| Application |
|---|

| Algorithm |
|---|

| Language |
|---|

Hardware complexity largely hidden via stable abstractions and interfaces

| Compiler |
|---|

| Architecture (I,S,N) |
|---|

| Microarchitecture |
|---|

Continuous improvements in all of the lower layers has created consistently large gains in performance

| Circuits |
|---|

| Devices |
|---|

# The glory of Moore's Law

## The experts look ahead

### Cramming more components onto integrated circuits

With unit cost falling as the number of components per circuit rises, by 1975 economics may dictate squeezing as many as 65,000 components on a single silicon chip

By Gordon E. Moore

Director, Research and Development Laboratories, Fairchild Semiconductor division of Fairchild Camera and Instrument Corp.

**Intel 4004**
2300 transistors
740 kHz clock
10um process
10.8 usec/inst

**Intel Core i7 980X**
1.17B transistors
3.33 GHz clock
32nm process
73.4 psec/inst

**%/year, Ratios:**
38%, 508000
23%, 4450
15%, 312
34%, 147000

# Moore's secret sauce: Dennard scaling

| Device or Circuit Parameter | Scaling Factor |
|---|---|
| Dimension, Tox, L, W | $1/k$ |
| Doping Concentration Na | $k$ |
| Voltage (V) | $1/k$ |
| Current (I) | $1/k$ |
| Capacitance (eA/t ) | $1/k$ |
| Delay time/circuit (VC/I) | $1/k$ |
| Power dissipation/circuit (VI) | $1/k^2$ |
| Power density (VI/A) | $1$ |

Historically, k ~= 1.4

[Dennard, Gaensslen, Yu, Rideout, Bassous, Leblanc, **IEEE JSSC**, 1974]



### Design of Ion-Implanted MOSFET's with Very Small Physical Dimensions

ROBERT H. DENNARD, MEMBER, IEEE, FRITZ H. GAENSSLEN, HWA-NIEN YU, MEMBER, IEEE, V. LEO RIDEOUT, MEMBER, IEEE, ERNEST BASSOUS, AND ANDRE R. LeBLANC, MEMBER, IEEE

2x transistor count

40% faster

50% more efficient

# Dennard scaling is dead



Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten
Dotted line extrapolations by C. Moore

# Multicore to the rescue?



[Esmaeilzadeh, Blem, St. Amant, Sankaralingam, Burger, ISCA 2011]

# A brief history of computer architecture

1940s | 1950s | 1960s | 1970s | 1980s | 1990s | 2000s | 2010s | 2020s

Era of discovery

Era of invention

Era of integration

Era of ILP

Era of multicore

Era of X

X = {Logic specialization, neural computing, cold computing, ?}

# Specialization: A path forward (?)

# More gains the lower you go

| | |
|---|---|
| Code specialization | 10x |
| Logic specialization | 100x |
| Circuit specialization | 1000x |
| Device specialization | 10000x |

# Logic synthesis as a platform

## Large gains in efficiency with direct software-to-logic

100x for FPGA/CGRAs, 1000x for ASICs

## Development and compilation is a huge challenge

Mix of cores, hard IP blocks, and tools to target them (AutoESL, OpenCL)

Map common operations and flows into libraries, compose them

## Will see growing adoption, increased tool investment

Initially FPGAs in the cloud, CGRAs/ASICS in the client

Generality
(CPUs)

Efficiency
(ASICs)

How do we resolve this tension?

# An end to Moore's Law

| High Volume Manufacturing | 2008 | 2010 | 2012 | 2014 | 2016 | 2018 | 2020 | 2022 |
|---|---|---|---|---|---|---|---|---|
| Technology Node (nm) | 45 | 32 | 22 | 16 | 11 | 8 | 6 | 4 |
| Integration Capacity (BT) | 8 | 16 | 32 | 64 | 128 | 256 | 512 | 1024 |

Source: Shekhar Borkar, Intel Corporation

# Approximate computing

## Changing workloads offer an opportunity

Large-scale machine learning

Computer vision

Bioinformatics

Mining big data
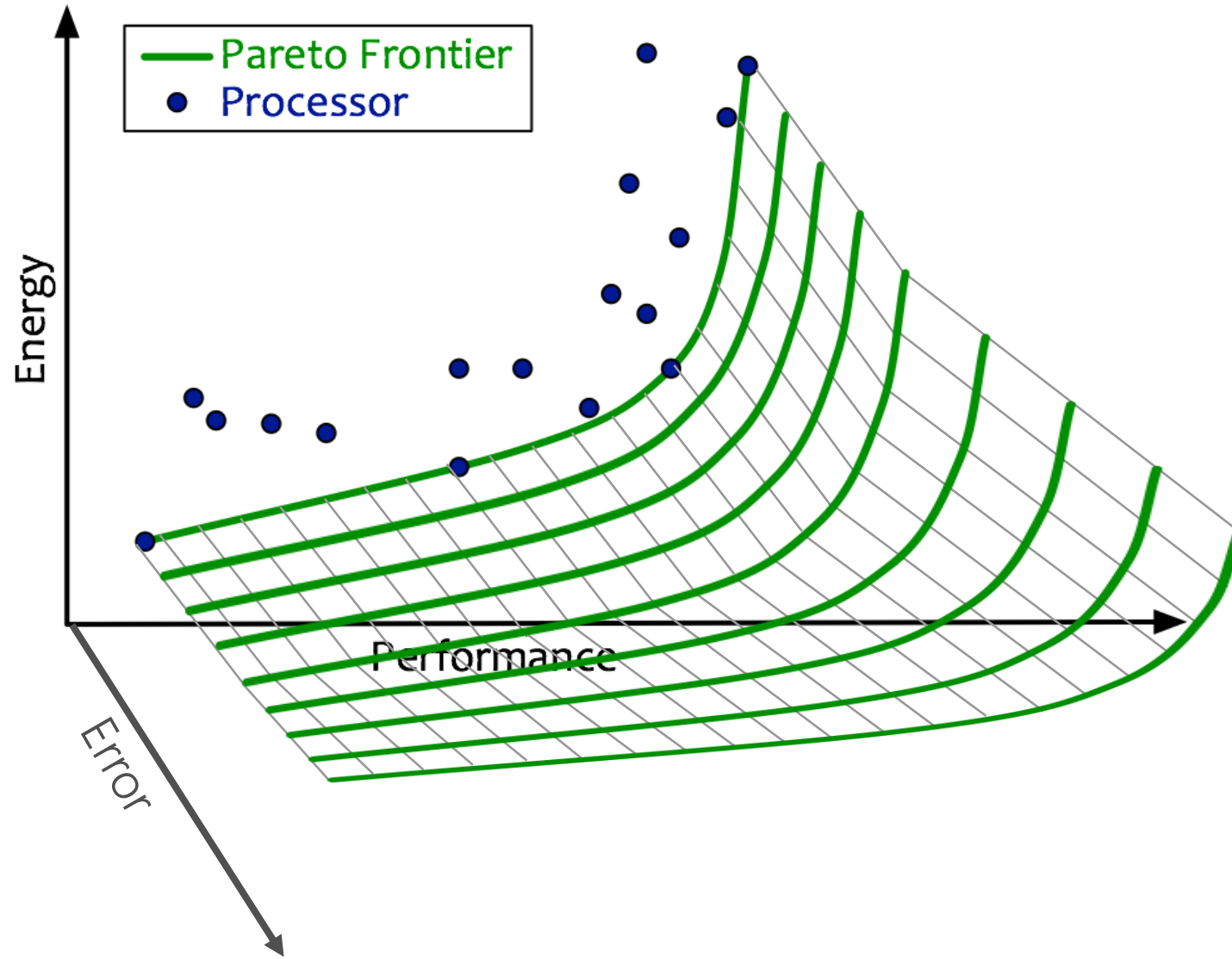
Speech and AI

## Robust to reduced precision

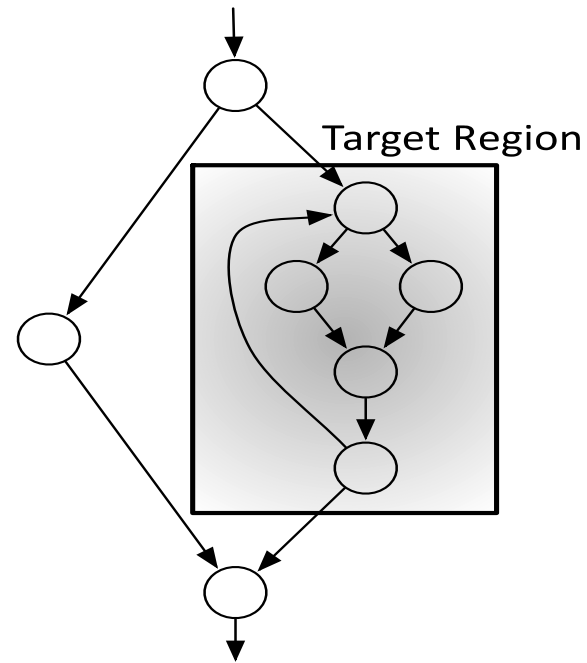Need formal semantics to reason about and bound error

Need to handle dynamic noise and variations

Thanks to Luis Ceze, Hadi Esmaeilzadeh, Adrian Sampson, and others
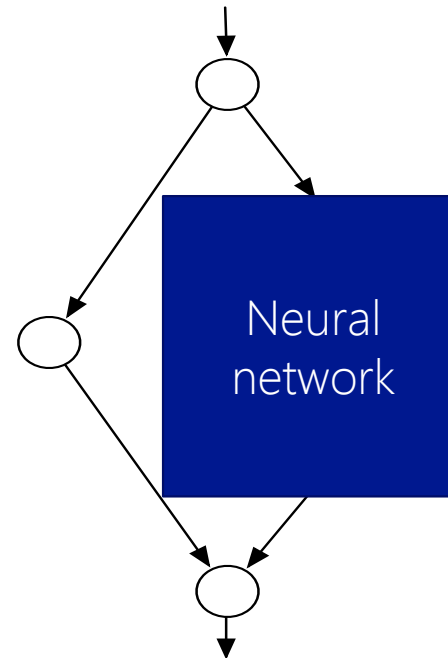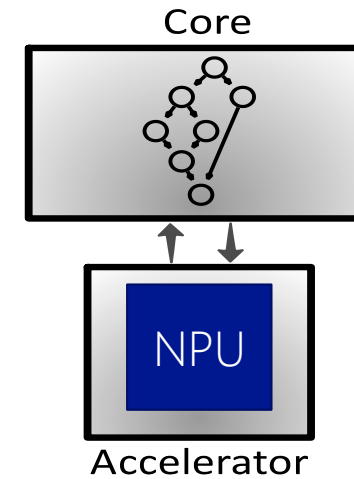
# Approximate computing

# Transforming von Neumann to NNs



Target Region

Core

Neural network

NPU

Accelerator

Imperative code

Transformed code

Accelerated execution

[Esmaeilzadeh, Sampson, Ceze, and Burger, MICRO 2012]

# Applications using Neural Transformation

Signal Processing
  fft
Robotics
  inverse2kj
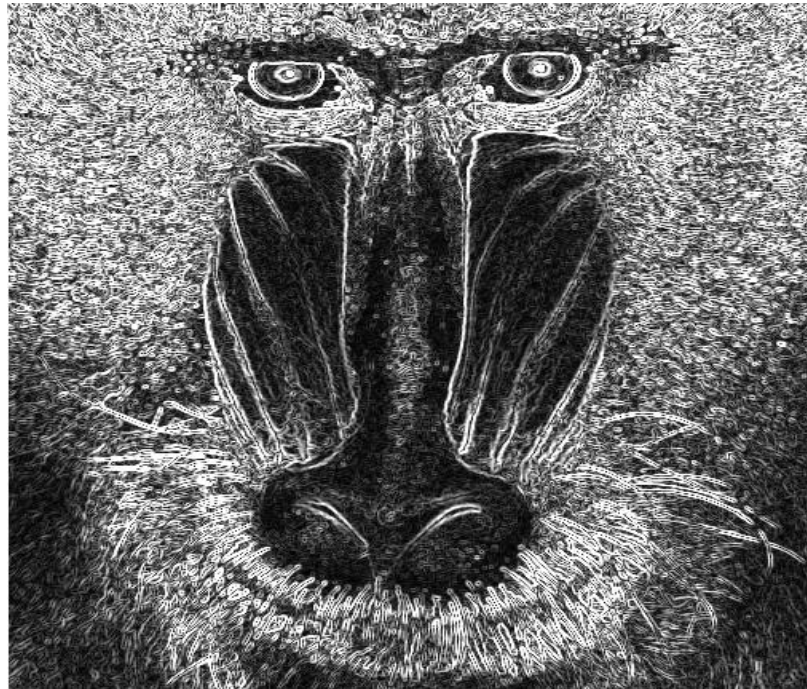3D graphics
  jmeint
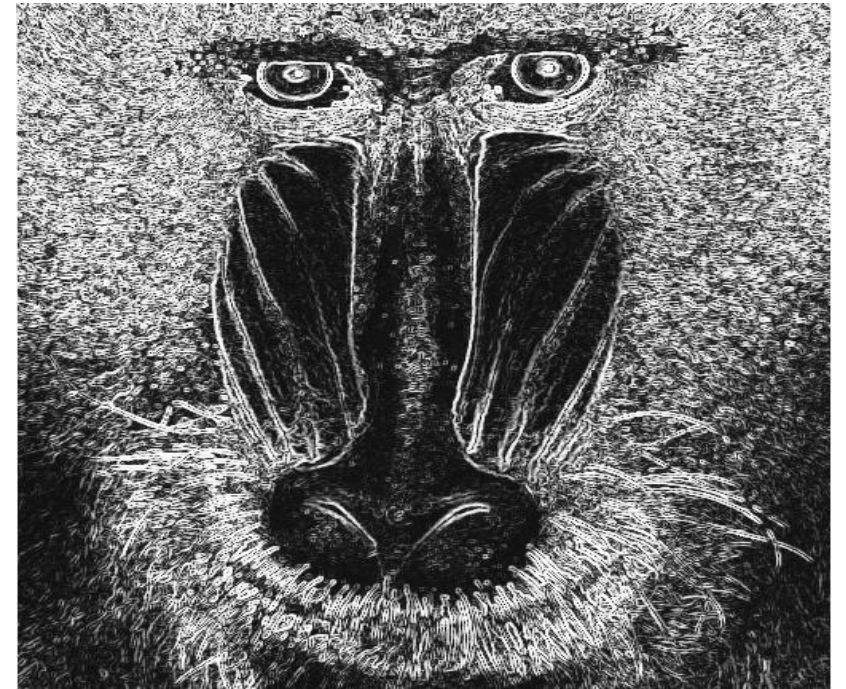Compression
  jpeg
Machine Learning
  kmeans
Image processing
  sobel

ORIGINAL CODE

NPU-TRANSFORMED CODE

Application

Algorithm

Language

Compiler
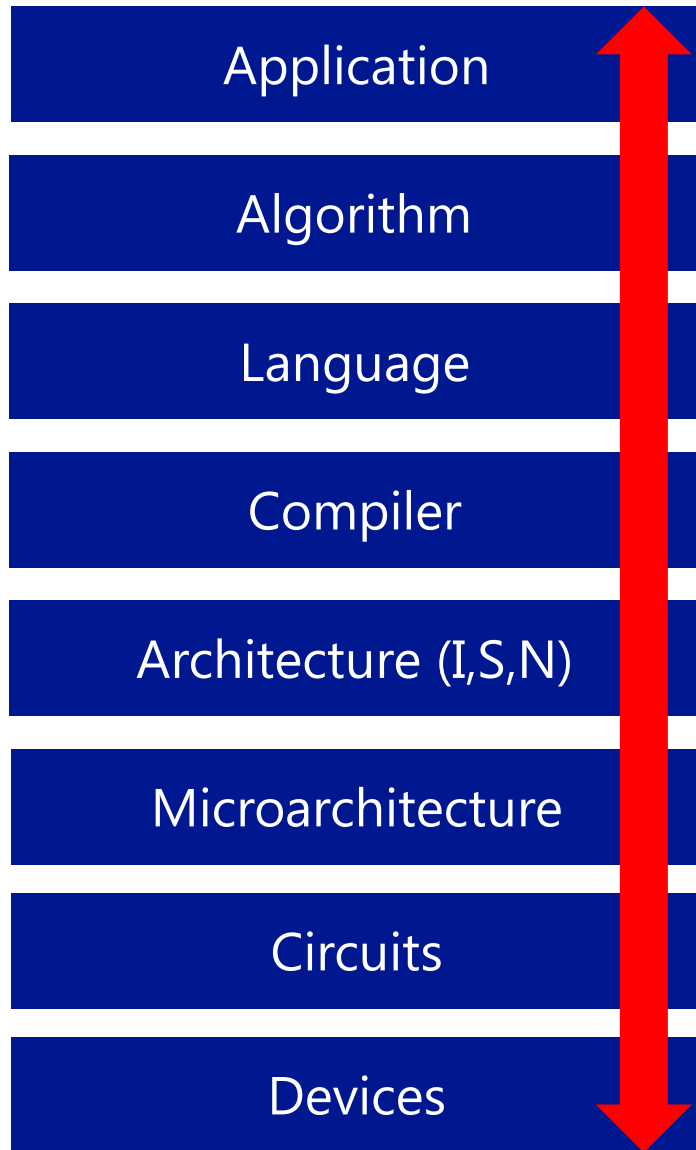
Architecture (I,S,N)

Microarchitecture

Circuits

Devices

Gains from NPUs still limited by being hidden behind standard abstractions

neural acceleration

# Neural Networks as a Platform

## The seeds of widening efforts are growing

## Many types of networks

Artificial Neural Networks

Bio-inspired Neural Networks (DNNs)

Bio-abstracted Neural Networks

## A recent quote from Jim Smith

"Giants are walking the Earth today, but we don't yet know whom they are"

# Five predictions for 2025

1. Moore's Law will be dead
2. Hardware/software compilation will be common
3. Neural execution will start to have a complete stack
4. Physics-based computation will be a hot topic
5. Machines will beat humans at many more tasks

# MOSFET (N-type) Transistor Operation