# **Unconventional** Compute Architectures for Enabling the Roll-Out of Deep Learning

Michaela Blott

Distinguished Engineer, Xilinx Research



**XILINX®**

# Background

> **Xilinx**

>> Fabless semiconductor company

>> Founded in Silicon Valley in 1984

>> Today:

– 3,500 employees

– $2.5B revenue

>> Invented the FPGA



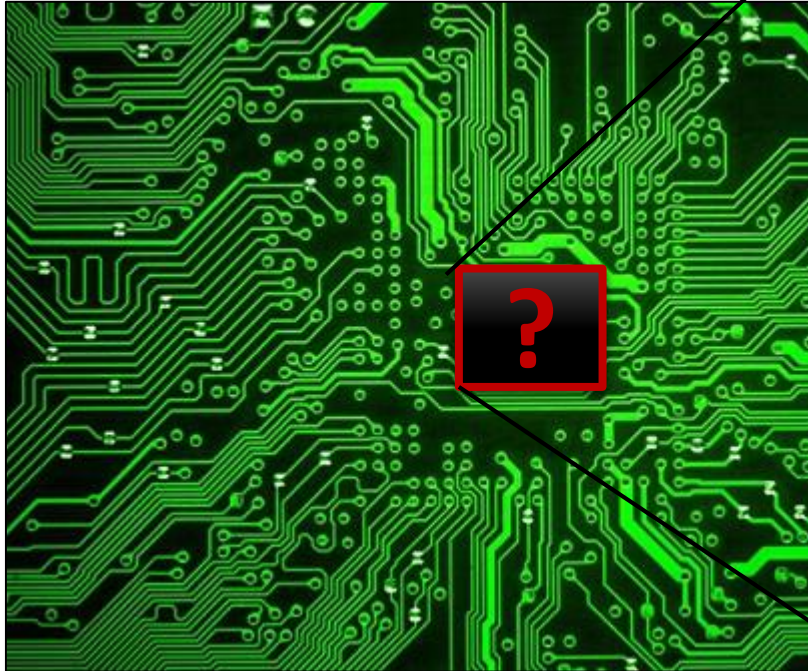1st FPGA in 1985: XC2064
128 3-input LUTs

XILINX.

# What are FPGAs?
## *Customizable, Programmable Hardware Architectures*

> The **chameleon** amongst the semiconductors…

>> Customizes IO interfaces, compute architectures, memory subsystems to meet the application

> **Classic use case:** Nothing else works, and you want to avoid ASIC implementation

> **Recent use cases:** Custom hardware architecture for performance or efficiency required

Non-standard IOs

Different functionality?

Higher performance or efficiency metrics?

**?**

HONEY, WAIT !
I CAN CHANGE !

XILINX.

# Trends Meet Technological Reality



**Roll-out of machine learning**
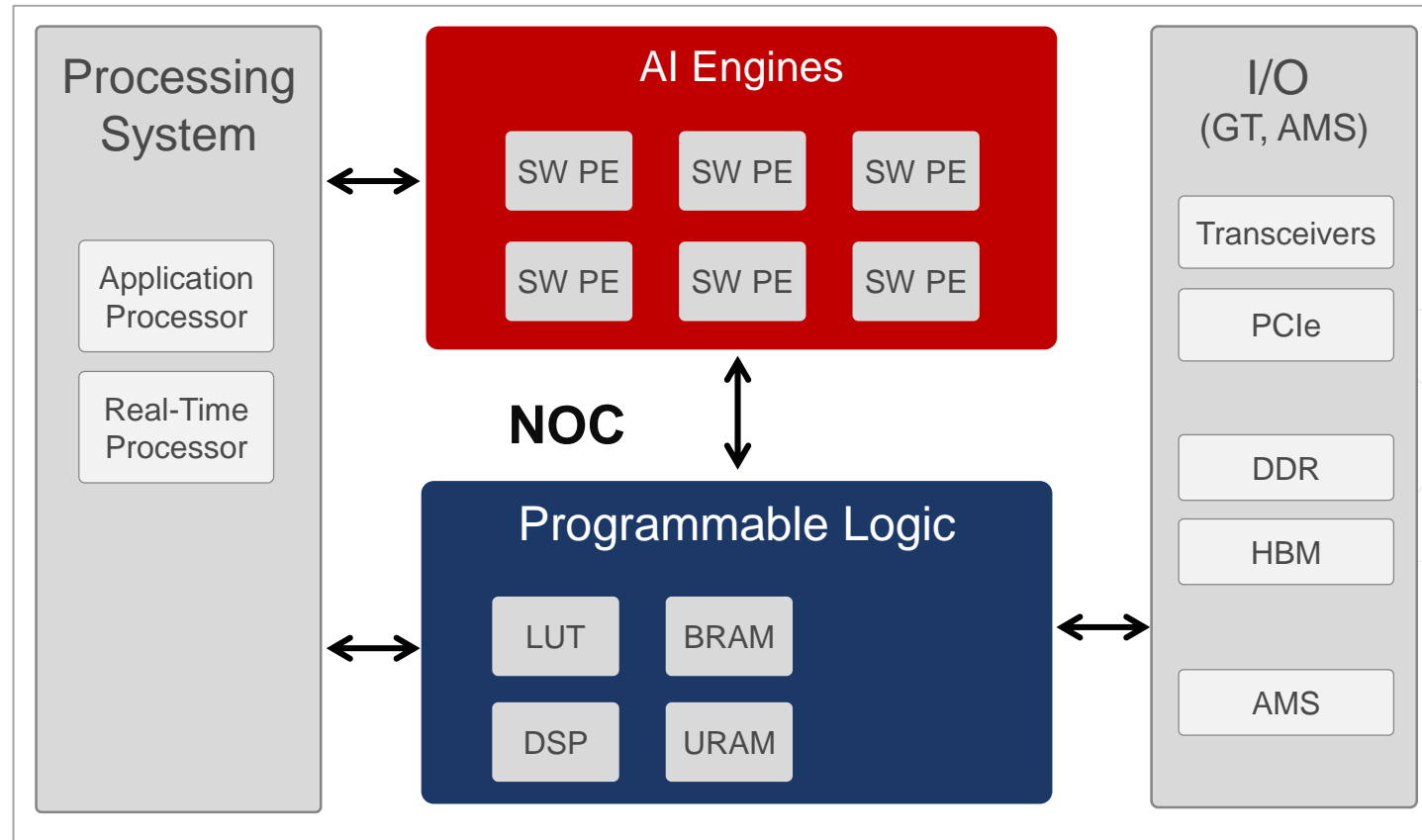"Potential to solve the unsolved"

**Explosion of Data**
"genomical"

**End of Moore's Law**
**End of Dennard Scaling**

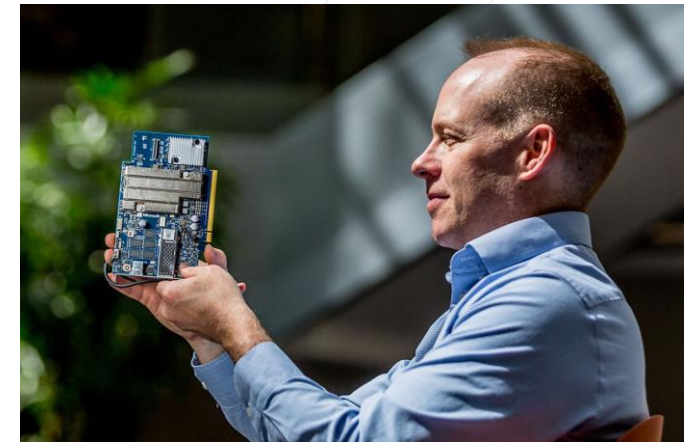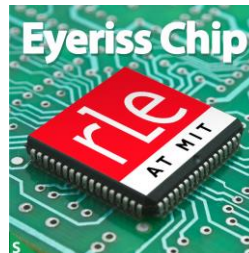## Era of heterogeneous computing has begun

- Diversification of increasingly heterogeneous "unconventional" devices
- Moving away from van Neumann architectures
- Archit5ectural and algorithmic innovation is needed

# Increasingly Heterogeneous Devices
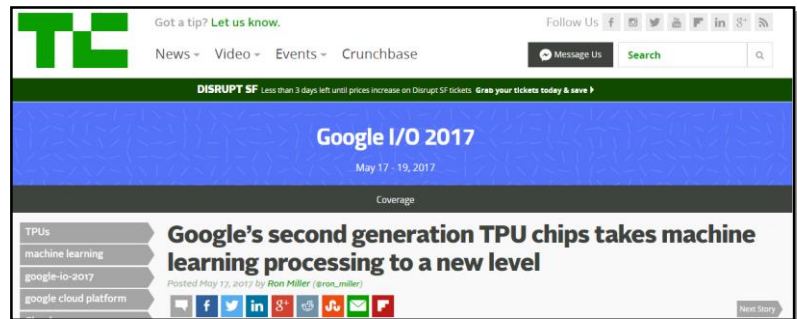# From the Xilinx World: Evolution of FPGAs to ACAPs

XILINX.

# More Unconventional:
# Customized Hardware for AI
## *DPU: Deep Learning Processing Unit*

> **Custom AI Accelerators (soft in FPGA and hard in ASICS)**



Microsoft Brainwave

# Popular DPU Architecture



CNN

DPU

DMA

MAC, Vector Processor

Onchip buffering

Matrix of Processing Engines

*"Layer by layer compute"*

XILINX

# *Even more unconventional:*
# Custom-Tailored Hardware Architectures (Macro-Level)
# *Synchronous Dataflow*



*"Hardware Architecture Mimics the NN Topology"*

> Customized feed-forward dataflow architecture to match network topology

> Higher compute and memory efficiency

# *Further unconventional at the Micro-Architecture, leveraging* Floating Point to Reduced Precision Neural Networks

# Reducing Precision
## *Scales Performance & Reduces Memory*

> **Reducing precision shrinks hardware cost**
>> Instantiate **100x** more compute within the same fabric
>> Thereby scale performance **100x**

> **Potential to reduce memory footprint**
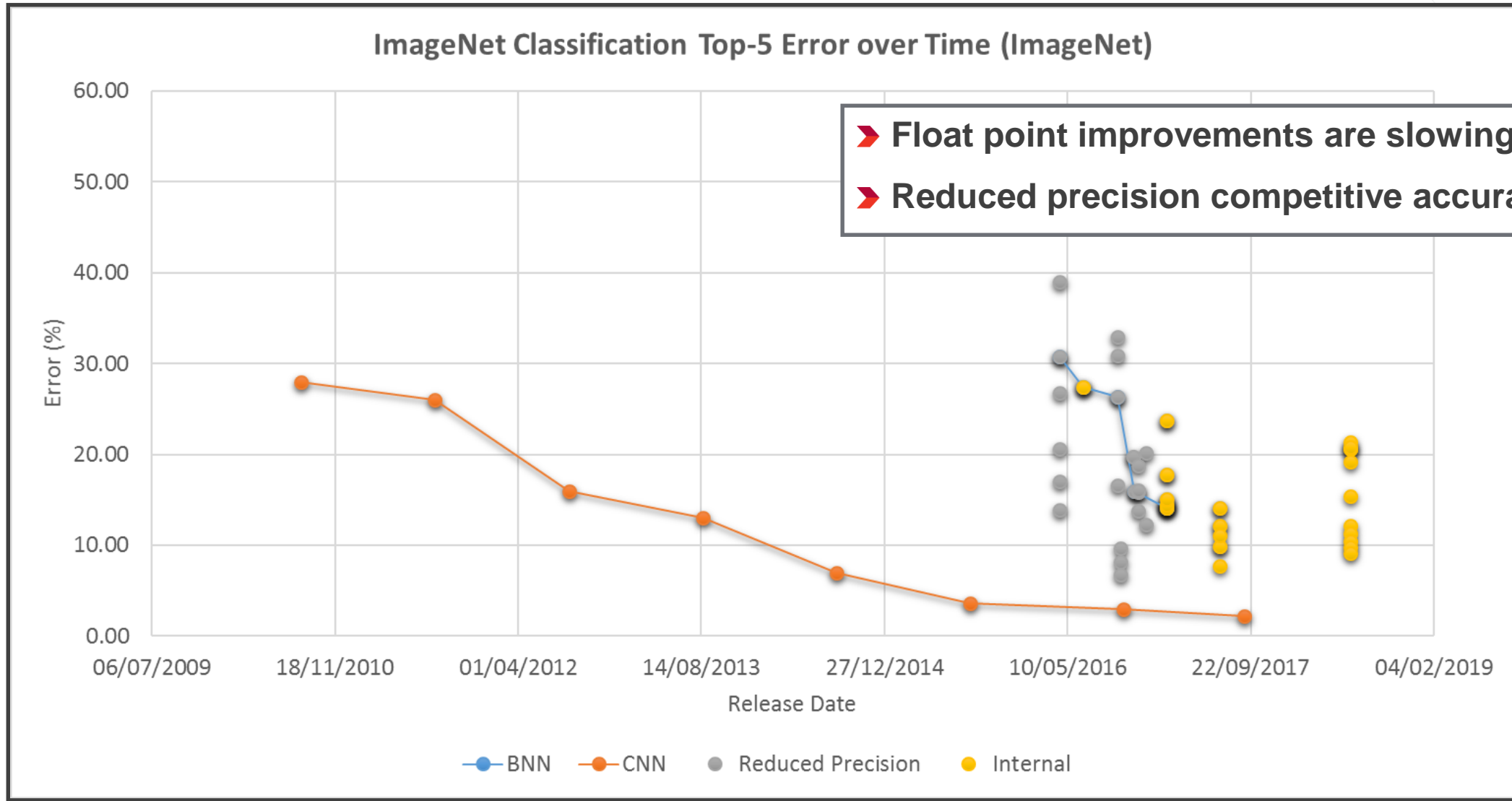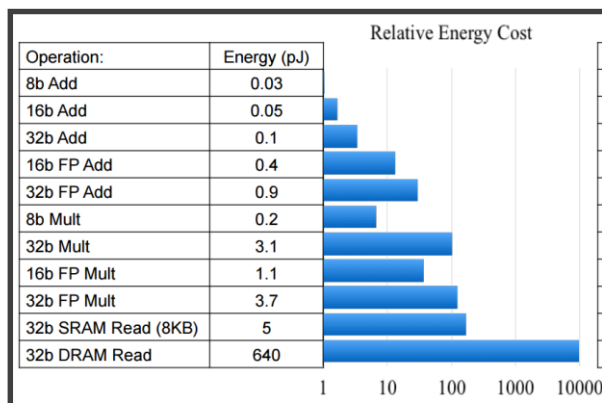>> NN model can stay on-chip => no memory bottlenecks

> **Reducing precision inherently saves power**

| Precision | Modelsize [MB] (ResNet50) |
|-----------|---------------------------|
| 1b        | 3.2                       |
| 8b        | 25.5                      |
| 32b       | 102.5                     |

**Relative Energy Cost**

| Operation: | Energy (pJ) |
|------------|-------------|
| 8b Add | 0.03 |
| 16b Add | 0.05 |
| 32b Add | 0.1 |
| 16b FP Add | 0.4 |
| 32b FP Add | 0.9 |
| 8b Mult | 0.2 |
| 32b Mult | 3.1 |
| 16b FP Mult | 1.1 |
| 32b FP Mult | 3.7 |
| 32b SRAM Read (8KB) | 5 |
| 32b DRAM Read | 640 |

*Source: Bill Dally (Stanford), Cadence Embedded Neural Network Summit, February 1, 2017*

XILINX.

# Design Space Trade-Offs



**IMAGENET CLASSIFICATION TOP5% VS COMPUTE COST F(LUT,DSP)**

Legend: ◆ 1b weights   ■ 2b weights   ✕ 5bit weights   ● 8bit weights   ✳ FP weights   ■ minifloat   + ResNet-50   — Syq

Y-axis: ERROR (%) — 0.00, 5.00, 10.00, 15.00, 20.00, 25.00, 30.00
X-axis: COMPUTE COST (LUTS + 100*DSPS) — 1.0, 10.0, 100.0, 1000.0, 10000.0, 100000.0, 1000000.0, 10000000.0, 100000000.0, 1000000000.0

Resnet18
8b/8b
Compute Cost 286
Error 10.68%

Resnet50
2b/8b
Compute Cost 127
Error 9.86%

Pareto-optimal solutions

Unconventional with reduced precision can
- reduce cost / resources
- save power
- scale performance

© Copyright 2018 Xilinx

XILINX

# Summary

- **Unconventional computing architectures emerge to help with the roll-out of deep learning**

- **Customized dataflow architectures and precisions provide dramatic performance scaling and energy efficiency**

- **Providing new exciting trade-offs within the design space**

**More information can be found at:**
**http://www.pynq.io/ml**

XILINX.

# Adaptable.
# Intelligent.

**More information can be found at:**
**http://www.pynq.io/ml**

**XILINX**