performance model of upper-bound performance and lower... 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 2024, only 7.9x average speedup is possible across commonly used 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

## Conservation Cores:
## Reducing the Energy of Mature Computations

Ganesh Venkatesh  Jack Sampson  Nathan Goulding  Saturnino Garcia
Jose Lugo-Martinez  Steven Swanson  Michael Bedford Taylor
Vladyslav Bryksin  Department of Computer Science & Engineering
University of California, San Diego
{gvenkatesh,jsampson,ngouldin,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. Conservation cores, or c-cores, are specialized processors that focus on reducing energy and energy-delay instead of increasing performance. This focus on energy makes c-cores poor candidates 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 support patching, a form of targeted reconfigurability, that allows ... to new versions of the software they target. Our re... ...nservation cores can reduce software energy consumption by up to 2.1x for whole application ... ...ations and by up to ... ...and the useful lifetime of individual ... processors.
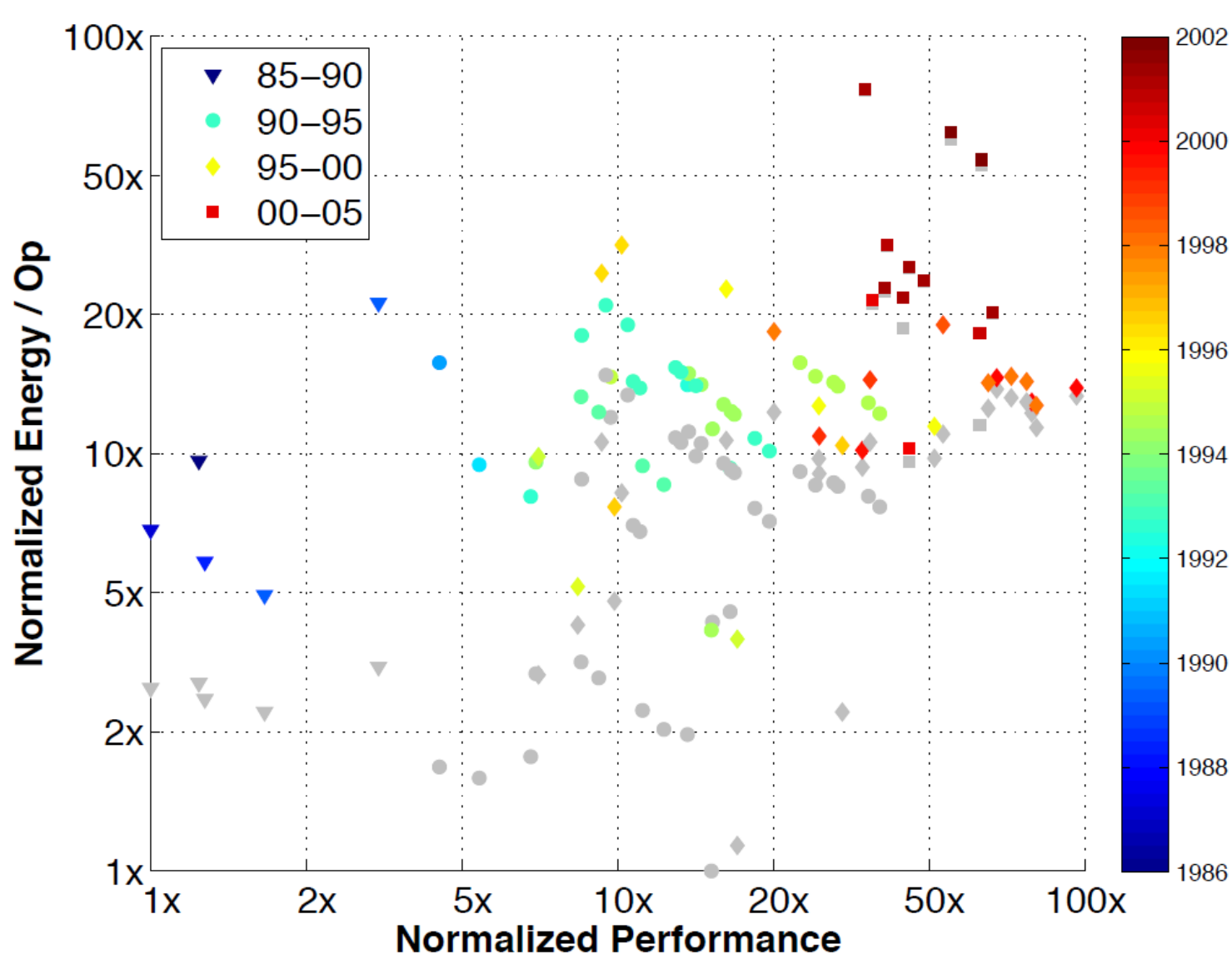
# Before Talking About Specialization

WE SHOULD CHECK
ONE MORE THING FIRST

# Don't Forget Memory System Energy

# Processor Energy w/ Corrected Cache Sizes

# Processor Energy Breakdown
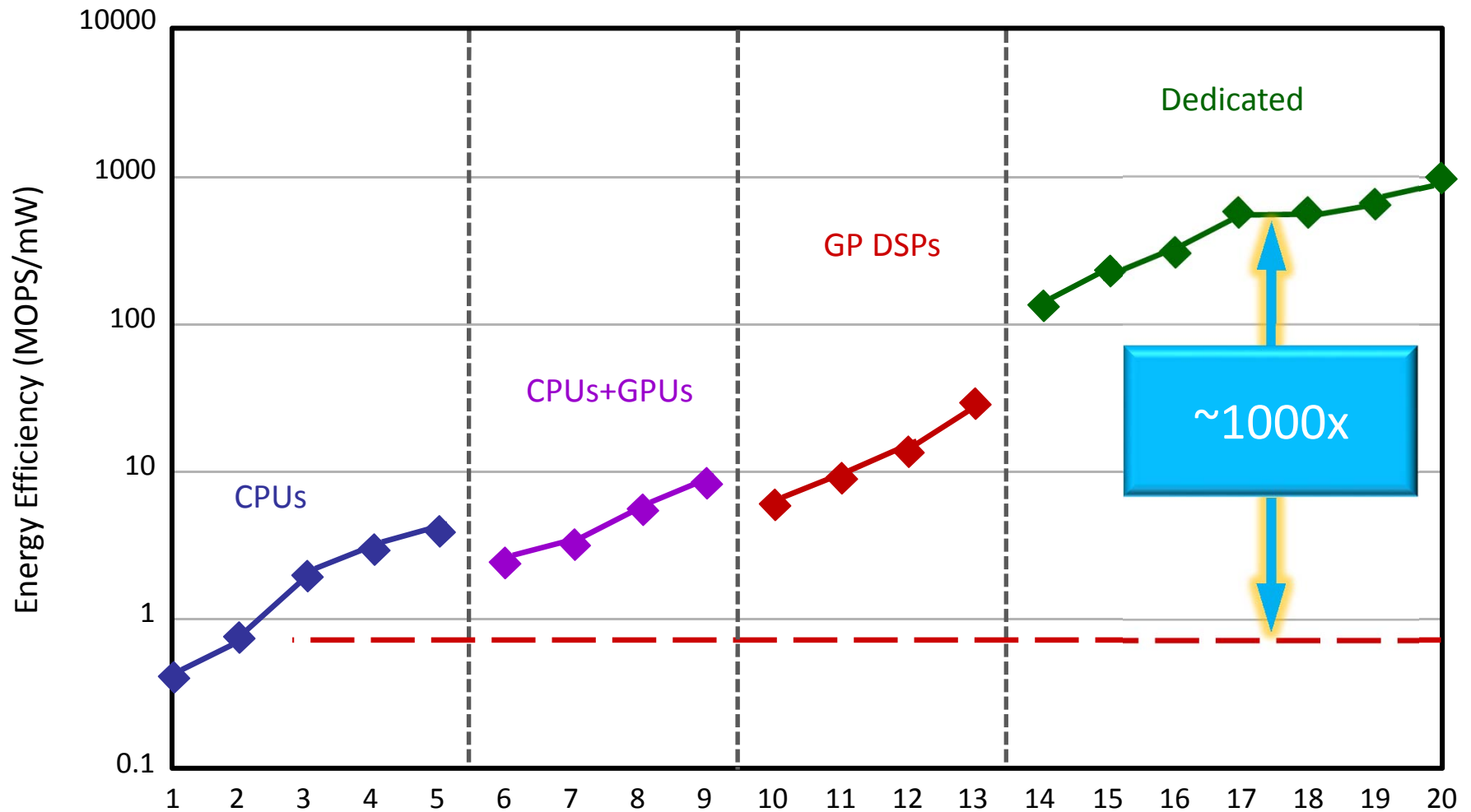


Legend:
- 8 cores
- L1/reg/TLB
- L2
- L3

# Data Center Energy Specs
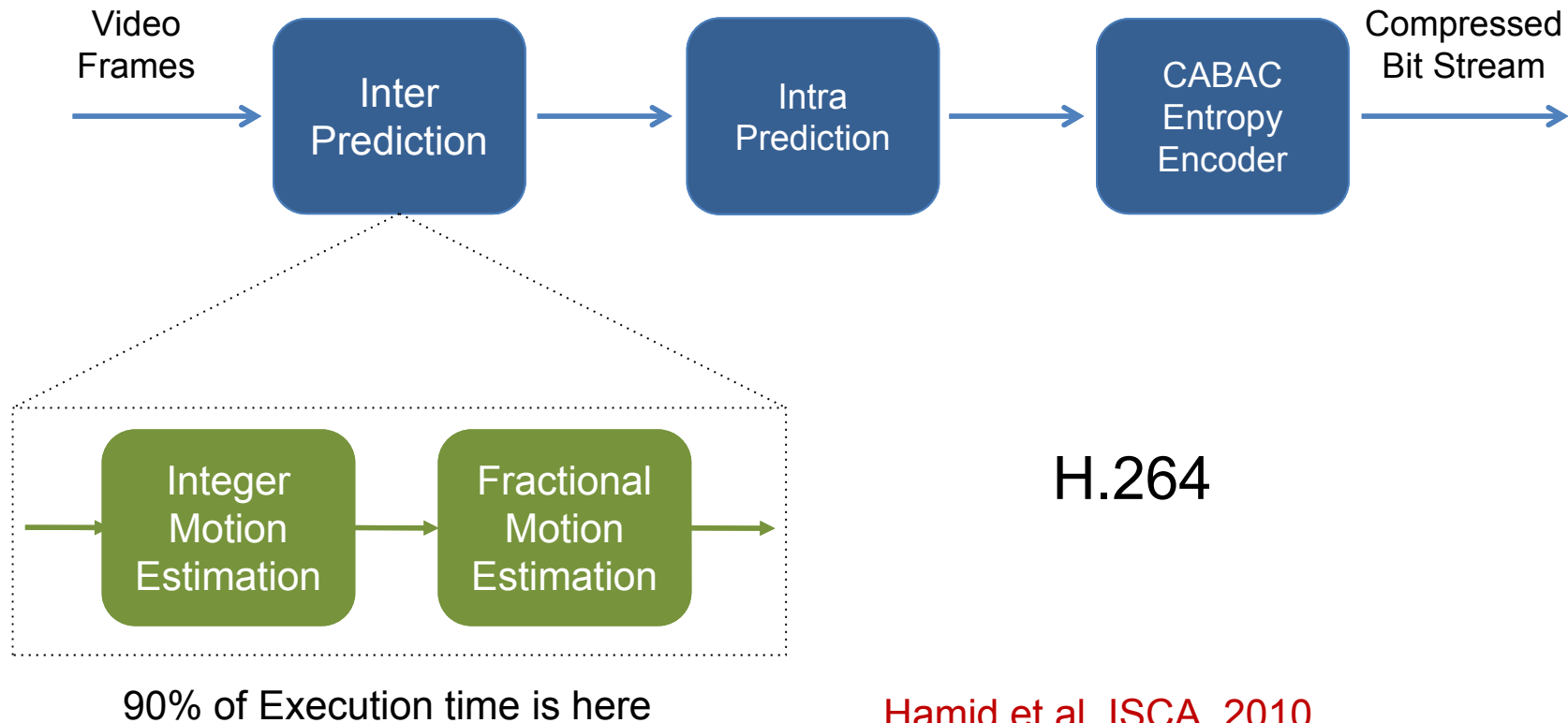


Malladi, ISCA, 2012

# 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



H.264

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

| 25pJ | 6pJ | Control | | 70 pJ |

I-Cache Access    Register File Access    Add

# The Truth:
## It's More About the Algorithm then the Hardware

All Algorithms

GPU Alg

# Highly Local Computation Model

# Highly Local Computation Model

# Highly Local Computation Model

# Highly Local Computation Model

# Compose These Cores into a Pipeline
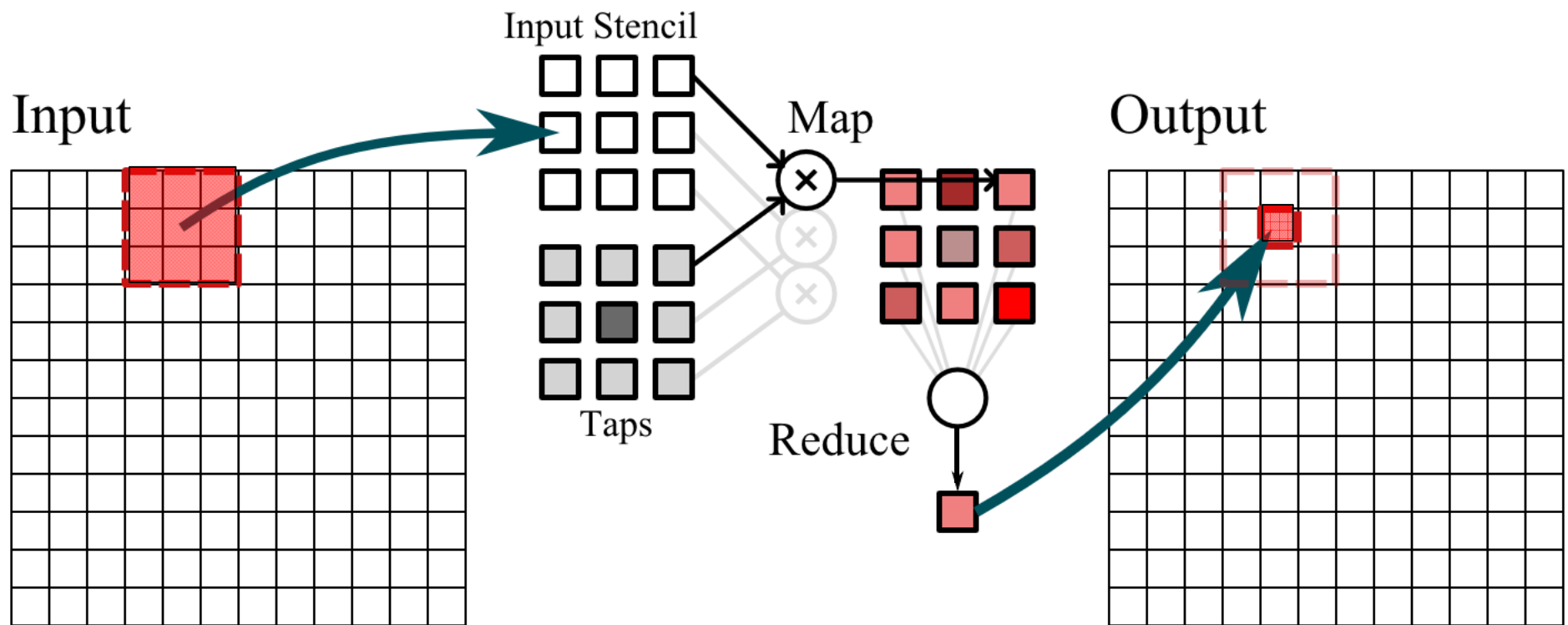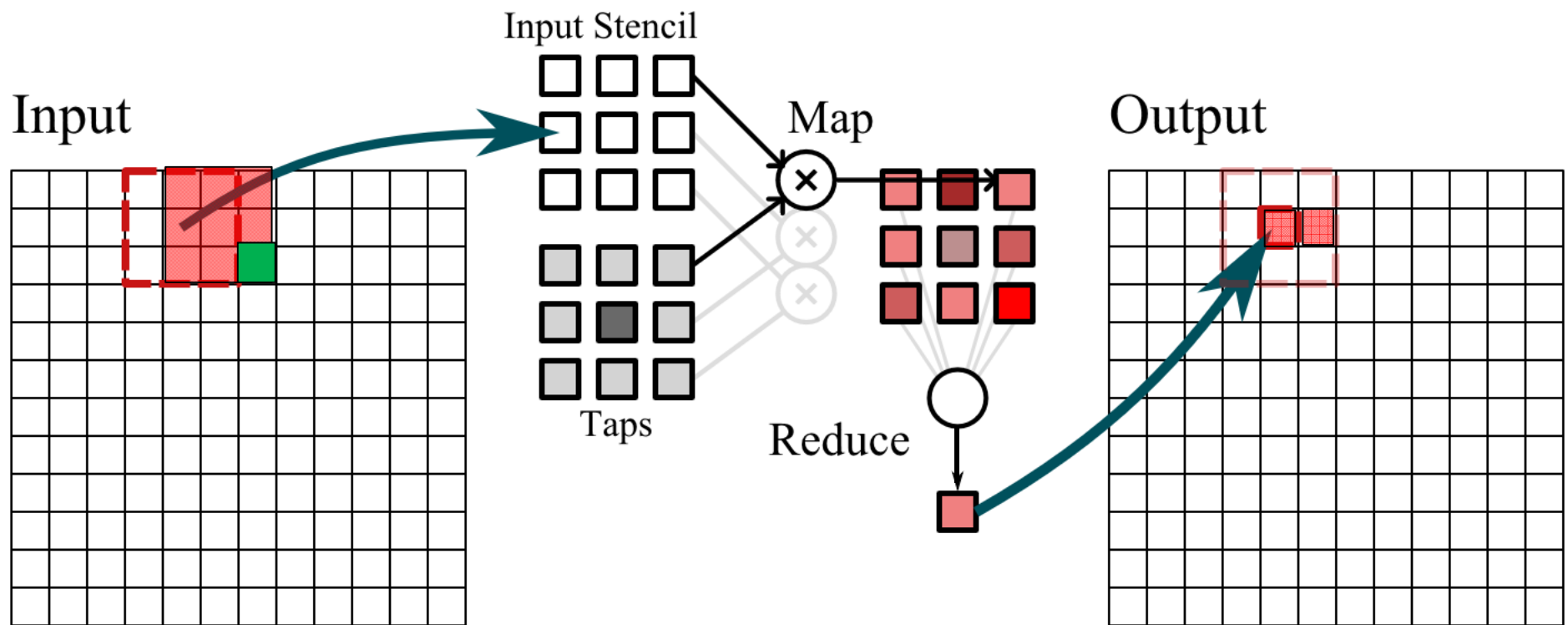


Intermediate Working Set

StencilA

StencilB

**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?



**Need to reduce cost to play**
- Building constructors, not instances

*(y-axis: Investment Risk; x-axis: Capital you need)*

http://genesis2.stanford.edu/

# Chip Constructor Idea



RTL, Verif Collateral, Firmware

http://genesis2.stanford.edu/

# Not All Systems Are On The Bleeding Edge

# App Store For Hardware

# Challenge



ARDUINO IS AN OPEN-SOURCE ELECTRONICS PROTOTYPING PLATFORM BASED ON FLEXIBLE, EASY-TO-USE HARDWARE AND SOFTWARE. IT'S INTENDED FOR ARTISTS, DESIGNERS, HOBBYISTS AND ANYONE INTERESTED IN CREATING INTERACTIVE OBJECTS OR ENVIRONMENTS.

## 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