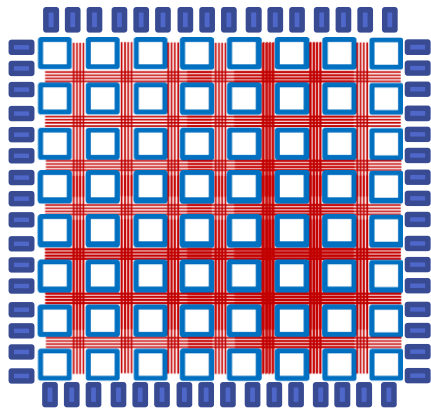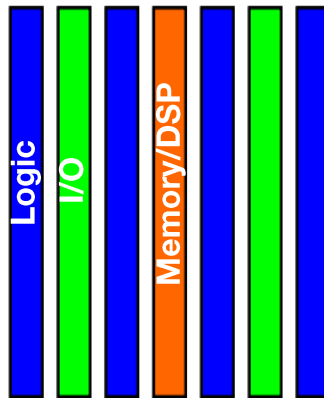# Introduction

- What is an FPGA
- Zynq MPSoC products for the embedded vision market
- Neural Networks on Zynq MPSoC products
- Benefits of Reduced Precision Neural Networks on Zynq MPSoCs
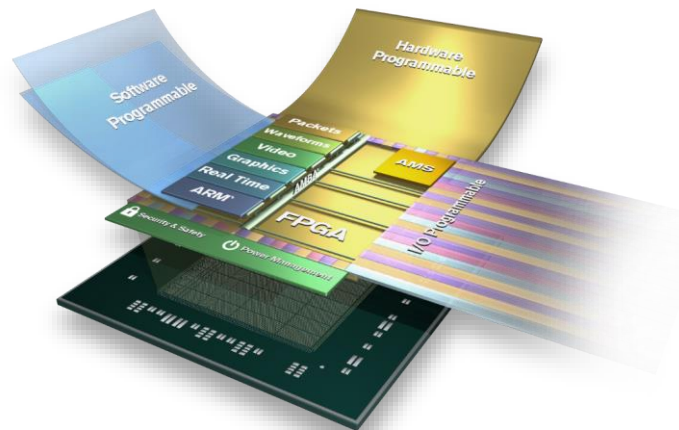- Programming environments

## LUTs

Basic bit-oriented logic
4-input, 6-input Lookup table

## Columns

Basic bit-oriented logic +
Word-oriented Multiply-
accumulate
Word-oriented Memory

## Processors + FPGA

Complete Processor systems
Dedicated programming environments
Heterogeneous MPSoC

Your program becomes a configuration that sets table values and switches via synthesis, Place and Route tools.

SW Programmability, Host code with Accelerator code, OpenCL, C/C++

High Level Synthesis (HLS), C/C++/OpenCL with Vivado IPI

Traditional HW design, Verilog or VHDL

# Zynq Ultrascale+ MPSOC

# Zynq® UltraScale+™ MPSoC: ZU7EV

- **Quad-core** ARM® Cortex™-A53 MPCore™ up to 1.5GHz

- **230 K LUTs**

- **11 Mb Block Ram**

- **27 Mb Ultra Ram**

- **1728 DSP slices**

- **Dedicated Video Encoder – Decoder hard block**

- **2 PCIexpress hard blocks (Gen3 x 16/ Gen4 x 8)**

# Neural Networks

# Neural Networks

- Exploit the Programmable Logic part of the MPSoC
- **Reduced Precision** is showing great promise
- The **trade-off** between precision and accuracy or error rate is essential.
- This presentation will show you the **pareto optimal** solutions!

# Multi-dimensional exploration space

**Trained Neural Network**

**End Device**

**Architecture**
Dataflow, systolic array, compression engines, sparse representations

**Hyperparameters**
Learning rates, regularization scheme, cost function, optimization scheme, pre/post processing

**Training data**
• ImageNet, C...
  VOC, GTSR...
  MNIST, CIFA...
  GSVHN, City...
  KITTY

**Training framework**
• Mxnet, Caffe,
  Tensorflow, T...
  Darknet

**Compute & data t...**
• Numerical represen...
  precision, quantizat...
  functions, activation...
  functions, non-linearity

**ML task (growing list)**
• Vision: image classification, recognition, sem...
  segmentation, S...
• Audio: voice co...
• Gaming strategy...
• Recommender s...

**Neural Network topology**
• AlexNet, GoogleNet, ResNet-X, Enet, Yolo

**Trained Neural Network**

Each combination yields to a different point in the multi-dimensional design space:
error, cost, throughput, latency, power

Data analysis required to understand the compromises and find optimal solutions



*ImageNet Classification: Error, compute cost, memory requirements, topology*

**XILINX**

# Training environment used for Reduced Precision

**TRAINING**

Input

labels

"cat"

Float. Model

= ?

"dog"

Many

Error

**Retraining**

**INFERENCE**

Input

Reduced Model

"dog" ✓

Fewer

## ResNet-50L ImageNet Top5 Error vs Hardware Cost

3b/5b

Effect of Retraining

4b/6b

6b/6b

1b/2b

2b/8b

8b/8b

Floating point baseline

Error (%)

Hardware Cost (LUT + 100*DSP)

○ Float    ○ Direct Quantization    ○ Retrained

Notation: 3b/5b: 3 bit weights/ 5 bit activation

- **<8bit: retraining**

- All code in C/C++
- Hardware Library is all HLS code
- Can execute on CPU and FPGA - No RTL needed

# Study of accuracy versus cost

- Several Networks were studied in detail, including **retraining.**
  - CNV on Cifar10
  - Resnet50 on IMAGENET

# CNV on CIFAR-10

- Topology;

  - Number of layers: 2 (3x3) Conv + Max Pool +
    2 (3x3) Conv + Max Pool + 2 Convolutional + 3 FC

  - Compute requirement: 112.5 MOPS/Frame

  - # Parameters: 1.54 M

- Parallelism adapted for BRAMs usage at 8 bits

  - Fixed parallelism for each precision

  - Estimated performances:

    - 385 FPS @ 200 MHz

    - Latency: 13 ms

    - Throughput: 43 GOPS
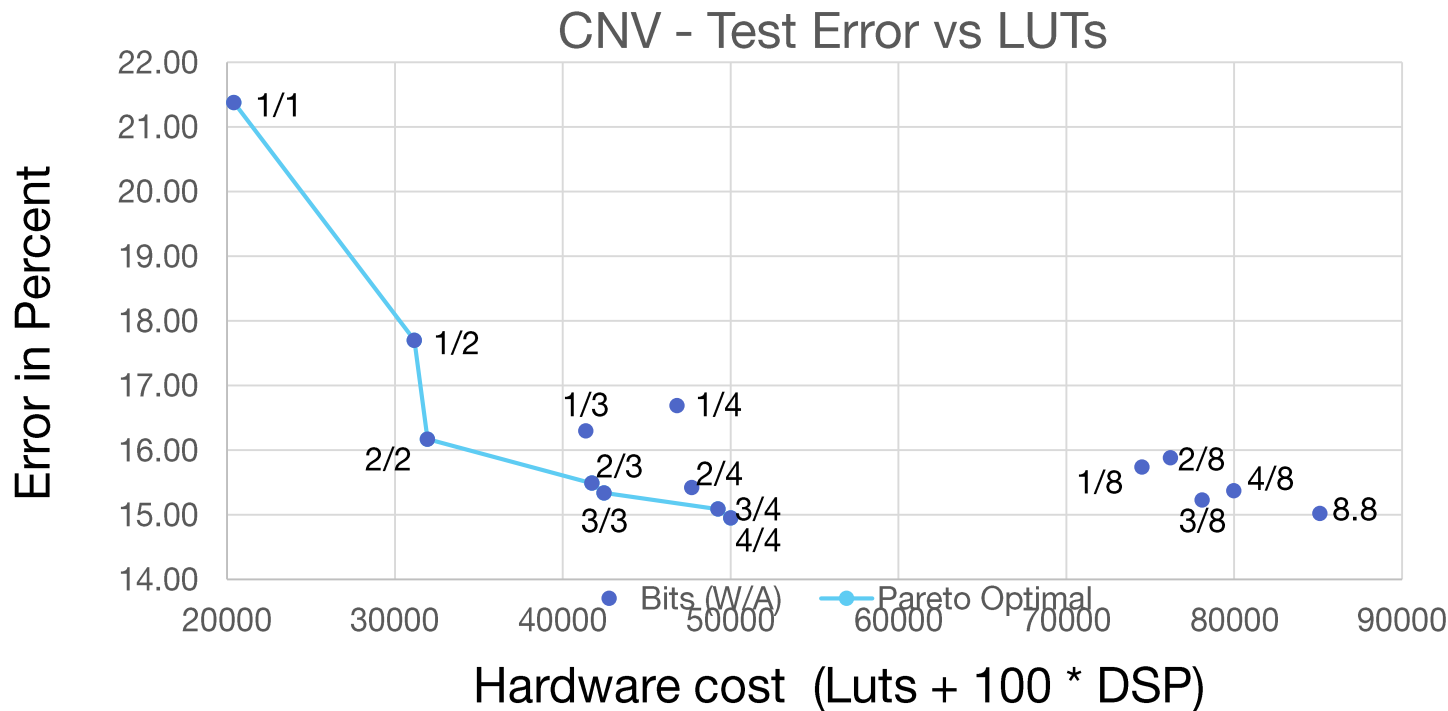
CNV - Test Error vs Memory Blocks

*Target Device ZU7EV ● Vivado 2017.3 tool suite ● 200 MHz target frequency ● Post-placed utilization ● #Blocks considered as BRAM36 + 4*URAM*

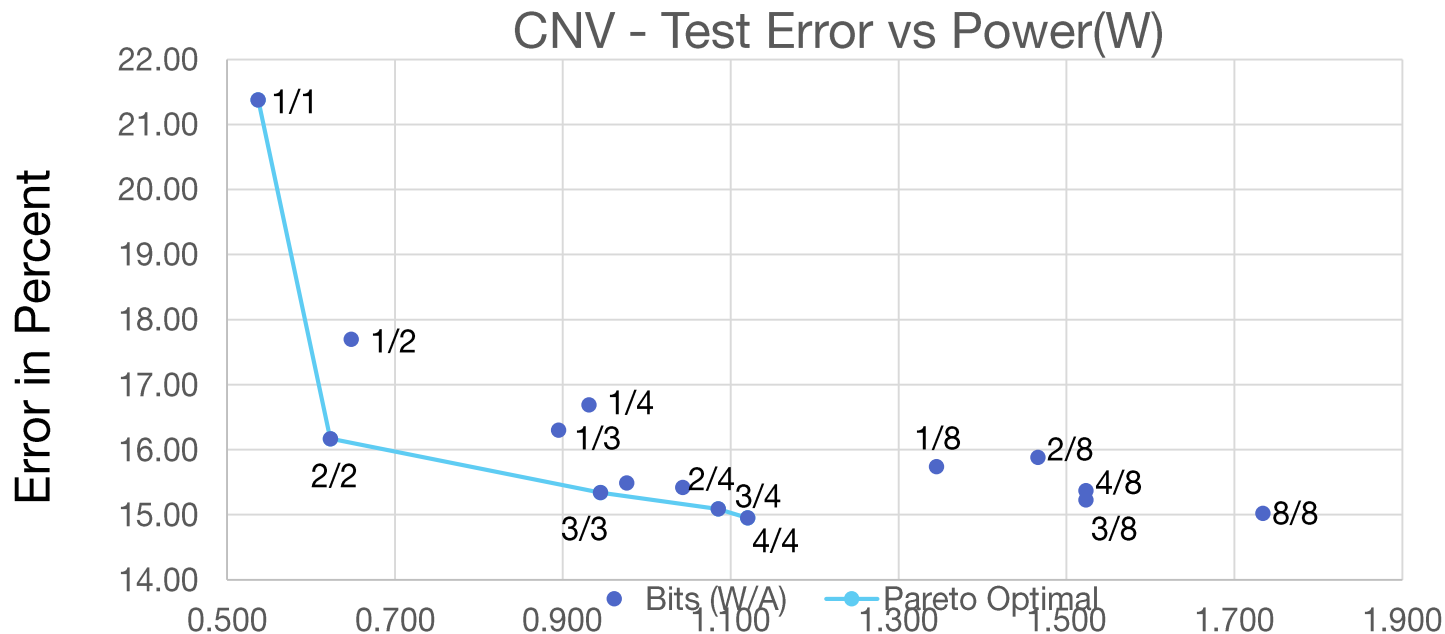CNV - Test Error vs LUTs

*Target Device ZU7EV ● Vivado 2017.3 tool suite ● Post-placed LUT utilization ● 200 MHz target frequency ●*
*Flow_PerfOptimized_high strategy for synthesis ● Performance_ExtraTimingOpt strategy for implementation*
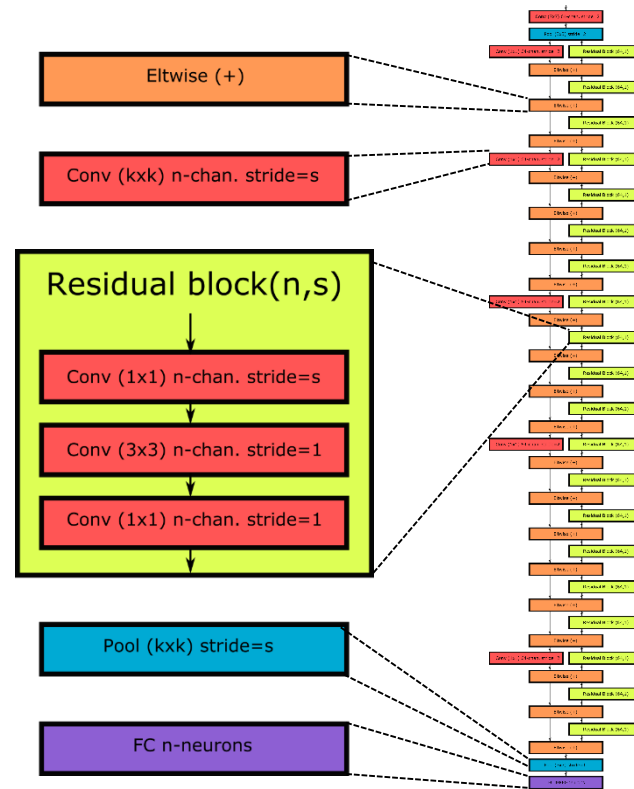
CNV - Test Error vs Power(W)

Target Device ZU7EV ● Ambient temperature: 25 ℃ ● 12.5% of toggle rate ● 0.5 of Static Probability ●
Power reported for PL accelerated block only

# Resnet50 with ImageNet accuracy study

XILINX

- Topology;
  - Number of layers: 53 Conv + 2 Pool + 1 FC
  - Compute requirement: 7.6 GOPS/Frame
  - # Parameters: 25.5 M
- Comparison between direct quantization and retraining
- Performance Model;
  - Combination of LUTs and DSPs used
  - Hardware cost:
    - a weighted sum of LUTs and DSPs required per operation (LUTs + 100*DSPs)

ResNet-50 - Top5 Validation Error

ResNet-50 - Top5 Validation Error

Legend: Float · Direct Quantization · Retrained · Pareto Optimal

# Conclusions

# Reduced Precision Conclusions on FPGAs

- **Pareto optimal solutions** show best implementations for a certain error
- Precisions **well below 8-bit** are very promising, benefits are:
  - Lower power
  - Less Hardware
  - At acceptable error rates
- **Re-training** is essential to exploit reduced precision implementations
- Xilinx FPGAs are an excellent implementation platform for reduced precision Neural Networks

Xilinx products:

- [www.xilinx.com](www.xilinx.com)

- [https://www.xilinx.com/video/application/revision.html](https://www.xilinx.com/video/application/revision.html)

- [https://www.xilinx.com/products/design-tools/embedded-vision-zone.html](https://www.xilinx.com/products/design-tools/embedded-vision-zone.html)

University support and open source:

- [https://www.xilinx.com/support/university/.html](https://www.xilinx.com/support/university/.html)

- [http://www.pynq.io/home.html](http://www.pynq.io/home.html)

- [https://github.com/Xilinx](https://github.com/Xilinx)

- Embedded Vision Alliance:

  [https://www.embedded-vision.com/summit](https://www.embedded-vision.com/summit)