

CSE 141: Introduction to Computer Architecture

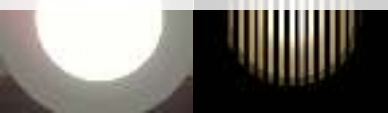
Pat Pannuto, UC San Diego

ppannuto@ucsd.edu

OInK
To:ck

Human
Perception

Camera
Perception



What is Computer Architecture and where does it fit in Computer (Science) Engineering?

- One view: what is an Architect and how do they fit in the creation of buildings?

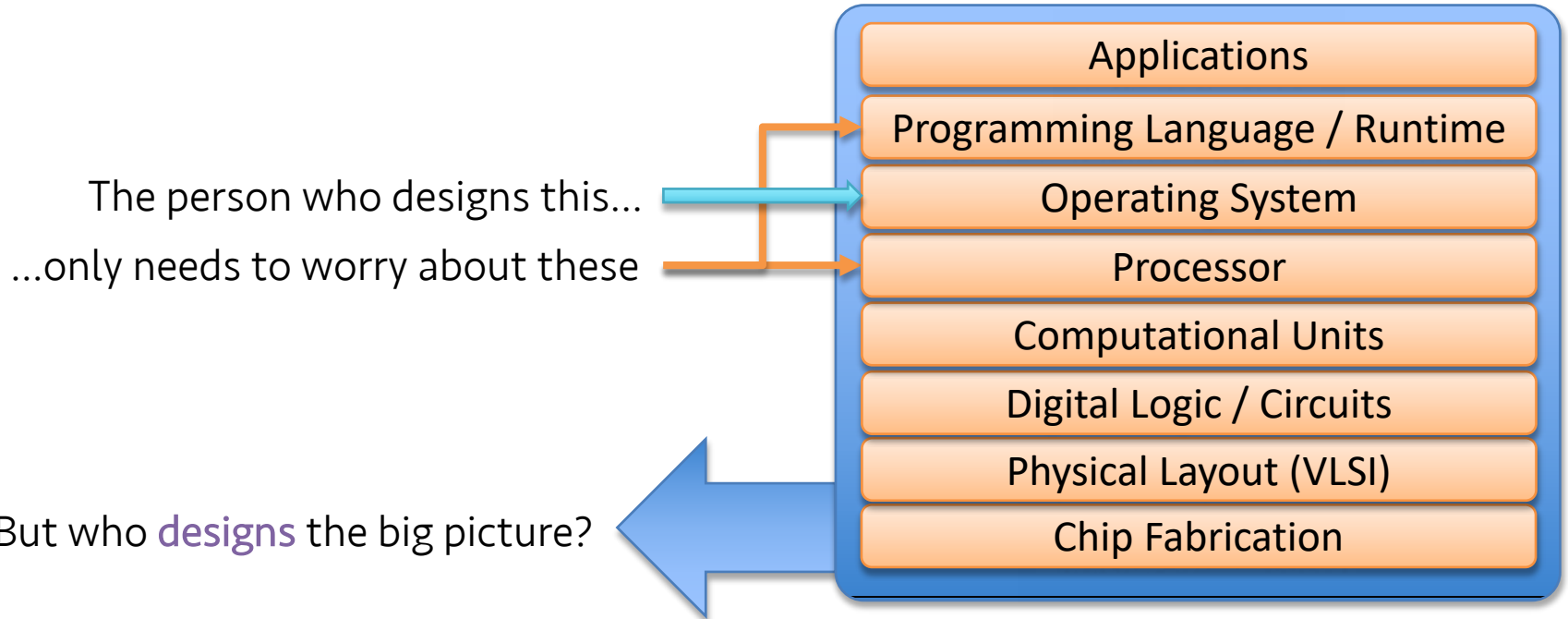


Zaha Hadid

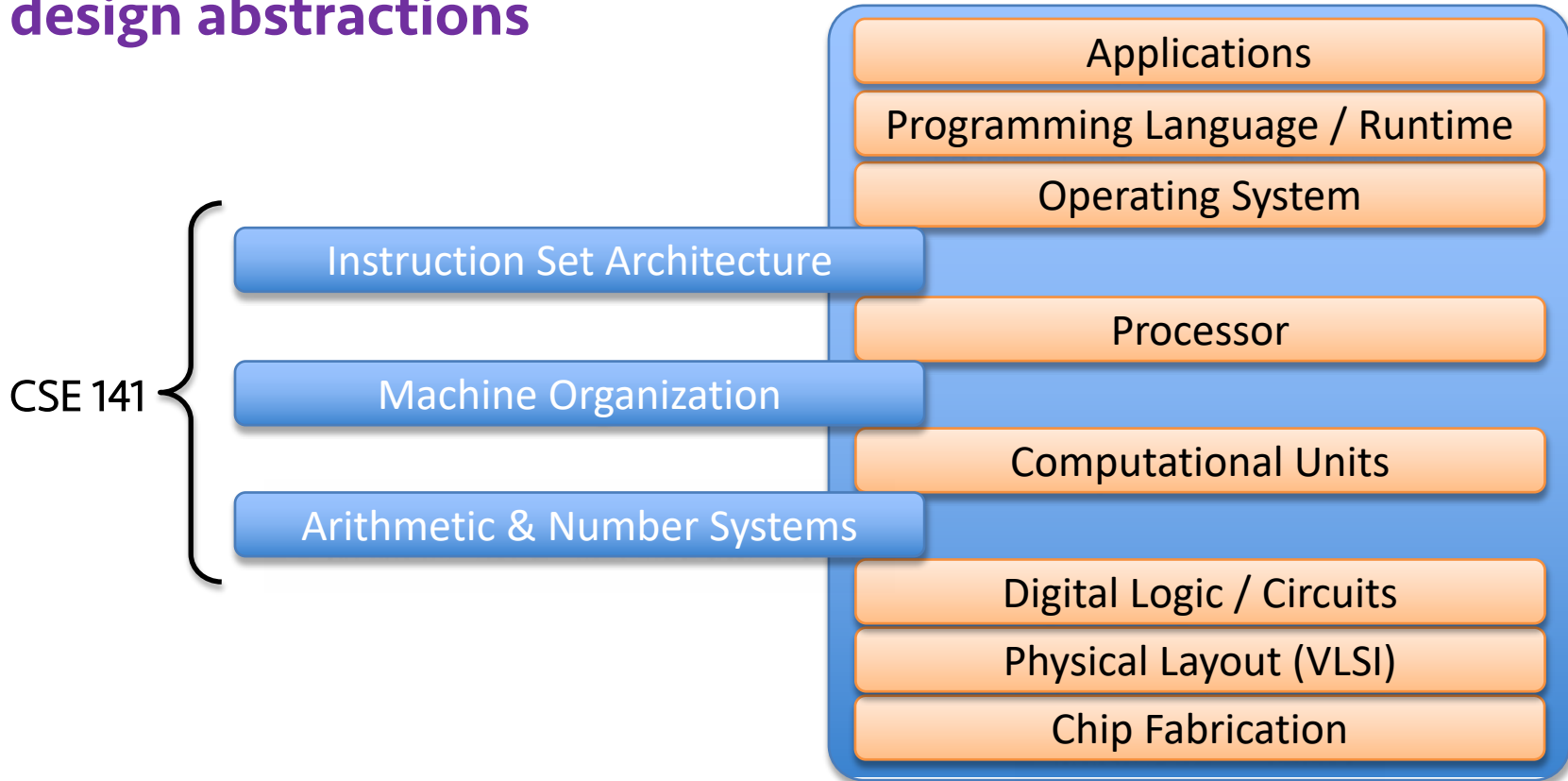


Port Authority Building in Antwerp, designed by Zaha

Computer science is all about abstractions



Computer architects look at the system as a whole and design abstractions



Good abstractions make it easier to focus on reasoning about one part of a large, complex system

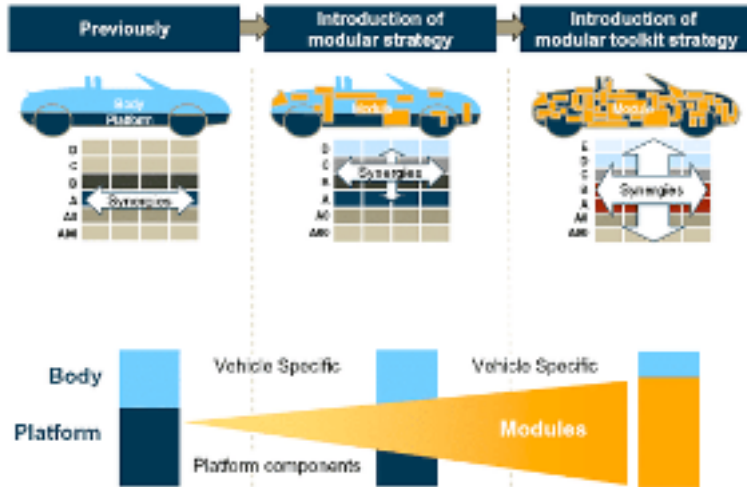
- Which of these maps is easier to use to plan a trolley trip?



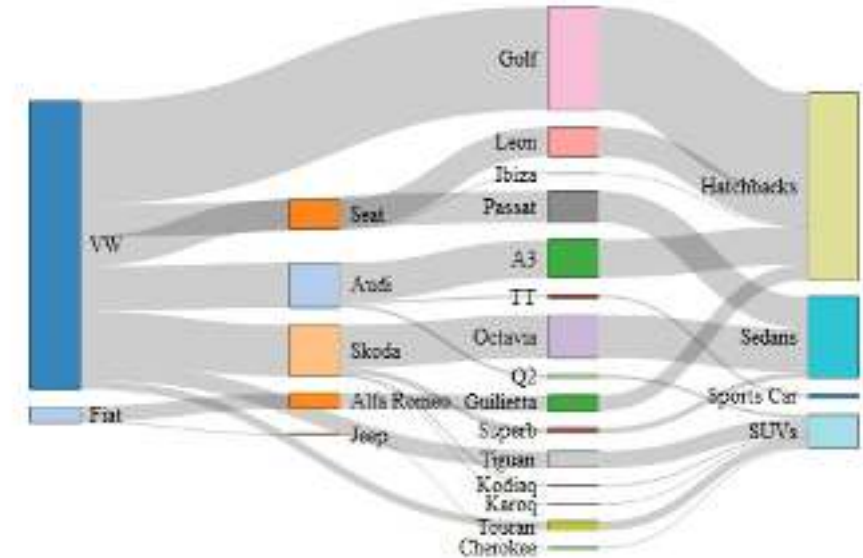
Good abstractions make it easier to focus on reasoning about one part of a large, complex system

- Modularization is fundamental to design in many domains

Volkswagen Group's Modular Toolkit Strategy



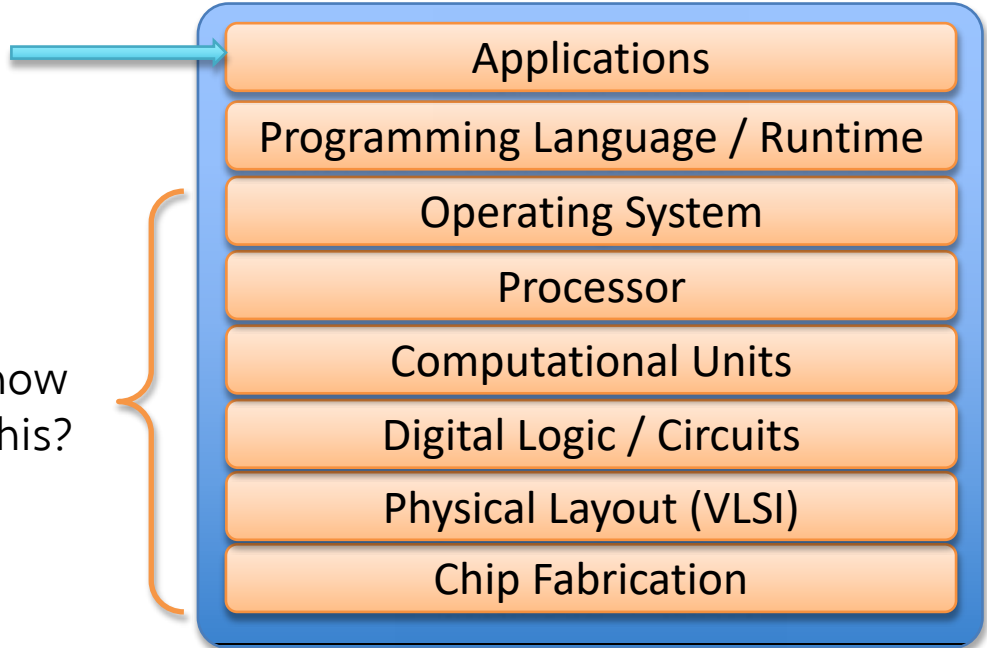
Modular Car Body Design and Optimization by an Implicit Parameterization Technique via SFE CONCEPT
Fabien Duddeck, Hans Zimmer



https://www.reddit.com/r/dataisbeautiful/comments/8m15g9/automobile_platform_sharing_work_in_progress/

But what if I'm not going to become a computer architect?

If I only want to build these...



...why do I need to know about any of this?

The real world is full of leaky abstractions

- Goal: Sum up all the entries of a two dimensional array
- Which of these implementations is faster?

```
int twoDarray[256][256];
int sum = 0;

for (int i=0; i<256; i++) {
    for (int j=0; j<256; j++) {
        sum += twoDarray[i][j];
    }
}
```

```
int twoDarray[256][256];
int sum = 0;

for (int i=0; i<256; i++) {
    for (int j=0; j<256; j++) {
        sum += twoDarray[j][i];
    }
}
```

Answer: "It depends"

Architects look across systems to find and improve inefficiencies – always valuable, sometimes critical...

- What happens when workload changes overnight, and there's no way to buy your way out of the problem?
 - Architects help to fix it

Microsoft Admits Supply Chain Issues Hit Cloud Server Supply

Microsoft
2020



Data Center Hardware Shortages in the Age of COVID-19

Coronavirus: Supply issues hit PC market

Despite strong demand, the pandemic caused problems in sourcing kit and forced the PC market to reverse three quarters of growth

By Simon Quicke, Microsoft Editor

Published: 14 Apr 2020 10:28

Coronavirus has had an impact on PC supply chains despite the hard work of the industry to counter the negative impacts.

work, and those changes put new demands on IT that have been hard to

more people work from home and remain connected with friends and family

increase in network capacity utilization.

Google and Amazon, have had outages as they rushed to add more data

But there's a basic problem afflicting everyone in the technology world. As the world searches for a light at the end of the COVID-19 tunnel, leading tech companies are working to overcome supply chain shortages of critical components for servers, storage, and networking products.

Course Administrivia



- Instructors
 - Sec A00: Pat Pannuto
 - Sec B00: Dean Tullsen
- The sections will run “in sync”
 - The lectures will cover the same content but in different ways
 - Assignments will be similar enough that folks should be able to work in groups across the sections, but may not be identical
 - Unified piazza, office hours
 - We will give exams at the same time
 - Note: The registrar assigned our final exam as **Saturday, December 12** – Plan now!!

Course Staff

- Four amazing TAs:

- Nitish Kulshrestha

- Shanti Modi

- Sumiran Shubhi

- Kazem Taram



- Discussions

- Sec A00: Wed from 11:00 to 11:50

- Sec B00: Tue from 14:00 to 14:50

- You may attend either, but A00 better when possible



Assessments & Workload

- Grading
 - Participation: 5%
 - Weekly **participation** quizzes of lecture material (10 weeks -> 0.5% / quiz!)
 - Homework: 20%
 - Midterm: 30%
 - Final Exam: 45%
 - (Inclusive final)
 - Our assigned final exam slot is **SATURDAY, December 12** — Plan for this now!

Repeated, active engagement is key to effective learning

- Pre-class reading is your first exposure
 - 5 minutes before class is better than not at all, but 5+ hours before is much better
 - **Read actively**, try writing notes for yourself of what you understood from readings
- Watching video is not a passive activity
 - Ask (or write down) questions about what you do not understand!
 - **Use checkpoints effectively**
- Discussions, office hours, and exercises are not passive activities
 - **Work through examples yourself** and ask the questions you have
- Homework is designed to help you solidify your understanding
- Study for tests “honestly to yourself” – **you** must engage with questions

Class is not a competition

- My philosophy
 - I care whether you learn the material
 - The purpose of a grade is to assess how well you know the material in 141
 - The purpose of a grade is not “rank” students
 - I am most successful if everyone in class **earns** an A
- My goal is not to curve
 - (But I reserve the right to)
 - The midterm and final may be “internally” curved

Academic Integrity

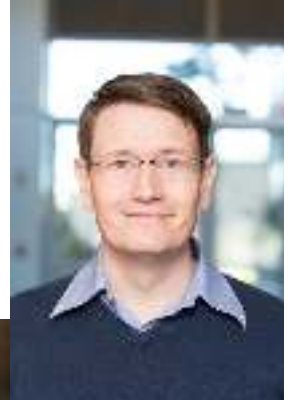
- Cheating will be taken very seriously
- Examples
 - Not cheating:
 - Discussing homework in groups, with **your own writeup, in your own words, done on your own, later**
 - Cheating:
 - Getting a walk-through from someone who has already done the homework
 - Looking at someone else's completed work (even "just to check")
 - Using solutions from the web, prior classes, or anywhere else
 - Receiving, providing, or soliciting assistance from another student during a test
- An experiment: **Regret Policy**

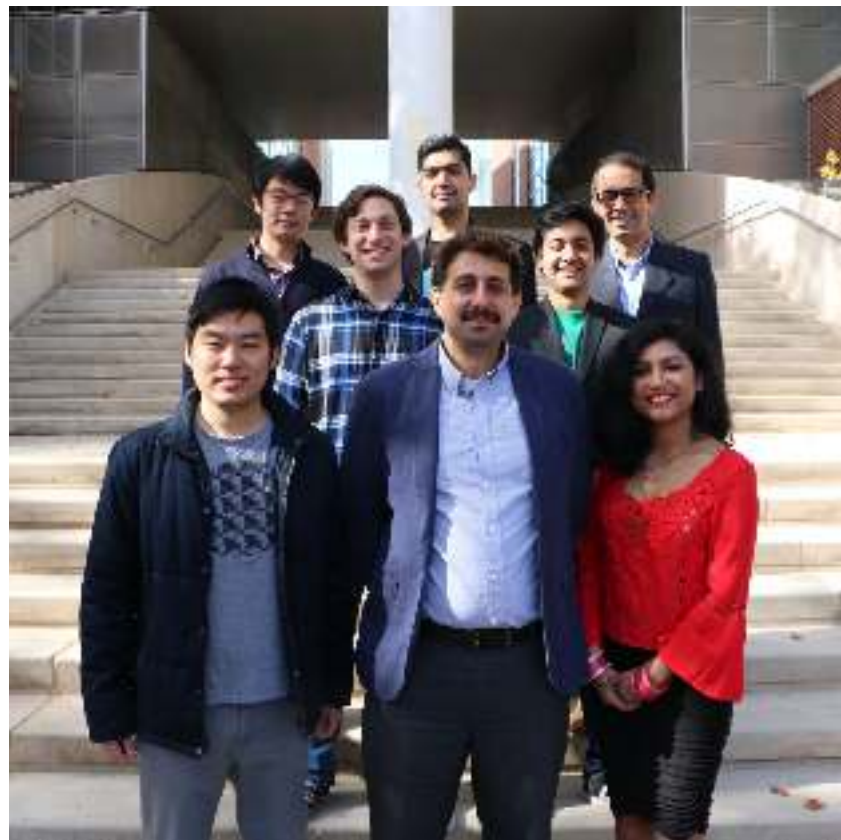
We'll take a short break here...

AND THEN SOME MODERN HIGHLIGHTS FROM HERE AT UCSD

But for the rest of today, I want to highlight the kinds of cool stuff that architects *do*

- UCSD has an amazing team of architecture faculty



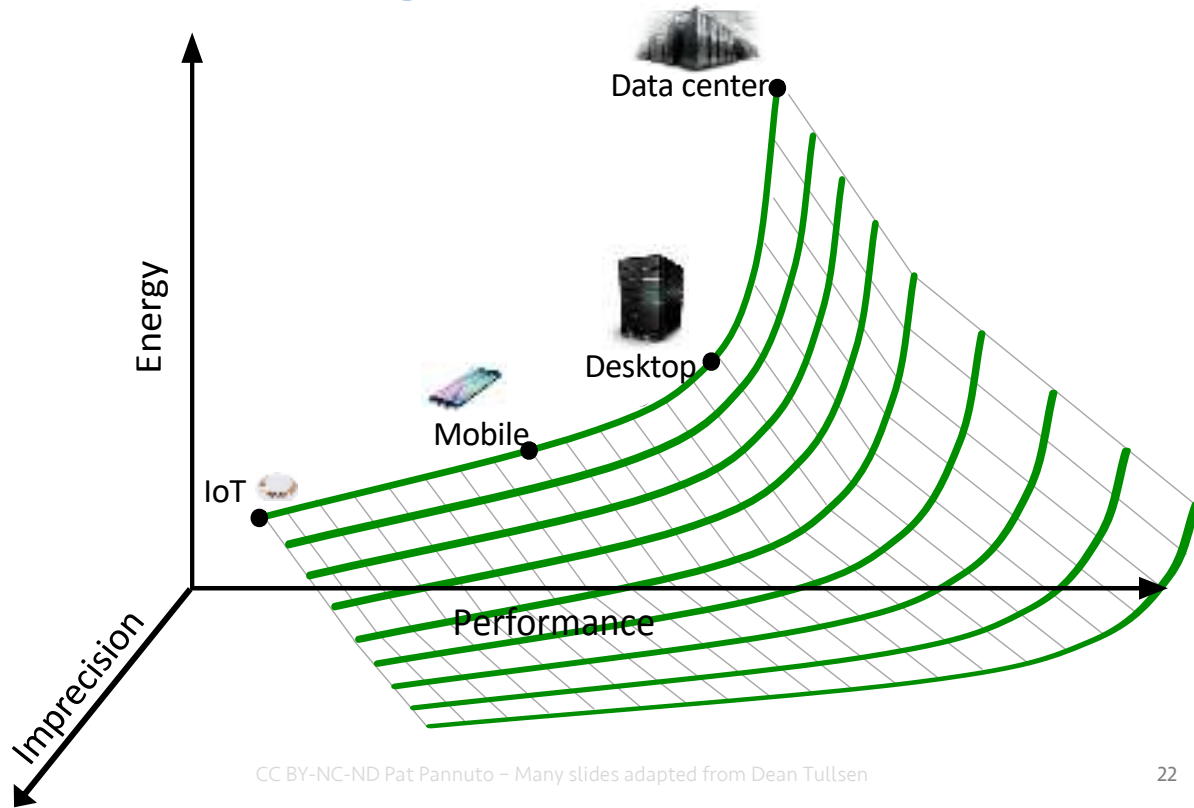


One wild idea: “Approximate Computing”

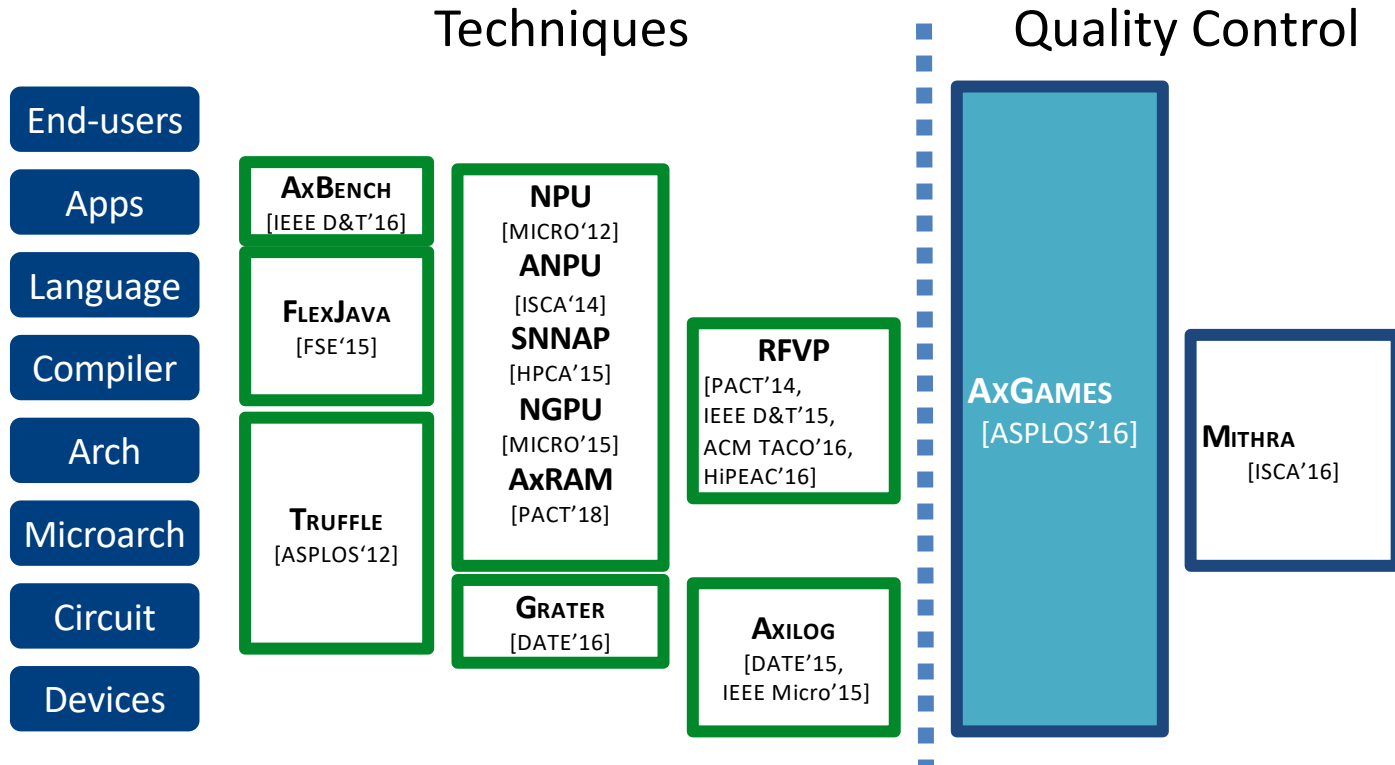
- Aka, what if $1 + 1$ doesn't *always* equal *exactly* 2?



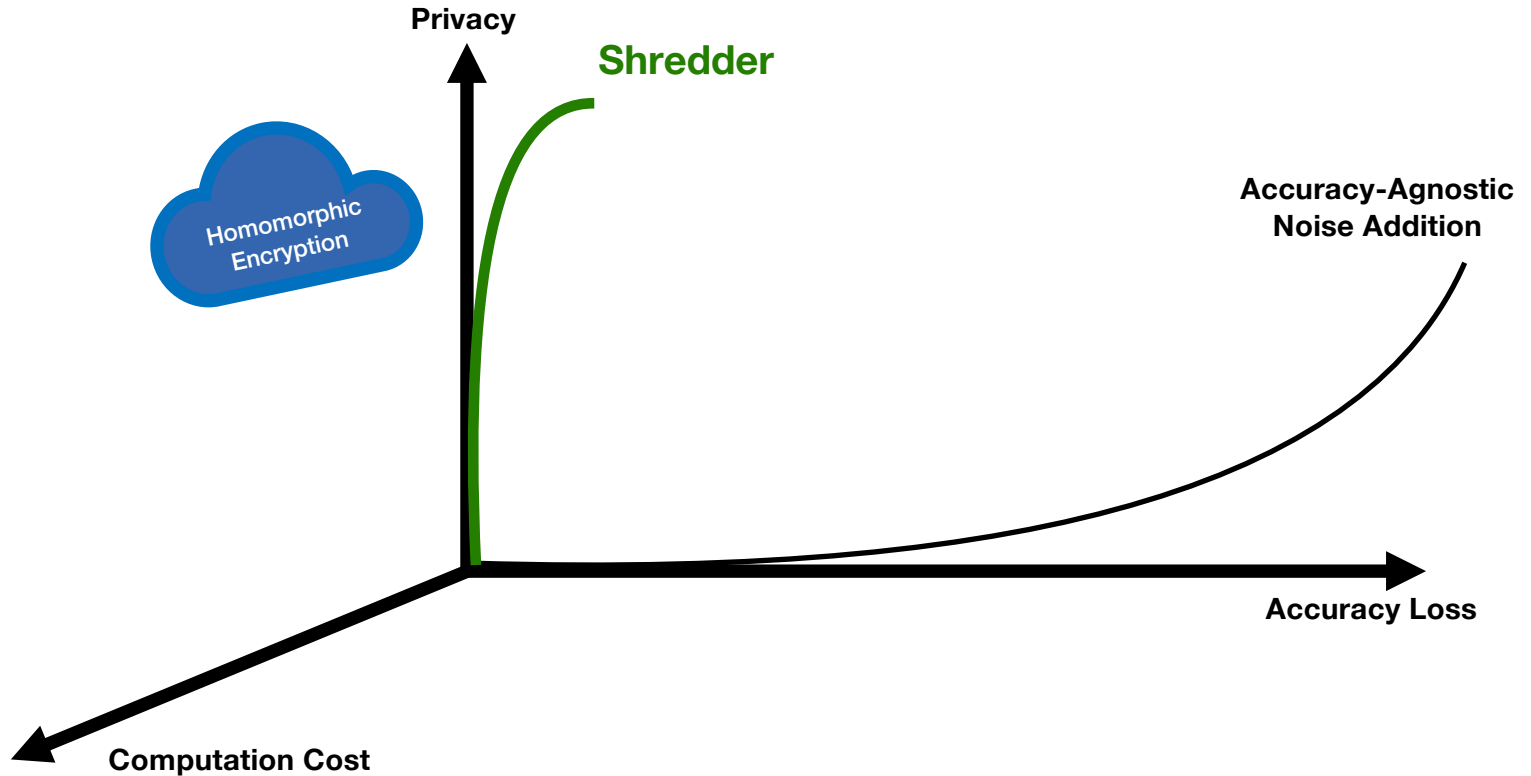
Embracing imprecision allows for major gains in performance and energy



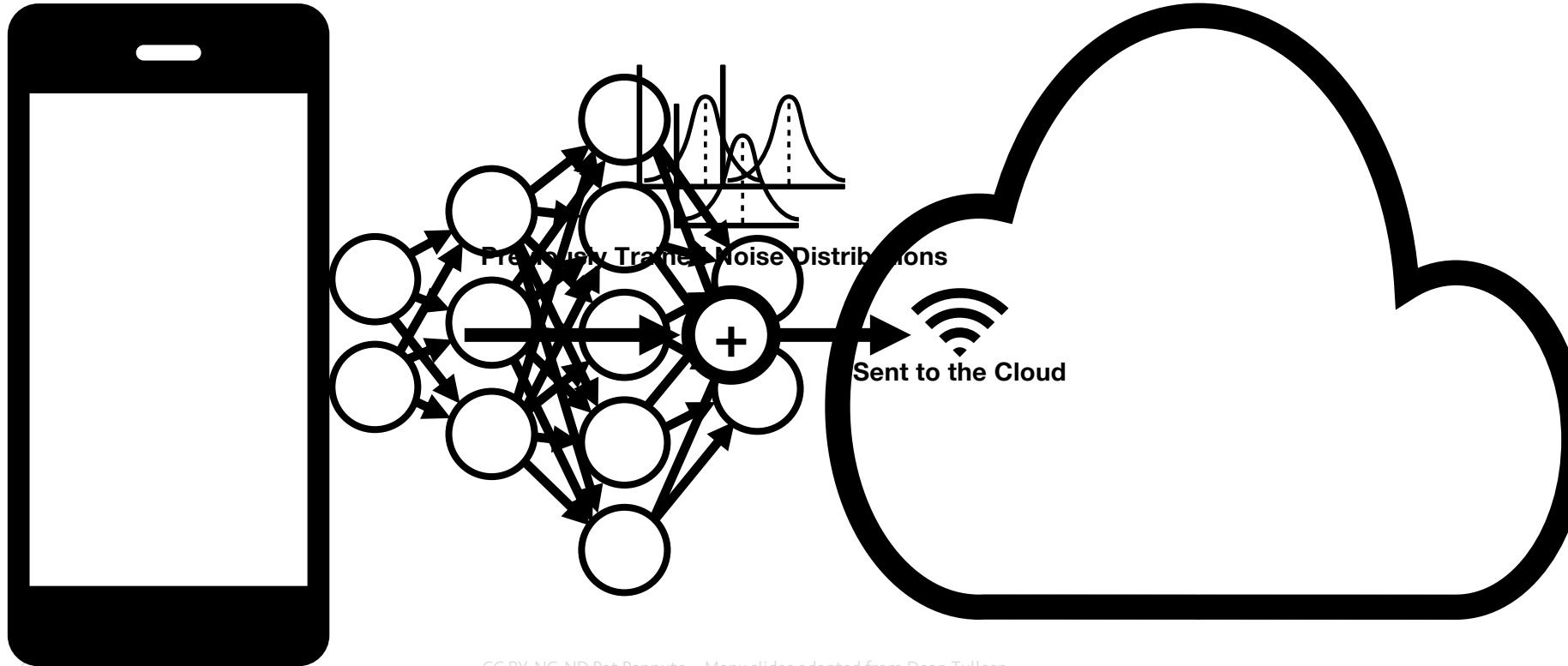
A cross-stack approach to enable approximation



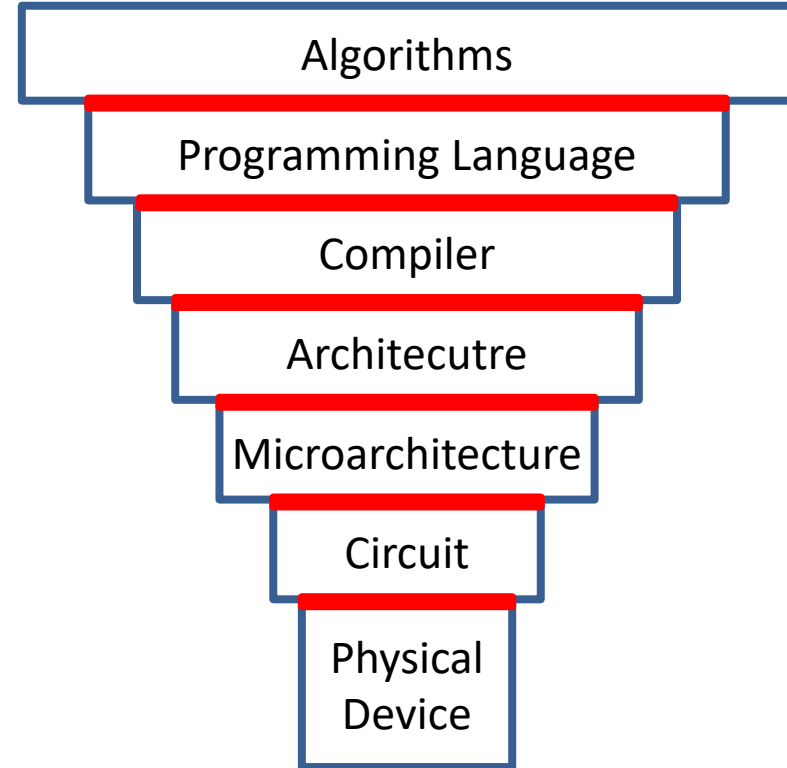
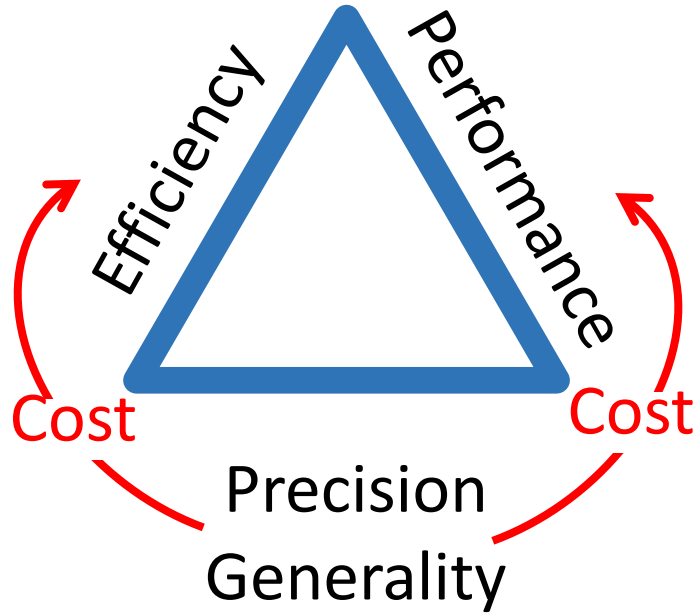
Privacy Preserving Techniques for Inference



Execution Model



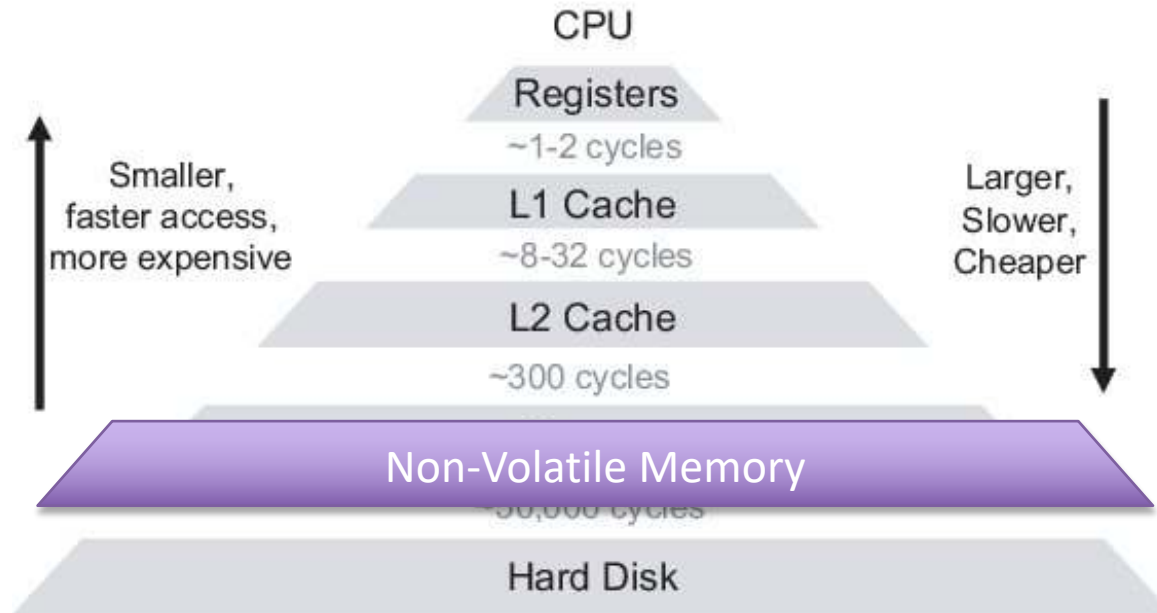
Rethinking the abstractions



Memory, Storage, Software, and Architecture in the NVSL



This is a slide you will encounter in many CE/CSE classes...



Applications

MARS

Willow

NOVA

Orion

Ziggurat

SubZero

NV-Heaps

Pangolin

Pronto

Tools

Libraries

Stacks

Operating Systems

Distributed Systems

Moneta

QuickSAN

3D XPOINT



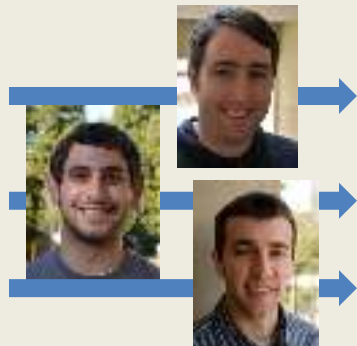
Core



NVSL Students Lead Industry

- We Built

- Opt. SSD interface (2009)
- Direct, remote SSD (2013)
- First PCM SSD (2011)
- PMEM prog. tools (2011)



- Industry Built

- NVMe (2011)
- NVMe over Fabrics (2016)
- Optane (2016)
- PMDK (~2014)

Mobilizing the Micro-Ops: Exploiting **Context Sensitive Decoding** for Security and Energy Efficiency



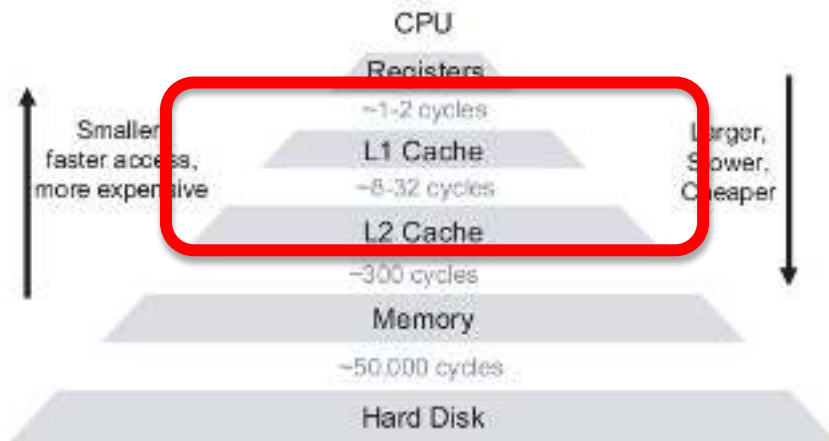
Leaky abstractions are not always just performance problems...

- This loop behaved differently because of how caches work

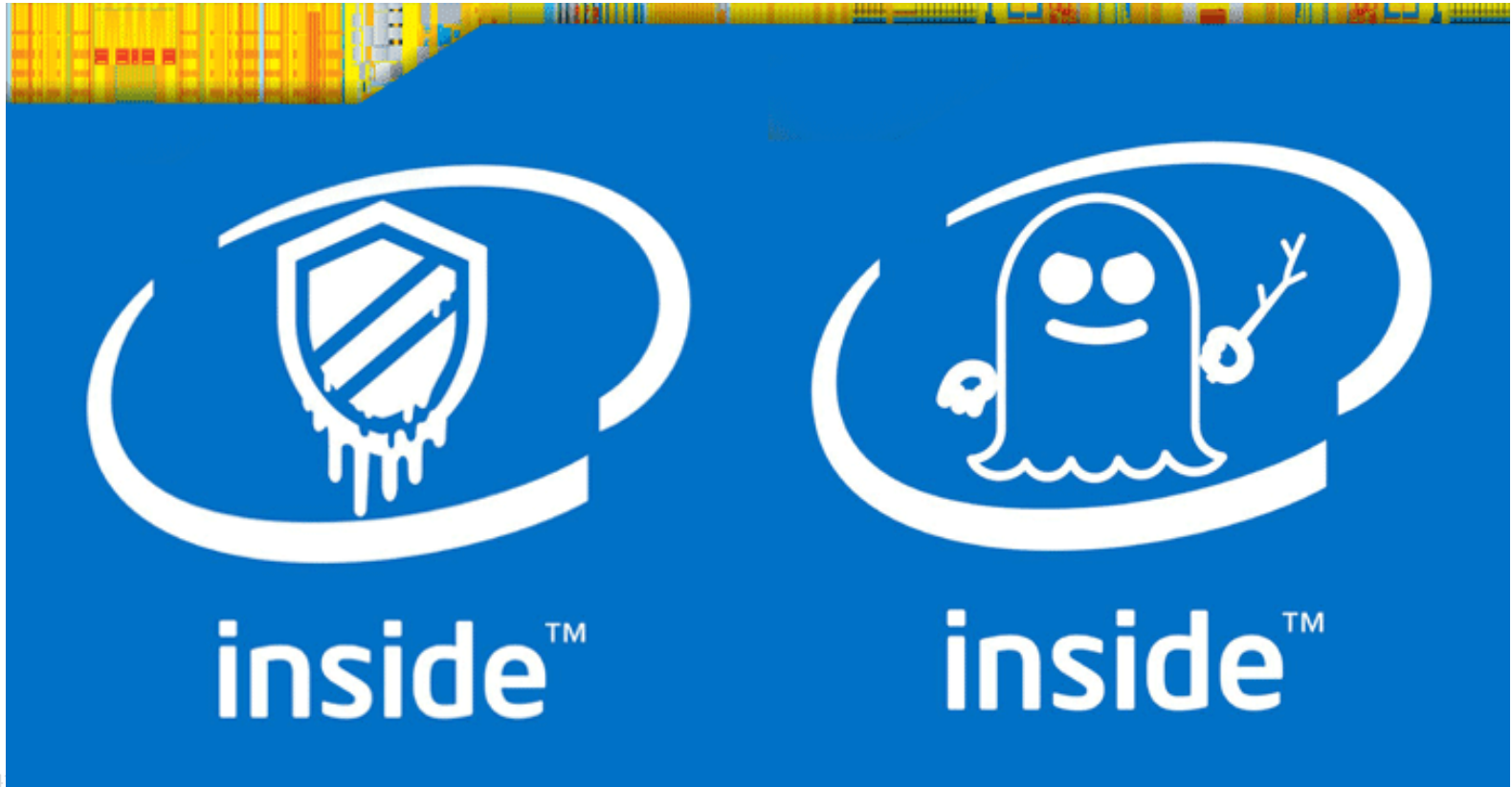
```
int twoDarray[256][256];
int sum = 0;

for (int i=0; i<256; i++) {
    for (int j=0; j<256; j++) {
        sum += twoDarray[i][j];
    }
}
```

Architects added “hidden” caches:
faster, intermediate memories

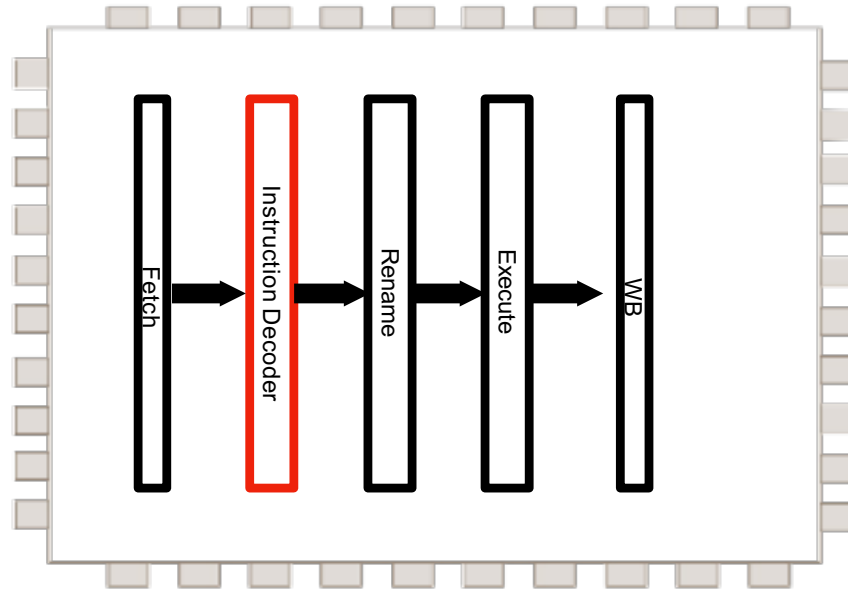


Leaky abstractions can be security threats!



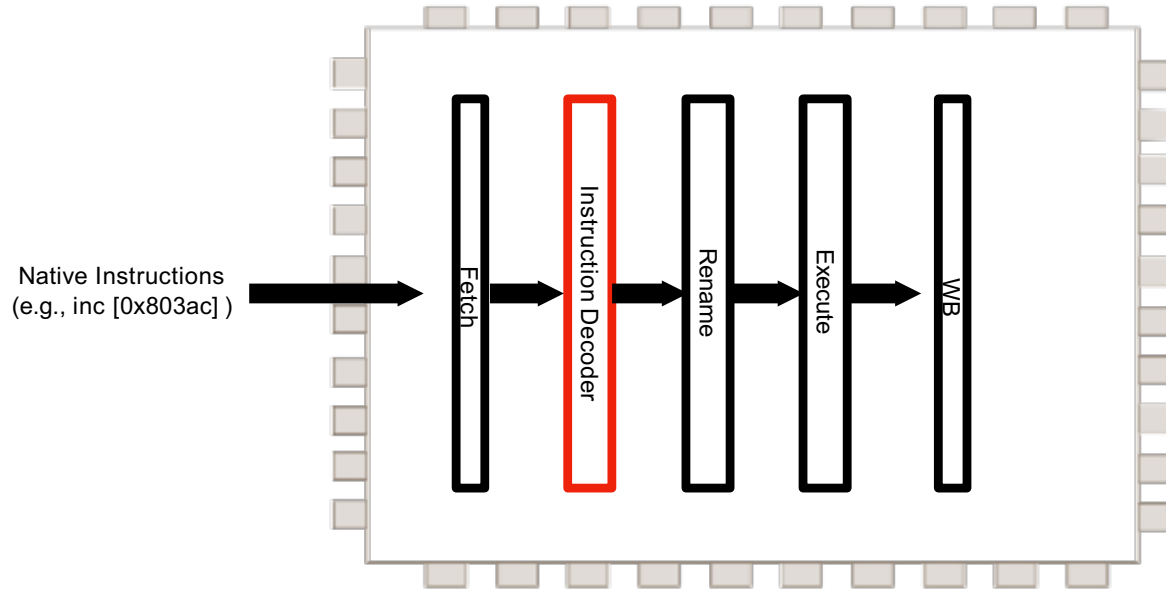
Mobilizing the Micro-Ops

Exploiting Translated ISAs



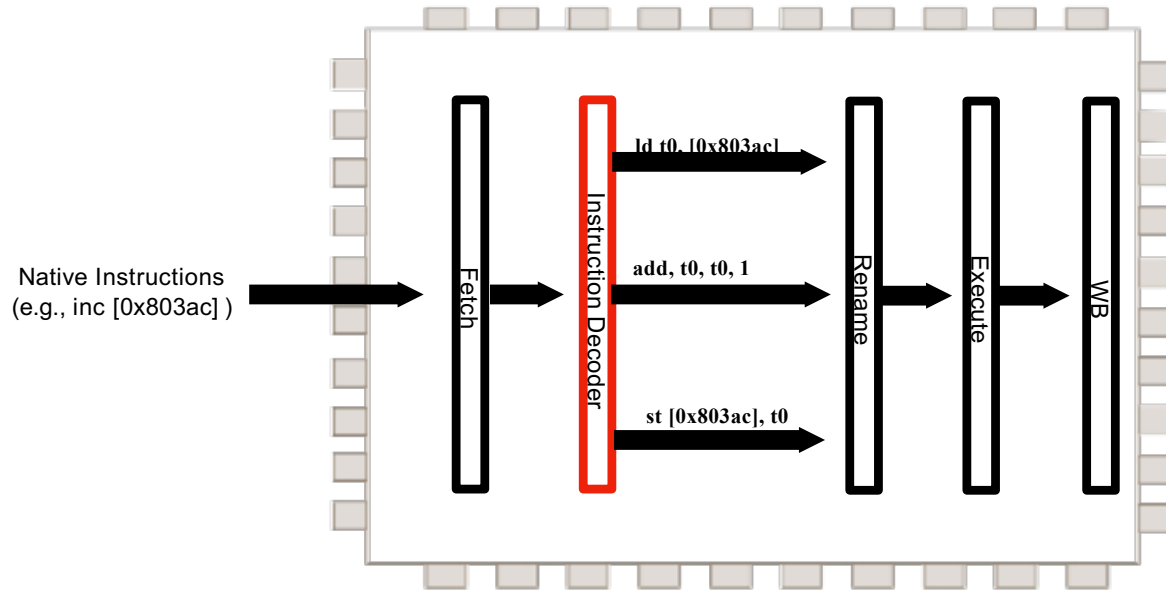
Mobilizing the Micro-Ops

Exploiting Translated ISAs



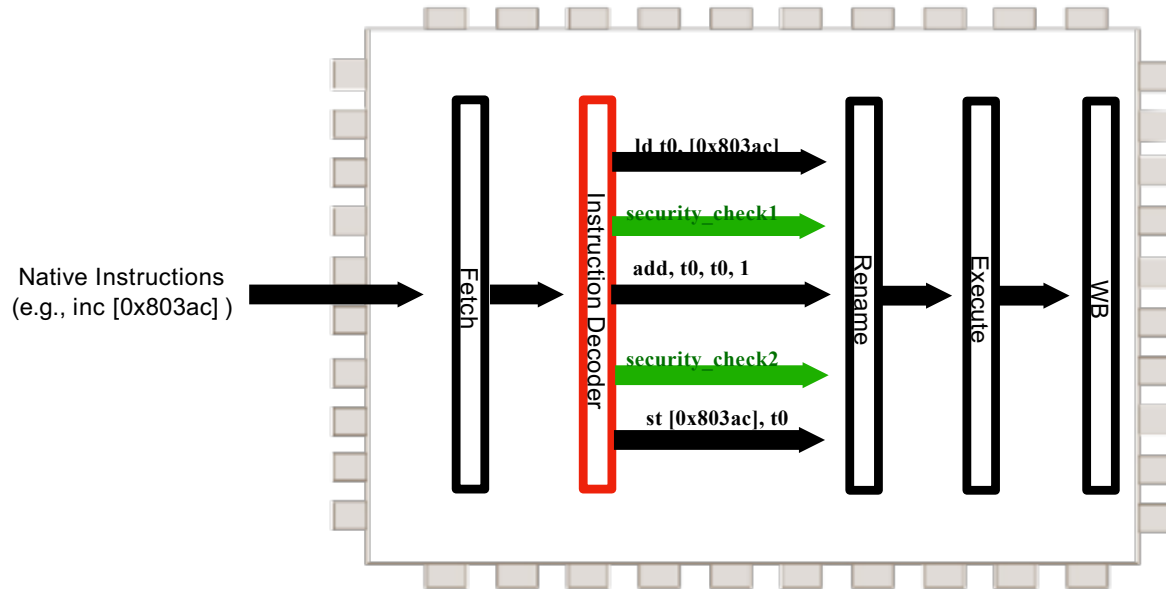
Mobilizing the Micro-Ops

Exploiting Translated ISAs



Mobilizing the Micro-Ops

Exploiting Translated ISAs

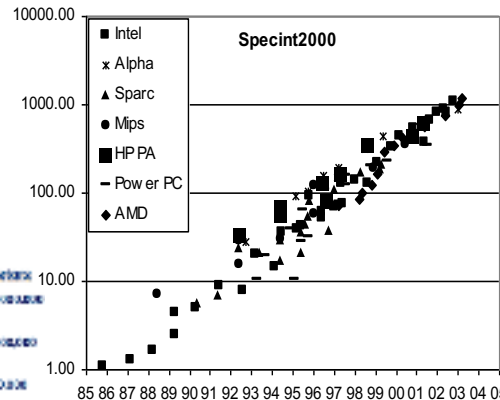
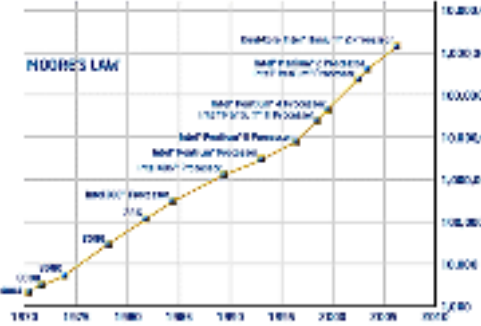
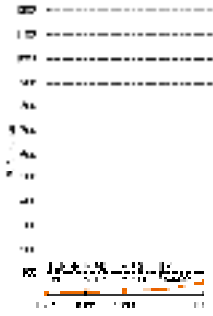


Context Sensitive Decoding fixes a leaky abstraction

- Eliminating cache side channels via cache obfuscation
- Energy and Performance optimization via selective devectorization
 - ISCA 2018
 - IEEE Micro Top Picks in Computer Architecture
- Spectre mitigation via targeted insertion of fence micro-ops (**Context Sensitive Fencing**)
 - ASPLOS 2019

Performance was king, until we unplugged computers

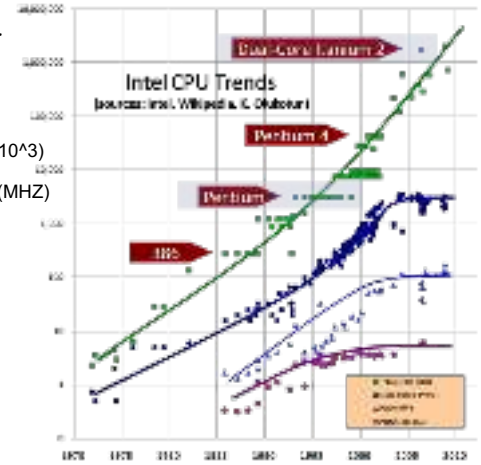
- A lot of “classic” architecture research is makes sure graphs continue to go up and to the right



Processor Design Trends

- Transistors ($\times 10^3$)
- Clock Speed (MHZ)
- Power (W)
- ILP (IPC)

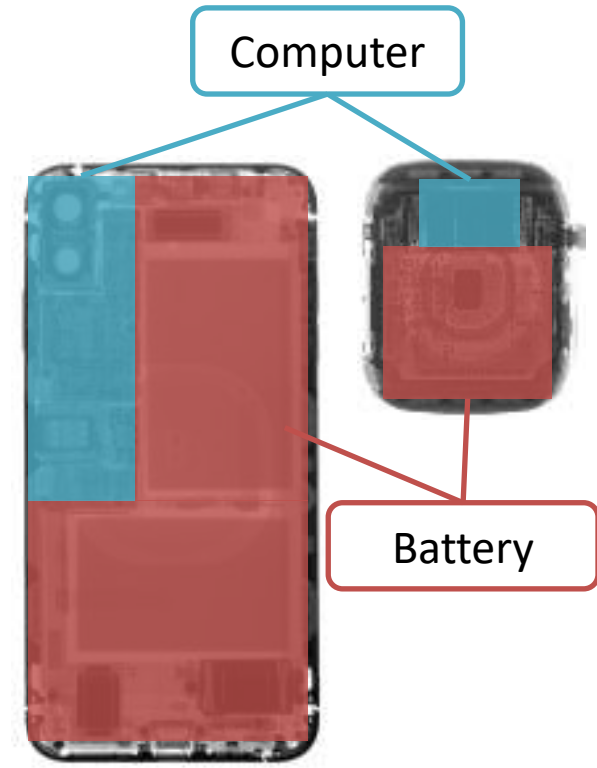
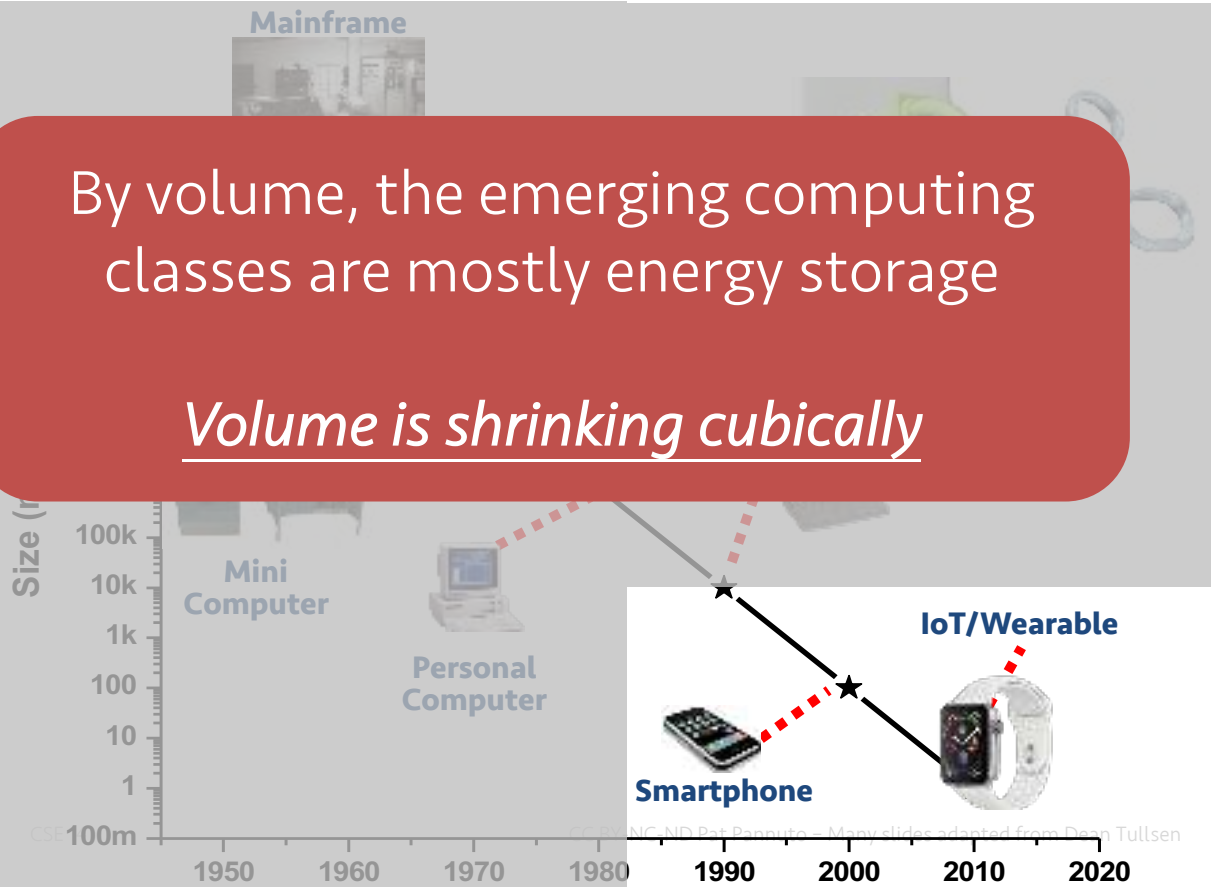
*From Herb Sutter, Dr. Dobbs Journal



I spend my time on graphs that go down and to the right

By volume, the emerging computing classes are mostly energy storage

Volume is shrinking cubically

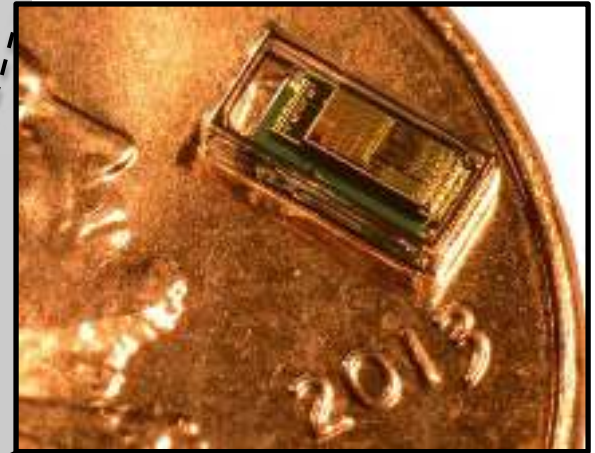


Computational platforms will continue to scale

The next generation of computing will only be a cubic millimeter in size

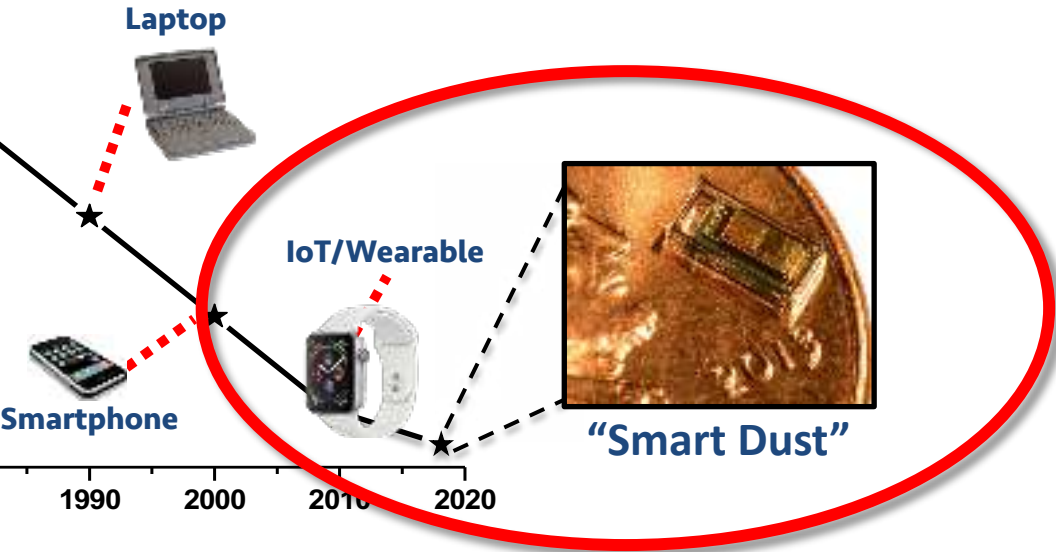
Millimeter-scale batteries have capacities around $5 \mu\text{Ah}$

(would power an idle iPhone for 0.6 s)



“Smart Dust”

Energy constraints will play a central role in the evolution of computing platforms



How must traditional paradigms change, adapt, or re-invent for the new computing classes?

One of the first challenges was re-thinking how we put together computers



Temperature Sensor
~10 pW standby, < 1 μ W active



CPU
~1 nW standby, ~5 μ W active

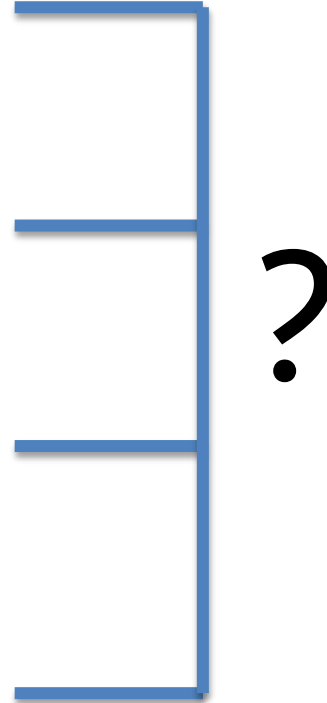


Radio
~10 pW standby, ~10 μ W active

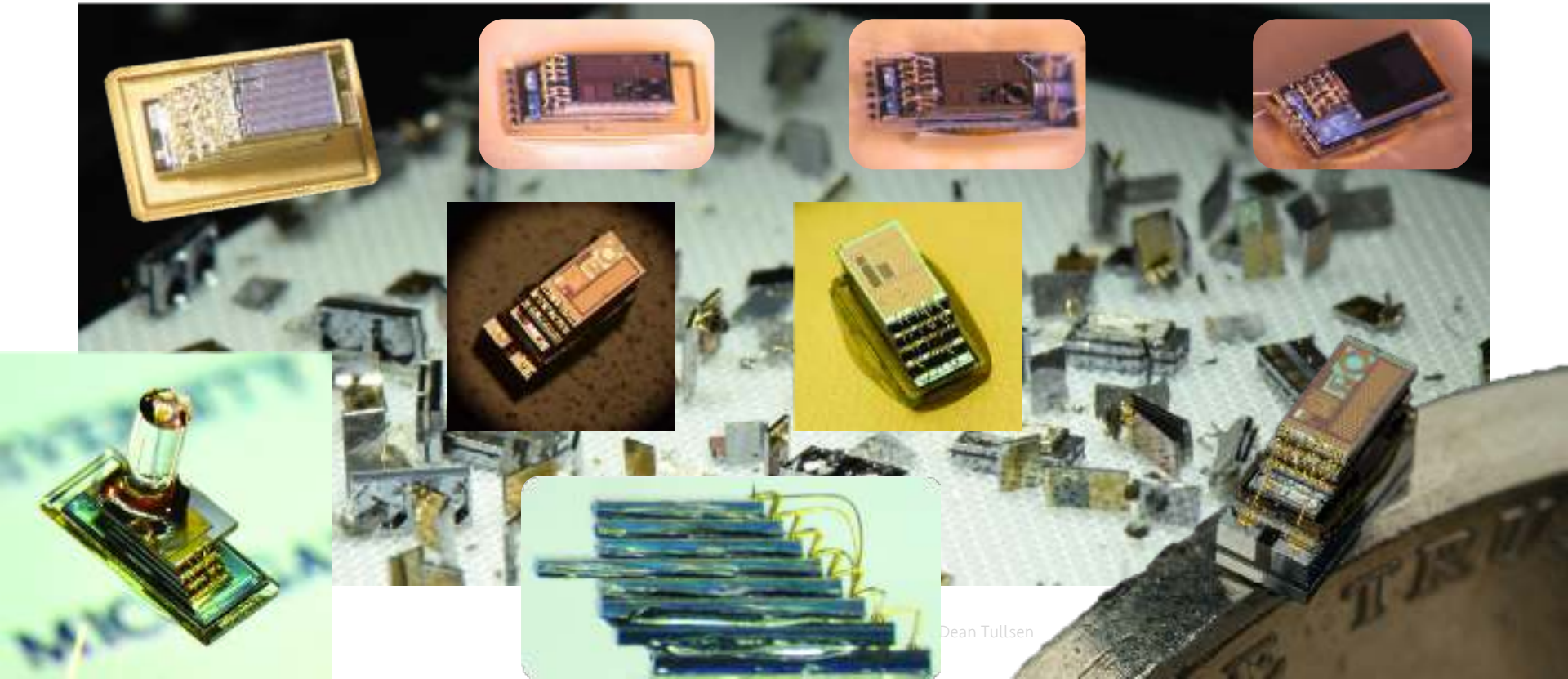


Energy Harvesting & Storage
1~10 nW indoors
2~10 μ Ah capacity

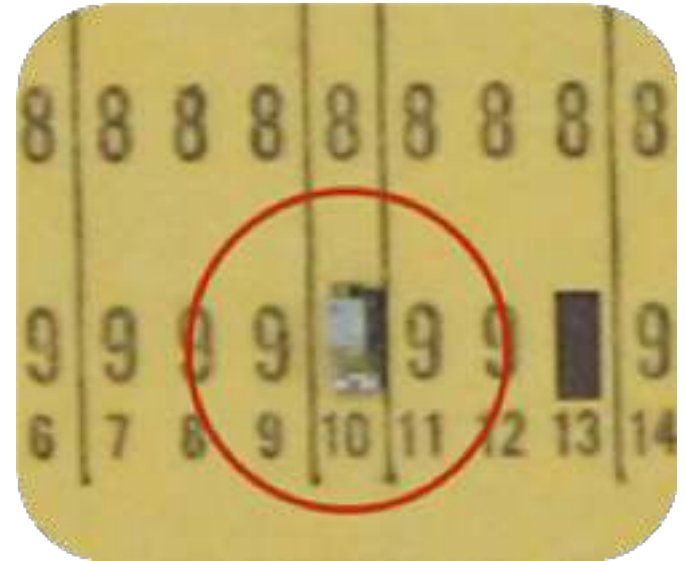
many slides adapted from Dean Tullsen



MBus enabled the development of dozens of millimeter-scale motes as part of the Michigan Micro Mote (M3) project



Check out the “World’s Smallest Computer” exhibit at Silicon Valley’s Computer History Museum!



Next week: Instruction Set Architectures (ISAs)

- For Monday:
 - Skim 1.1 [7 pages]
 - Read 1.2, 1.3 [6.5 pages]

