# Application Specific Architectures

## Introduction and Motivation

Todd Austin

EECS 573
Fall 2016
University of Michigan
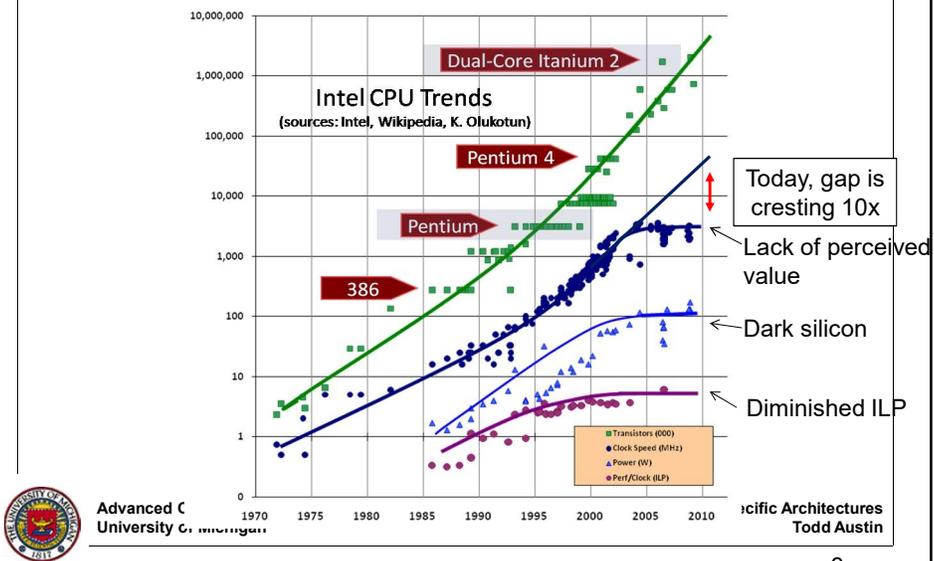
---

# Architecture's Diminishing Return

- Staples of value we strive for…
  - High Speed
  - Low Power
  - Low Cost
- Tricks of the trade
  - Faster clock rates, via pipelining
  - Higher instruction throughput, via ILP extraction
  - Homogeneous parallel systems
- Strong evidence of diminishing return, PIII vs. P4
  - PIII vs. P4: 22% less P4 throughput (0.35 vs. 0.45 SPECInt/MHz)
  - Parallel resources not fully harnessed by today's software
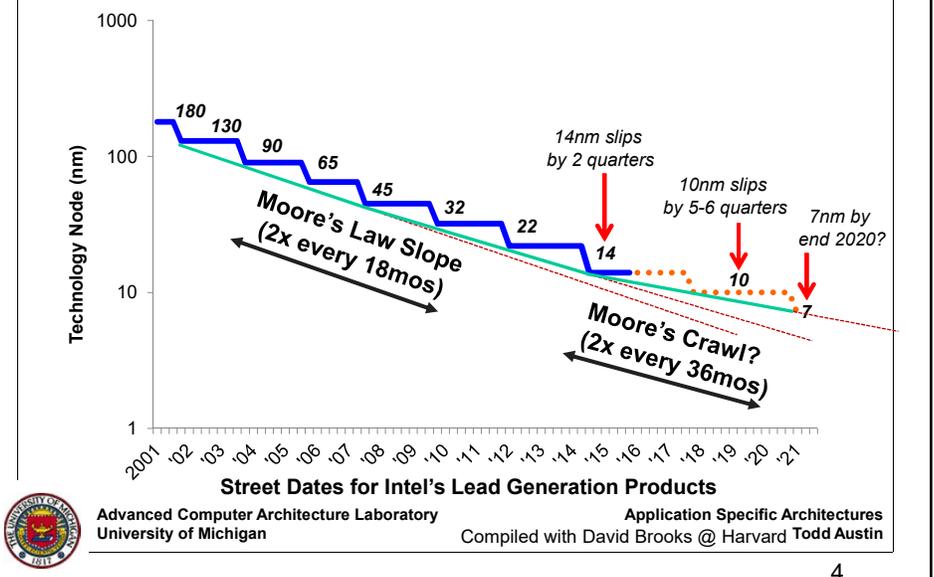- Less return $\Rightarrow$ less value $\Rightarrow$ ☹

## Moore's Law Performance Gap

Intel CPU Trends
(sources: Intel, Wikipedia, K. Olukotun)

Dual-Core Itanium 2

Pentium 4

Pentium

386

Today, gap is cresting 10x

Lack of perceived value

Dark silicon

Diminished ILP

■ Transistors (000)
● Clock Speed (MHz)
▲ Power (W)
● Perf/Clock (ILP)

1970  1975  1980  1985  1990  1995  2000  2005  2010

---

## Is Density Still Scaling?

Technology Node (nm)

180  130  90  65  45  32  22  14  10  7

Moore's Law Slope
(2x every 18mos)

Moore's Crawl?
(2x every 36mos)

14nm slips by 2 quarters

10nm slips by 5-6 quarters

7nm by end 2020?

2001  '02  '03  '04  '05  '06  '07  '08  '09  '10  '11  '12  '13  '14  '15  '16  '17  '18  '19  '20  '21

**Street Dates for Intel's Lead Generation Products**

**Advanced Computer Architecture Laboratory**          **Application Specific Architectures**
**University of Michigan**          Compiled with David Brooks @ Harvard **Todd Austin**

4

# Performance Demands Continue to Grow: Speech Recognition

**Words per Minute** (y-axis: 0, 50, 100, 150, 200, 250)

- Excited Speech (~250 words per minute)
- Unexcited Speech (~150 words per minute)
- Lifetime on 1AA battery

Bar chart data:
- SA-1110 - 206Mhz: 6 hrs
- Xscale - 400Mhz: 2 hrs
- PIII - 600Mhz: 14 min
- PIII - 900Mhz: 7 min
- PIII - 1Ghz: 6 min

**Processor Type** (x-axis)

# Remedy #1: Chip Multiprocessors

Memory Controller

Misc IO

Core Core Queue Core Core

PCIe

Shared L3 Cache

The Dark Silicon Dilemma

Advanced Scaling:
   Dennard: "Computing Capabilities Scale by $S^3$ = 2.8x"
   If S=1.4x …

S = 1.4x
Faster Transistors

$S^2$ = 2x
More Transistors
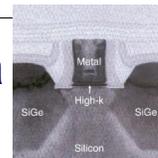
$S^3$
$S^2$
S
1
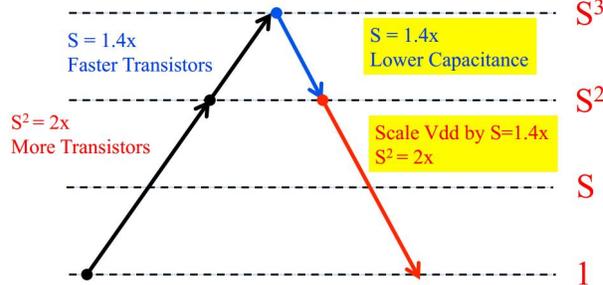
University of Michigan          Courtesy Michael Taylor @ UCSD          Specific Architectures          Todd Austin

7



The Dark Silicon Dilemma

Dennard:
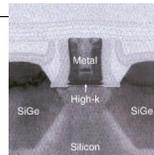   "We can keep power consumption constant"

S = 1.4x
Faster Transistors

S = 1.4x
Lower Capacitance

$S^2$ = 2x
More Transistors

Scale Vdd by S=1.4x
$S^2$ = 2x

$S^3$
$S^2$
S
1

University of Michigan          Courtesy Michael Taylor @ UCSD          Specific Architectures          Todd Austin
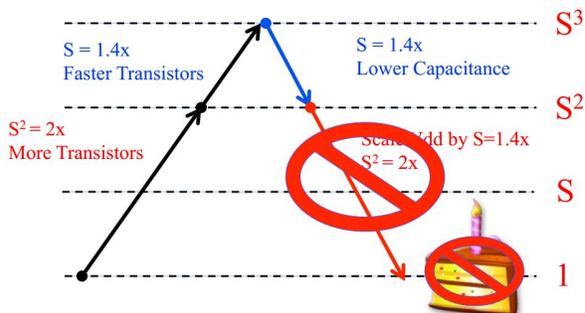
8

4

The Dark Silicon Dilemma

**Fast forward to 2005:**
*Threshold Scaling Problems due to*
*Leakage Prevents Us From Scaling Voltage*

$S = 1.4x$
Faster Transistors

$S = 1.4x$
Lower Capacitance

$S^2 = 2x$
More Transistors

Scale Vdd by S=1.4x
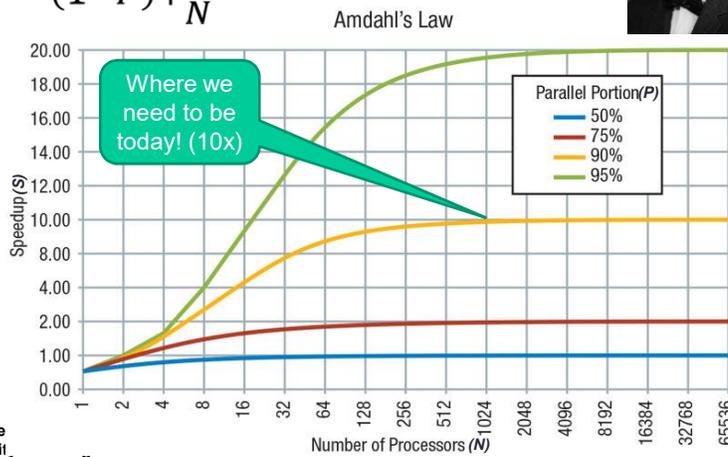$S^2 = 2x$

$S^3$

$S^2$

$S$

$1$

Advanced Computer Architecture Laboratory
University of Michigan

Courtesy Michael Taylor @ UCSD
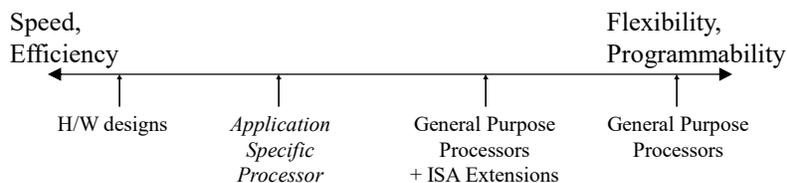
Application Specific Architectures
Todd Austin

9



The Tyranny of Amdahl's Law

$$S(N) = \frac{1}{(1-P)+\frac{P}{N}}$$

Amdahl's Law

Where we need to be today! (10x)

Parallel Portion (*P*)
— 50%
— 75%
— 90%
— 95%

Speedup (*S*)

Number of Processors (*N*)

Advance
Universit

10

5

# A Powerful Solution: Eschew Generality

Speed, Efficiency ← → Flexibility, Programmability

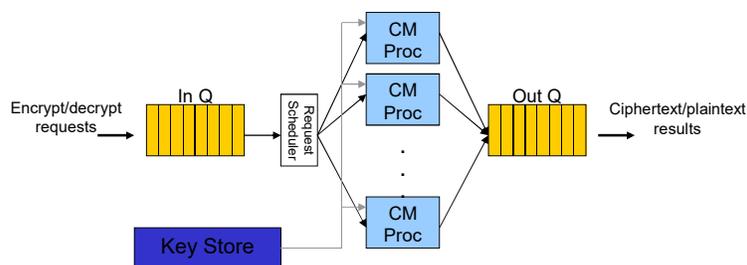| H/W designs | *Application Specific Processor* | General Purpose Processors + ISA Extensions | General Purpose Processors |

- Specialization limits the scope of a device's operation
  - Produces stronger properties and invariants
  - Results in higher return optimizations
  - Programmability preserves the flexibility regarded by GPP's
- A natural fit for embedded designs
  - Where application domains are more likely restrictive
  - Where cost and power are 1st order concerns
- Overcomes growing silicon/architecture bottlenecks
  - Concentrated computation overcomes dark silicon dilemma
  - Customized acceleration speeds up Amdahl's serial codes

Advanced Computer Architecture Laboratory
University of Michigan

Application Specific Architectures
Todd Austin

---

# First Case Study: CryptoManiac [ISCA'01]

Encrypt/decrypt requests → In Q → Request Scheduler → CM Proc / CM Proc / ... / CM Proc → Out Q → Ciphertext/plaintext results

Key Store

- A highly specialized and efficient crypto-processor design
  - Specialized for performance-sensitive *private-key* cipher algorithms
  - Chip-multiprocessor design extracting precious inter-session parallelism
  - CP processors implement with 4-wide 32-bit VLIW processors
  - Design employs crypto-specific architecture, ISA, compiler, and circuits
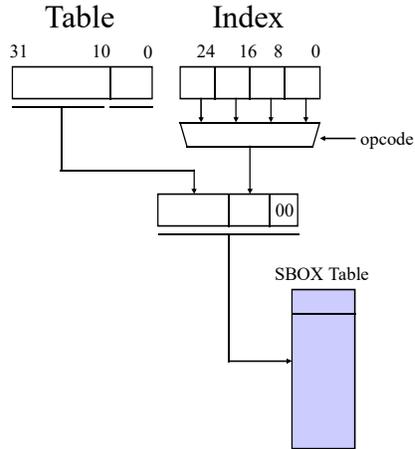
Advanced Computer Architecture Laboratory
University of Michigan

Application Specific Architectures
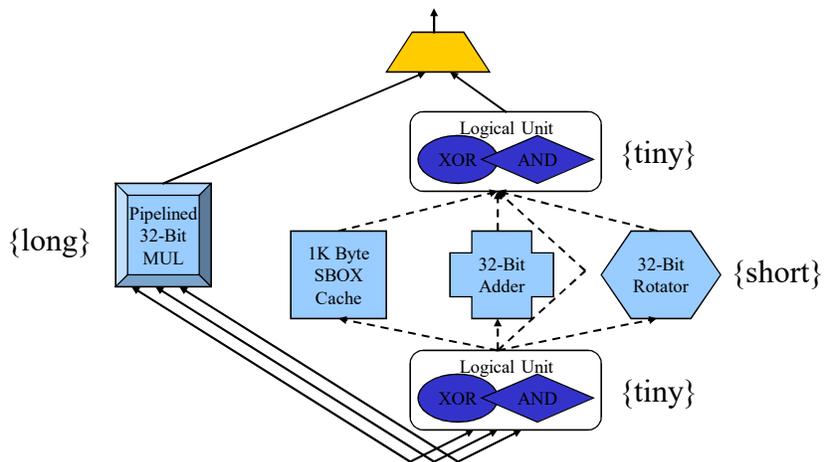Todd Austin

# Crypto-Specific Instructions

- frequent SBOX substitutions
  - X = sbox[(y >> c) & 0xff]
- SBOX instruction
  - Incorporates byte extract
  - Speeds address generation through alignment restrictions
  - 4-cycle Alpha code sequence becomes a single CryptoManiac instruction
- SBOX caches provide a high-bandwidth substitution capability (4 SBOX's/cycle)

Table      Index

31    10    0      24  16   8    0

← opcode

00

SBOX Table

# Crypto-Specific Functional Unit

Logical Unit
XOR   AND      {tiny}

{long}

Pipelined 32-Bit MUL

1K Byte SBOX Cache      32-Bit Adder      32-Bit Rotator      {short}

Logical Unit
XOR   AND      {tiny}

# Crypto-Specific Circuits



- Overclock design until decryption check fails
  - Demonstrated approach with dual SA-1110 IPAQs
- 26% performance increase at room temperature
  - Chill for more improvements, ~10% per 30 degree C

# CryptoManiac Results

- Design implemented in 0.25um physical design flow
  - All components synthesized with Synopsys tools
  - Evaluated with timing analysis and high-level simulation
- Encryption Speed
  - Nearly 1.5x faster than a 600Mhz Alpha 21264 (both 0.25um)
  - 2.25x fast for AES encryption standard
- Design Cost
  - 2 mm$^2$ total area for a single CryptoManiac processor
  - Less than 1/100$^{th}$ the size of an Alpha 21264 (205 mm$^2$)
- Power Characteristics
  - Less than 750 mW total power dissipation
  - Nearly 1/100$^{th}$ the power dissipation of an Alpha 21264 (72 W)

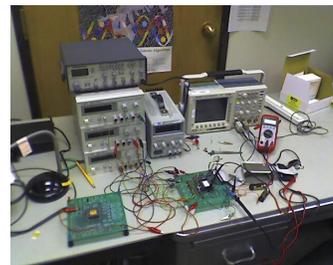## Second Case Study: Subliminal Systems [ISCA'05]

- ◆ Project goals
  - • Explore area-constrained low-energy systems
  - • Develop 100% silicon platforms
  - • Target form factors below 1 mm$^3$



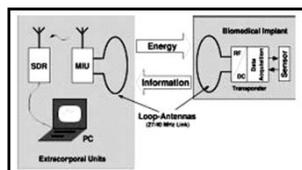| I/O |
| Memory / Sensors |
| CPU / Power |
| I/O |

< 0.5 mm

- ◆ Technology Developments
  - • Subthreshold-voltage processors and memories
  - • Robust subthreshold circuit/cell designs
  - • Compact integrated wireless interfaces
  - • Energy scavenging technologies
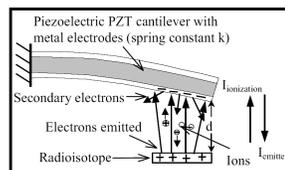  - • Sensor designs

---

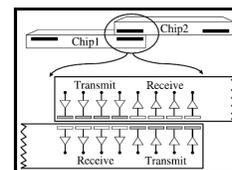## Energy Efficiency: A Key Requirement

- ◆ They live on a limited amount of energy generated from a small battery or scavenged from the environment.

- ◆ Traditionally the communication component is the most power-hungry element of the system. However, new trends are emerging:
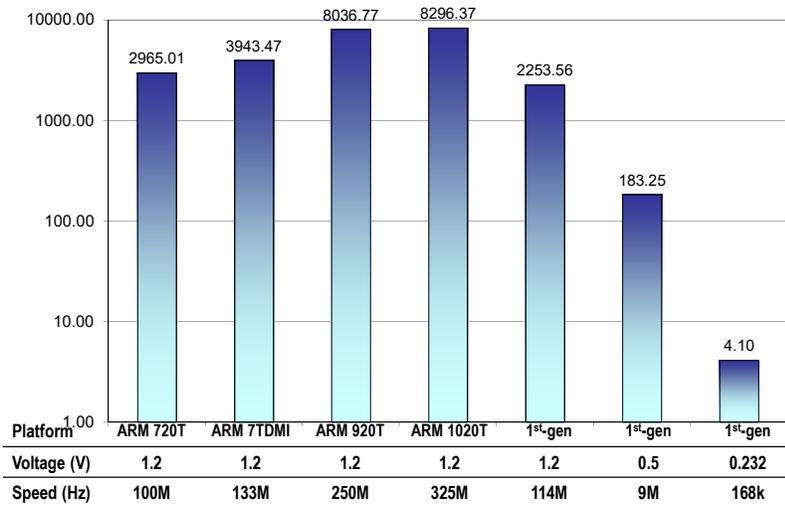


Passive telemetry      Self-powered RF      Proximity comm.

# *Performance of Various Platforms*

| Platform | ARM 720T | ARM 7TDMI | ARM 920T | ARM 1020T | 1st-gen | 1st-gen | 1st-gen |
|---|---|---|---|---|---|---|---|
| Voltage (V) | 1.2 | 1.2 | 1.2 | 1.2 | 1.2 | 0.5 | 0.232 |
| Speed (Hz) | 100M | 133M | 250M | 325M | 114M | 9M | 168k |

Chart values (xRT rating): ARM 720T: 2965.01, ARM 7TDMI: 3943.47, ARM 920T: 8036.77, ARM 1020T: 8296.37, 1st-gen: 2253.56, 1st-gen: 183.25, 1st-gen: 4.10
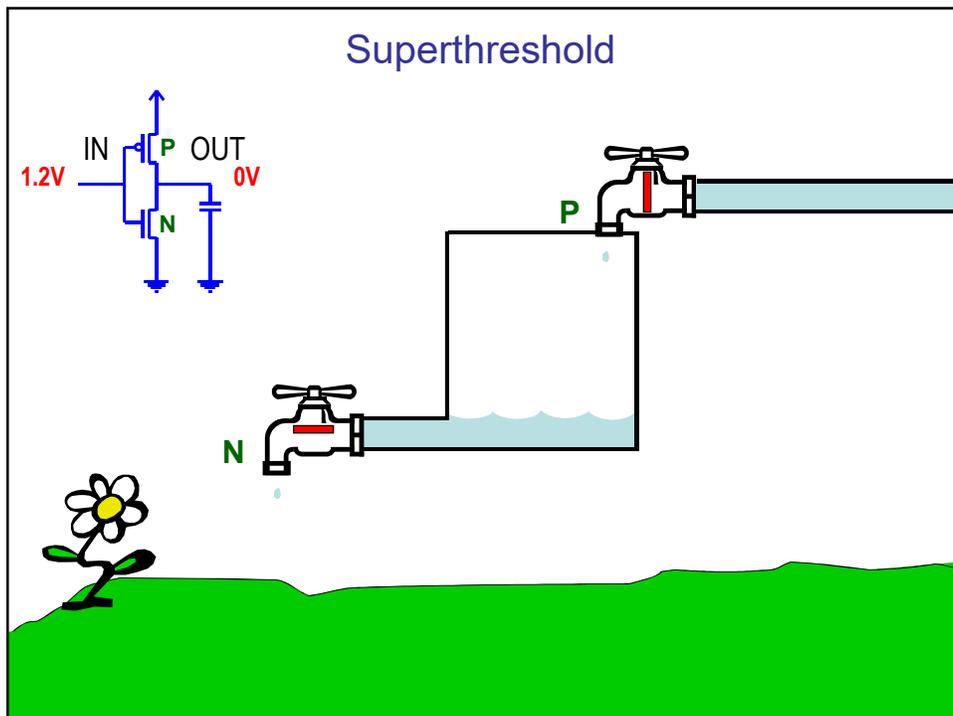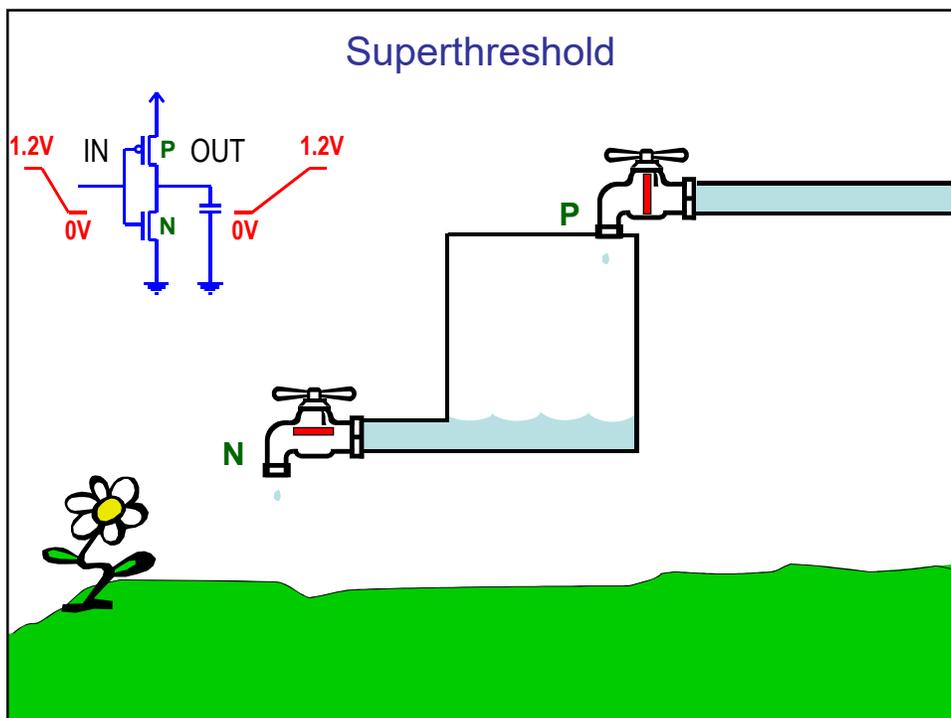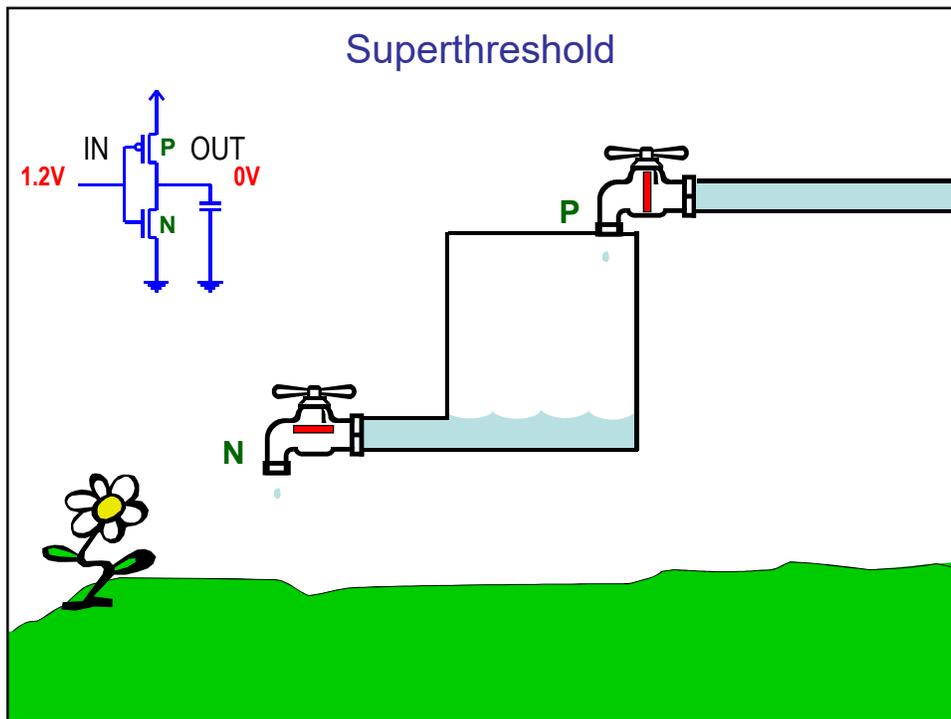
**xRT rating**: how many times faster than real-time the processor can handle the worst-case data stream rate on the most computationally intensive sensor benchmark
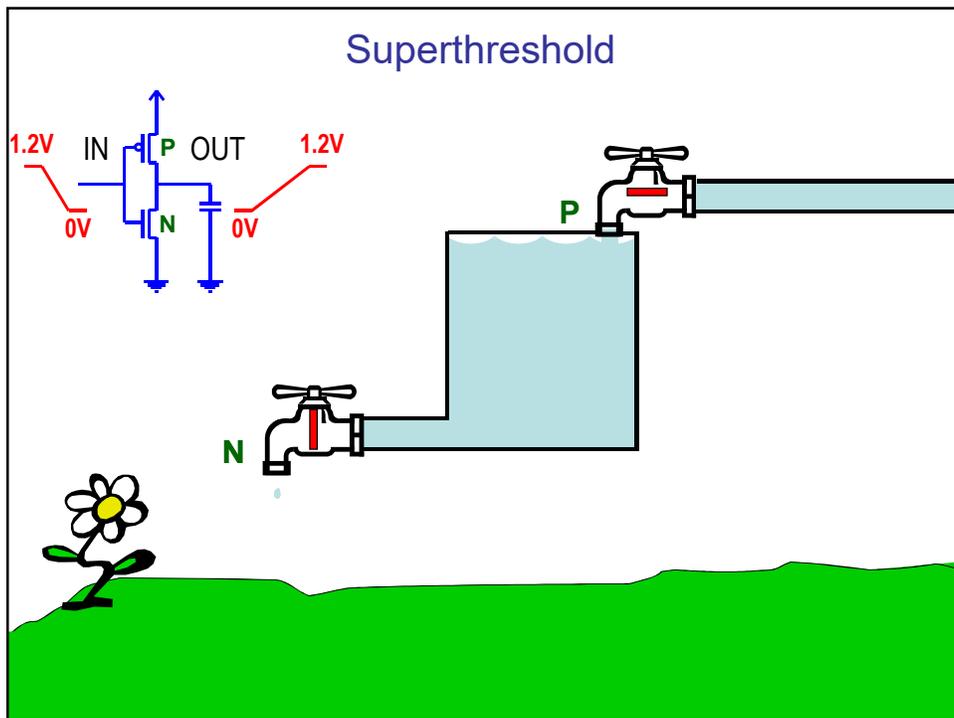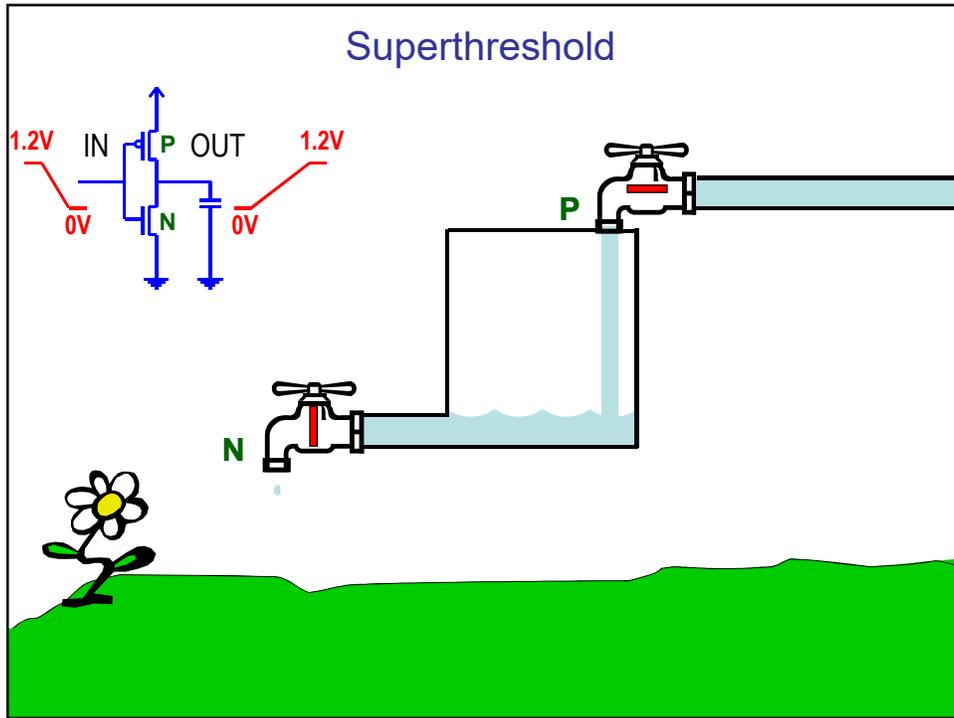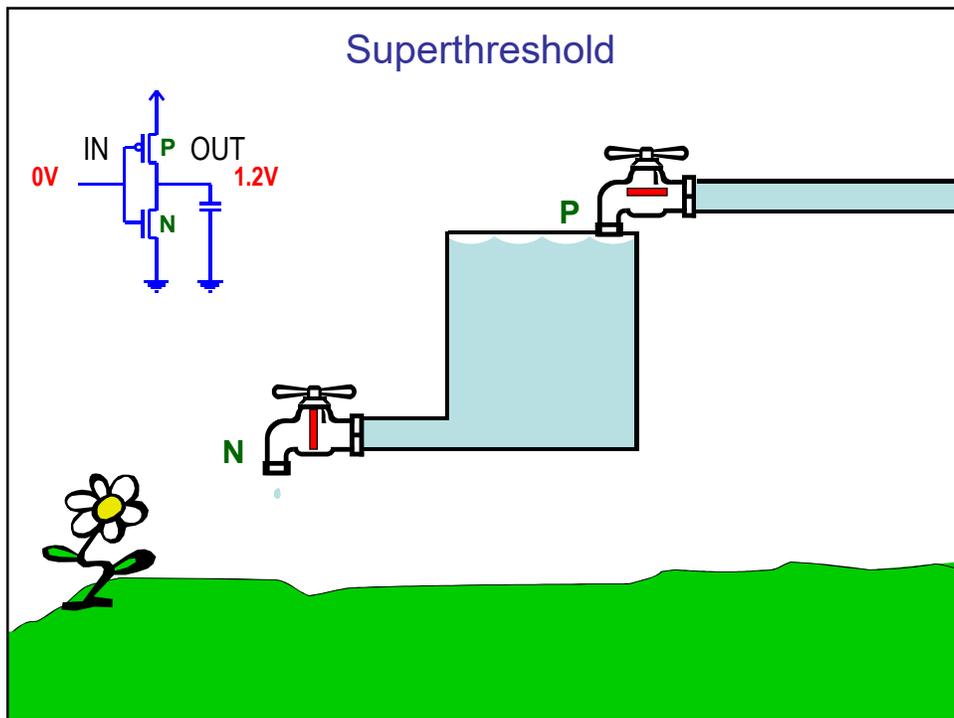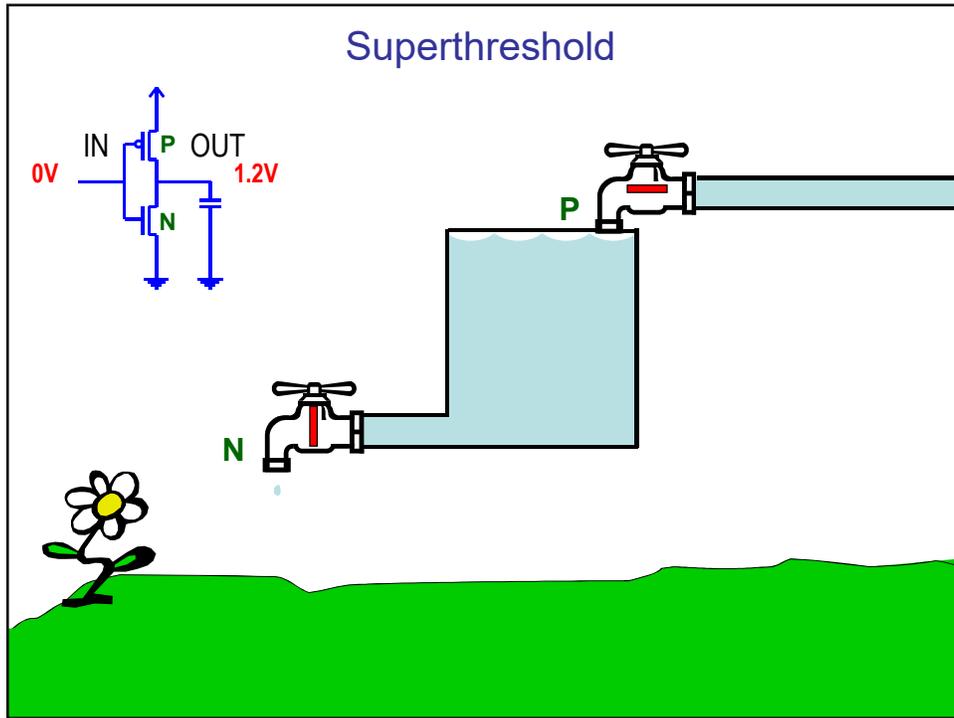
---

# *The Basics of Subthreshold Circuit Operation*

*A Short Animation by Leyla Nazhandali*

*Episode 1: Inverter operation in superthreshold domain*



Superthreshold

IN  P  OUT
1.2V        0V
N

P

N

Superthreshold



Superthreshold

Superthreshold



Superthreshold

13

Superthreshold



Superthreshold

# Episode 2: Inverter operation in subthreshold domain

## Subthreshold

IN P OUT
0.2V 0V
N

P

N

## Subthreshold

IN   P   OUT
0.2V ─────  0V

P

N

## Subthreshold

0.2V   IN   P   OUT        0.2V

0V            N   0V

P

N

16

Subthreshold



Subthreshold

Subthreshold



Subthreshold

18

Subthreshold


Subthreshold

## Subthreshold



## *Summary from Architecture Study*

- We studied 21 different processors experimenting with following options:
  - Number of stages
  - w/ vs. w/o instruction prefetch buffer
  - w/ vs. w/o explicit register file
  - Harvard vs. Von-Neumann architecture
- To minimize energy at subthreshold voltages, architects must:

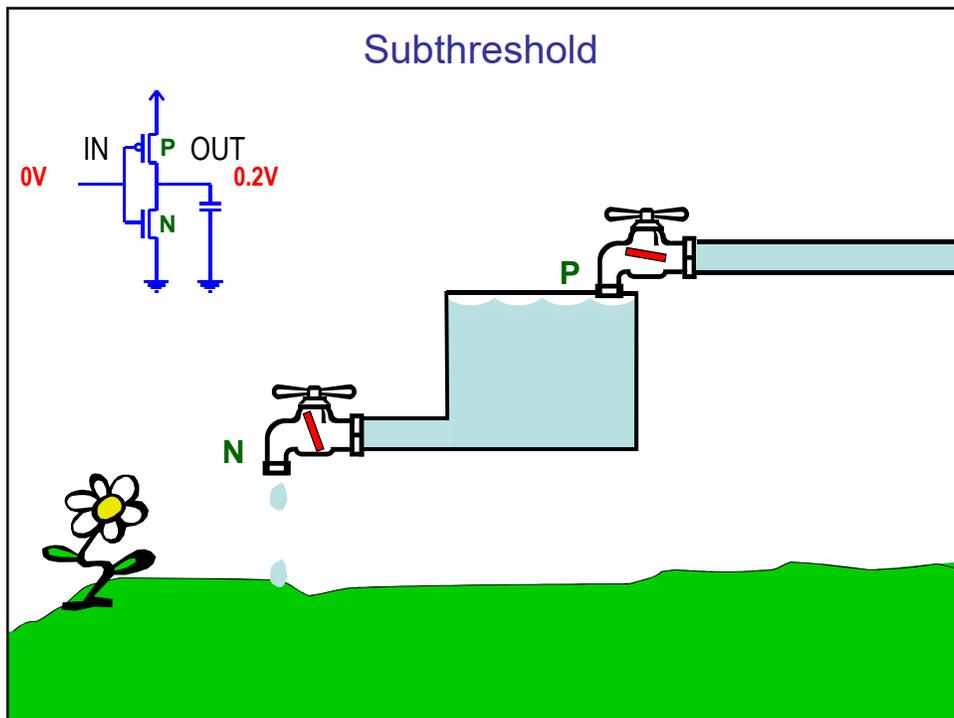| | | |
|---|---|---|
| **Minimize area** | ⇨ | **To reduce leakage energy per cycle** |
| **Maximize Transistor utility** | ⇨ | **To reduce $V_{min}$ and energy per cycle** |
| **Minimize CPI** | ⇨ | **To reduce Energy per instruction** |

- The memory comprises the single largest factor of leakage energy, as such, efficient designs must reduce memory storage requirements.

## Microarchitecture Overview



## First Subliminal Chip

# *Pareto Analysis for Several Processors*



**Energy (J/inst.)** (y-axis), values: 3.00E-12, 2.80E-12, 2.60E-12, 2.40E-12, 2.20E-12, 2.00E-12, 1.80E-12, 1.60E-12, 1.40E-12

**Inst Latency (1/perf == s/inst.)** (x-axis), values: 5.00E-06, 1.00E-05, 1.50E-05, 2.00E-05, 2.50E-05, 3.00E-05, 3.50(, 00E-05

Labels within plot:
- 2s_h_08w_r → w/ explicit register file
- 2s_v_08w_r
- 2s_v_32w
- 2.33 / 4.39
- # of stages = 3
- 3s_h_16w_r
- 2.66 / 3.59
- 3s_h_08w_r
- architecture: Von Neumann (vs. Harvard)
- 2s_h_32w_r
- 2s_h_16w_r
- 3s_h_32w_r
- 2s_v_16w
- 3s_v_16w
- 1.77 / 5.17
- Better
- 3s_h_32w
- 3s_h_08w
- 3s_v_08w
- ALU width
- 2s_h_32w
- 3s_h_16w
- Implemented design
- 3s_v_08w
- Area = 2.14
- CPI = 2.88
- 2s_h_16w
- 2s_h_08w
- 1.78 / 3.62
- 1.37 / 4.99
- 1.10 / 6.14

---

# *Third Case Study: Taking Computer Vision Mobile*

- ◆ Embedded mobile computation on the rise
  - • Smart Phones, Tablets
  - • Improved sensors
    - ♦ High megapixel cameras, HD video
  - • New capabilities from new sensors



  - • There is a need for near real time computation
    - ♦ Users don't want to wait
- ◆ Why not use the cloud?
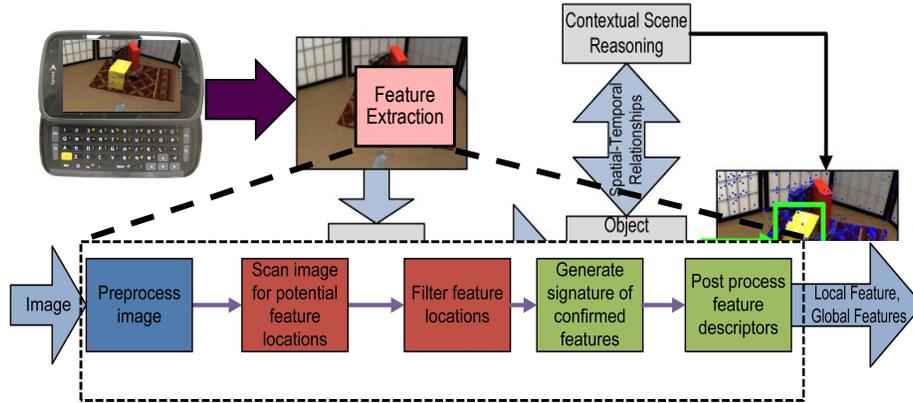  - • High latency
  - • Bandwidth Limits
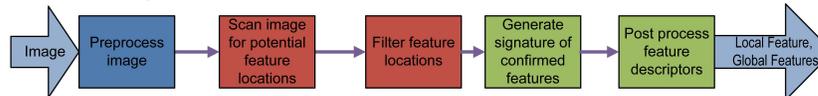  - • Reliability

# Computer Vision

**Typical computer vision pipeline**



---

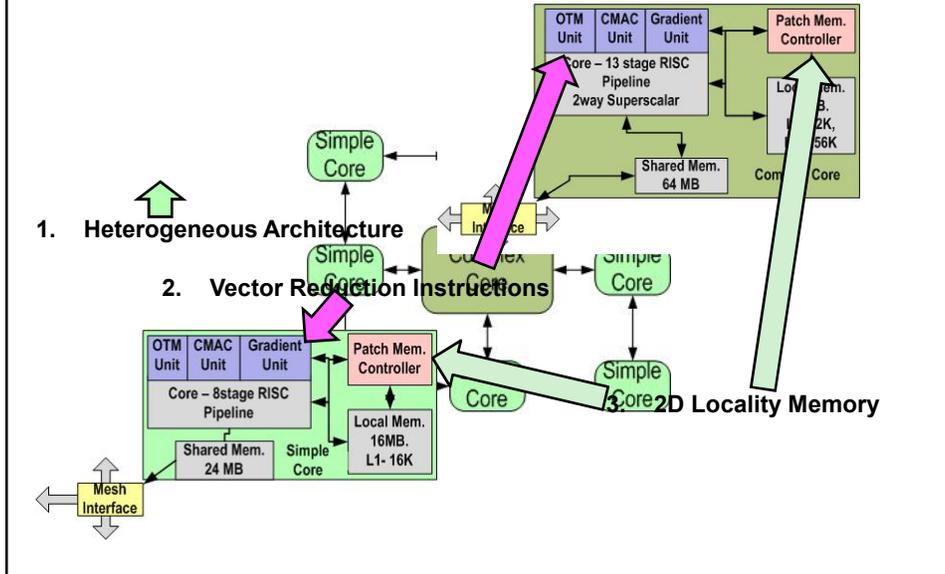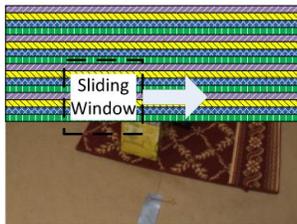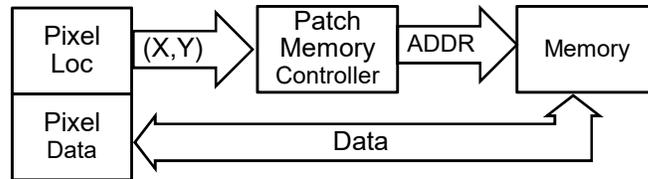# Feature Extraction Characteristics

- 3 Algorithms
  - FAST – corner detection
  - HoG – general object shape detector
  - SIFT – specific object/blob detector
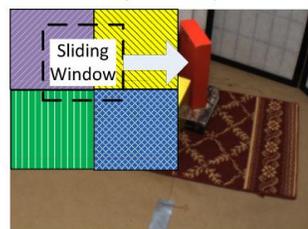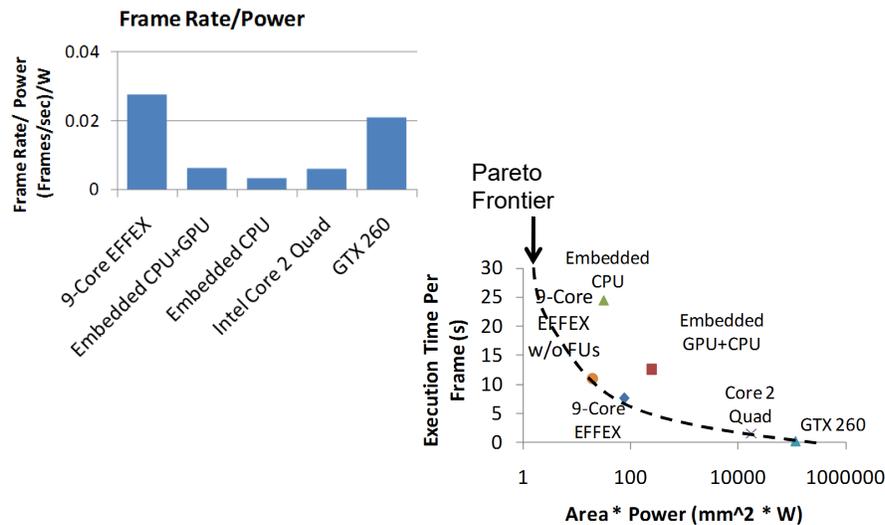
# Efficient Fast Feature EXtraction

1. **Heterogeneous Architecture**

    2. **Vector Reduction Instructions**

        3. **2D Locality Memory**

# Patch Memory



Traditional image storage

Patch memory storage

## A Taste of the Results

**Frame Rate/Power**

(bar chart: Frame Rate/ Power (Frames/sec)/W — y-axis 0, 0.02, 0.04)

Categories: 9-Core EFFEX, Embedded CPU+GPU, Embedded CPU, Intel Core 2 Quad, GTX 260

(scatter plot)
- Y-axis: Execution Time Per Frame (s) — 0, 5, 10, 15, 20, 25, 30
- X-axis: Area * Power (mm^2 * W) — 1, 100, 10000, 1000000

Pareto Frontier

Labels: Embedded CPU, 9-Core EFFEX w/o FUs, Embedded GPU+CPU, 9-Core EFFEX, Core 2 Quad, GTX 260

---

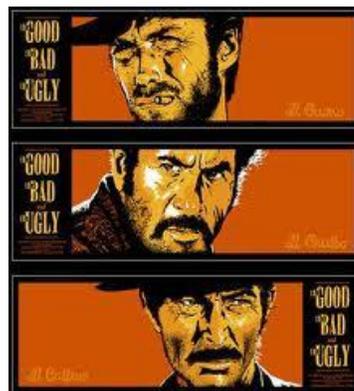## Outlook for App-Specific Design is Unsure: The Good, the Bad and the Ugly

- ◆ **The Good**: Moore's law will continue for the near future
  - It won't last forever, but that another problem

- ◆ **The Bad**: Dennard scaling has all but stopped, leaving innovation to fill the performance/power scaling gap
  - E.g., app-specific design, custom accelerators

- ◆ **The Ugly**: Hardware innovation requires design diversity, which is ultimately *too expensive to afford*
  - Skyrocketing NREs will necessitate broadly applicable (vanilla and slow) H/W designs

50

**Design Costs Are Skyrocketing**

*Source: International Business Strategies*

51



**High Costs Will be a Showstopper**

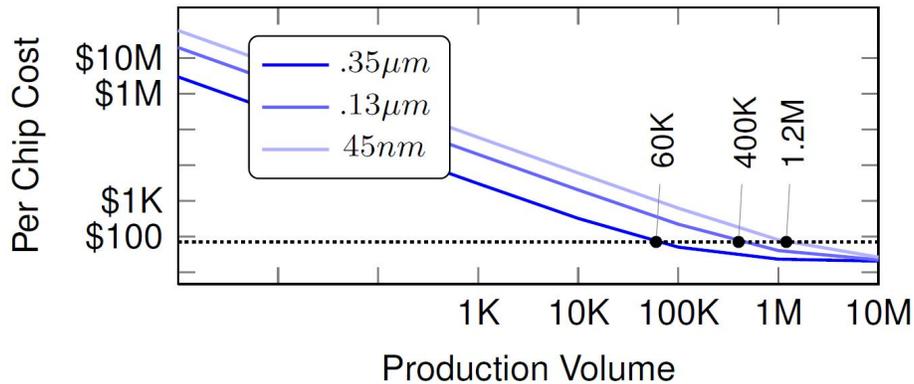◆ Heterogeneous designs often serve smaller markets

52

## Outcome: "Nanodiversity" is Dwindling



Expensive development costs demand BIGGER markets, this trend works against customized designs.

*Source: Gartner Group*

53

---

## The Remedy: Scale Innovation

◆ Ultimate goal: ***accelerate system architecture innovation*** and make it efficient and inexpensive enough that ***anyone can do it anywhere***

  • Approach #1: Embrace system-level innovation

  • Approach #2: Leverage technology advances on CMOS silicon

  • Approach #3: Reduce the cost to design custom hardware

  • Approach #4: Widen the applicability of custom hardware

  • Approach #5: Reduce the cost of manufacturing custom H/W

54

# 1) Embrace system-level innovation



"Give me 15% speedup and I'll accept your paper"
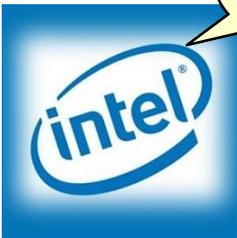
"I need 1% speedup for 1% area"

"Your **system-level** ideas needs to **deliver 2x or more**, or someone else should fund it"

55

---

# HELIX-UP Unleashed Parallelization

David Brooks @ Harvard

- ◆ Traditional parallelizing compilers must honor *possible* dependencies

- ◆ HELIX-UP manufactures parallelism by profiling which deps do not exist and *which are not needed*
  - ▪ Based on user supplied *output distortion function*

- ◆ Big step for parallelization
  - ▪ **2x speedup** over parallelizing compilers, 6x over serial, < 7% distortion

Thread 0 — Iteration 0
Thread 1 — Data — Iteration 1
Thread 2 — Data
Thread 3 — Data



**Nehalem 6 cores, 2 threads per core**

Output distortion 2.1%, 3%, 6.9%, 3.8%, 0%, 2%

56

28

## Association Rule Mining with the Automata Processor

Kevin Skadron @ UVA

- Micron's Automata processor
  - Implements FSMs at memory
  - Massively parallel with accelerators

- Mapped data-mining ARM rules to memory-based FSMs
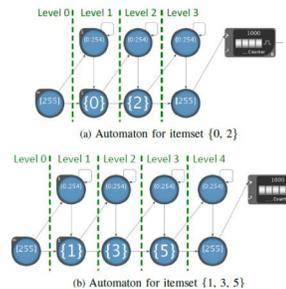  - ARM algorithms identify relationships between data elements
  - Implementations are often memory bottlenecked

- Big-data sets had big speedups
  - 90x+ over single CPU performance
  - **2-9x+ speedups** over CMPs and GPUs

- Joint effort with UVA and Micron



(a) Automaton for itemset {0, 2}

(b) Automaton for itemset {1, 3, 5}

57

---

## 2) Leverage technology advances on CMOS silicon

- Recent success: the reduced leakage and transient fault protection of FinFETs

- Upcoming: the density and durability of Intel/Micron's XPoint memory technology

- Many additional opportunities possible: TFETs, CNTs, spin-tronics, novel materials, analog accelerators, etc…

- Key challenge: integration of non-silicon technologies

- Advice: to maximize benefits of these devices, architects need to work with device and materials researchers



58

## Top 10 Technology Plays that Would Make Architects REALLY Excited

- Reduced leakage for memory
  - Helps with low power sleep states, allows lower computational power states
- Reduced leakage for computation
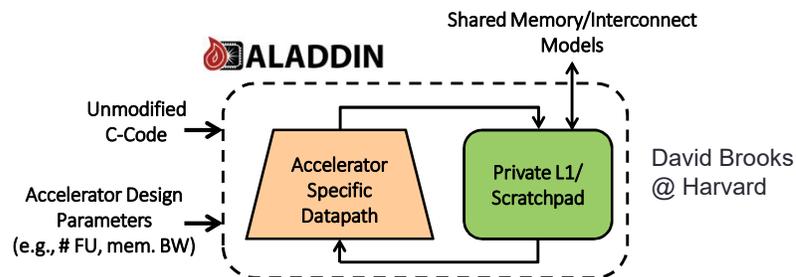  - Re-balances the power-parallelism tradeoff in favor of more performance/watt
- Controllable and recognizable analog functions
  - Allow computation to be replaced with potentially fast and efficient analog compute
- Ultra-cheap fabrication technologies
  - Re-balances the specialization-cost tradeoffs, making system-level optimization more valuable
- Emerging technologies that deliver additional traditional value at low fault rates
  - We have many low-cost system-level fault tolerance technologies, let's use them!, limit faults to < 0.1%
- Emerging technologies that are not too fiddly, unless they deliver *significant* value
  - We need clean productive abstractions, CMOS is the benchmark, compare to asynch and CUDA
- Faster, more energy efficient, less destructive writes for nonvolatile storage
  - Allows for simpler, denser, more efficient memory designs, supports ultra-low power states
- Computation/memory capabilities with no power/electrical/*etc.* signature
  - Today's systems are fraught with side channels, this is needed as a basis for establishing H/W trust
- More energy efficient communication that doesn't overtly exacerbate latency
  - Allows for more system scalability – both scale-in and scale-out
- More energy efficient computation that is dense and cheap
  - Allows for more T-flops, since almost all computational capabilities today are energy bounded

59

---

## 3) Reduce the cost to design custom hardware



- Better tools and infrastructure
  - Scalable accelerator synthesis and compilation, *generate code and H/W for highly reusable accelerators*
  - Composable design space exploration, *enables efficient exploration of highly complex design spaces*
  - Well put-together benchmark suites to drive development efforts

60

## CortexSuite:
## A Synthetic Brain Benchmark Suite

Michael Taylor @ UCSD



61

## Embrace Open-Source Concepts to Reduce Costs



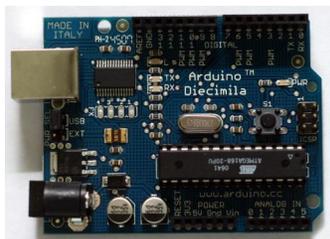**Red** = non-free IP, **Green** = free IP

62

## Embrace Open-Source Concepts to Reduce Costs



Red = non-free IP, Green = free IP

---

## *Open-Source H/W is Growing*

## 4) Widen the Applicability of Customized H/W

Krste Asanovic @ UC-Berkeley



**Applications**: Computer Vision, Multimedia Analysis, Machine Learning

**Computational Patterns**: Dense, Sparse, ... Graph

**Specializers with custom implementations and autotuning**

ESP Code: Glue Code, Dense Code, Sparse Code, Graph Code

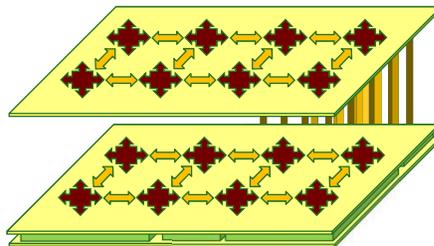ESP Core: ILP Engine, Dense Engine, Sparse Engine, Graph Engine

- ◆ ESP: Ensembles of Specialized Processors
  - ◆ Ensembles are algorithmic-specific processors optimized for code "patterns"
  - ◆ Approach uses *composable customization* to deliver speed and efficiency that is widely applicable to general purpose programs
  - ◆ Grand challenges remain: *what are the components* and *how are they connected*?

65

## 5) Reduce the cost of manufacturing customized H/W

Martha Kim @ Columbia

- ◆ A breakthrough experiment: what if building custom silicon were like fabricating a chip?

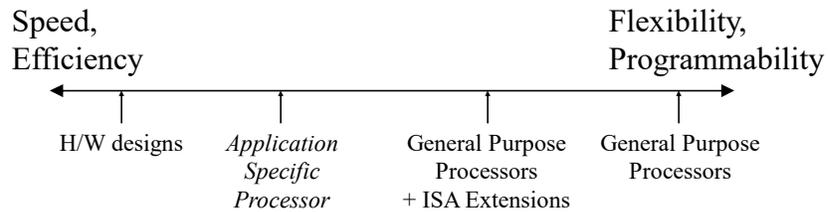- ◆ Rich and near-term research agenda in customization, i.e., MCMs + 3D + eFPGA interfaces



**Brick-and-mortar silicon design flow:**
1) Assemble brick layer
2) Connect with mortar layer
3) Package assembly
4) Deploy software

- ◆ Diversity via brick ecosystem & interconnect flexibility
- ◆ Brick design costs amortized across all designs
- ◆ Robust interconnect and custom bricks rival ASIC speeds
- ◆ Facilitates non-silicon integration and mature design strategies

66

33

# *Summary: Benefits of App-Specific Design*

Speed,
Efficiency

Flexibility,
Programmability

| H/W designs | *Application Specific Processor* | General Purpose Processors + ISA Extensions | General Purpose Processors |
|---|---|---|---|

- ◆ Specialization limits the scope of a device's operation
  - • Produces stronger properties and invariants
  - • Results in higher return optimizations
  - • Programmability preserves the flexibility regarded by GPP's
- ◆ A natural fit for embedded designs
  - • Where application domains are more likely restrictive
  - • Where cost and power are 1st order concerns
- ◆ Overcomes growing silicon/architecture bottlenecks
  - • Concentrated computation overcomes dark silicon dilemma
  - • Customized acceleration speeds up Amdahl's serial codes