Fotakis Tzanis
Electrical & Computer Engineering School
Technical University of Crete

# Analysis and Design Methodology of Convolutional Neural Networks mapping on Reconfigurable Logic
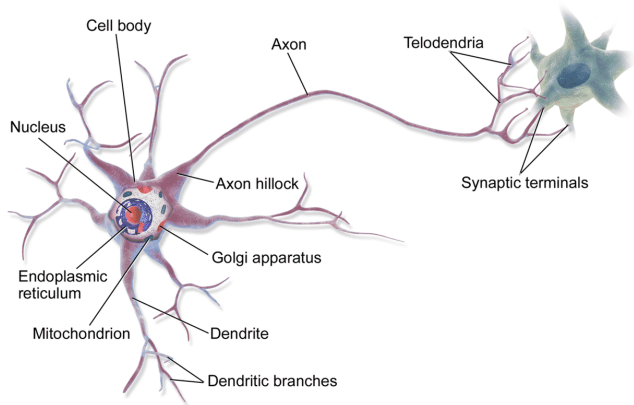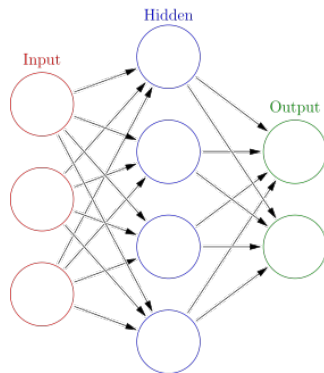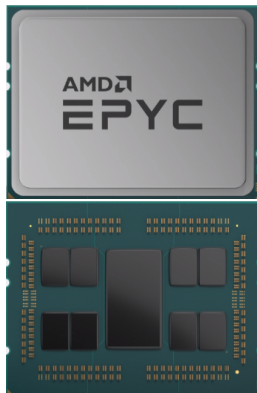
Diploma Thesis

# Table of Contents

# What are Neural Networks?

# Representation of Artificial Neural Networks

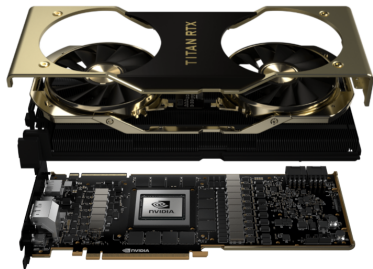## ANNs on CPUs



AMD Epyc 7002 series chip

+ Easy development
+ High clock frequency
+ Advanced Vector Extensions (AVX)
+ Streaming SIMD Extensions (SSE)

- Costly
- Non-scalable
- Energy inefficient
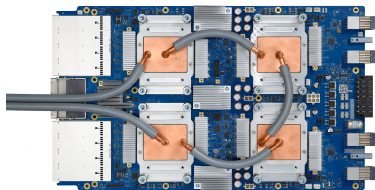- Traditional Low Bandwidth Memory - up to 50 GB/s

## ANNs on GPUs



NVIDIA Titan RTX card

+ Relatively easy development
+ Very high parallelism - Thousands of cores
+ Specialized Tensor Cores
+ Vector processing - Streaming Multiprocessors
+ High Bandwidth Memory - up to 750 GB/s
+ Multiple GPUs in a system

- Very power hungry
- Costly to scale up
- Increased latency

## ANNs on ASICs



Google TPU v3

+ Best parallelism
+ Lowest power & energy consumption
+ High Bandwidth Memory
- Extremely expensive to design and produce
- Serve a single purpose
- Can become deprecated fast - AI field is still developing

## ANNs on FPGAs



FORTH QFDB

+ Flexible
+ Low power & energy consumption
+ Low latency
+ Standalone

- Difficult to develop
- Constrained resources
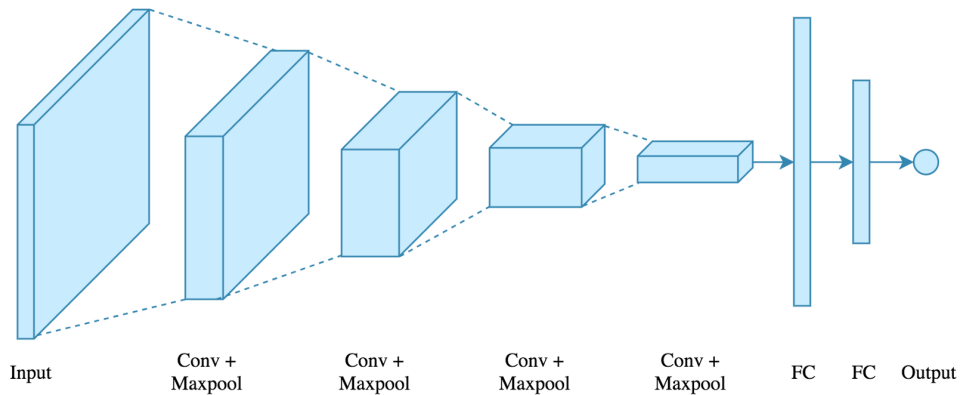
## What is Reconfigurable Logic

- Look-Up Tables (LUTs)
- Flip-Flops (FFs)
- Block-RAM (BRAM)
- Ultra-RAM (URAM)
- Digital Signal Processing (DSP) blocks
- Hard Processor cores (SoC/MPSoC devices)
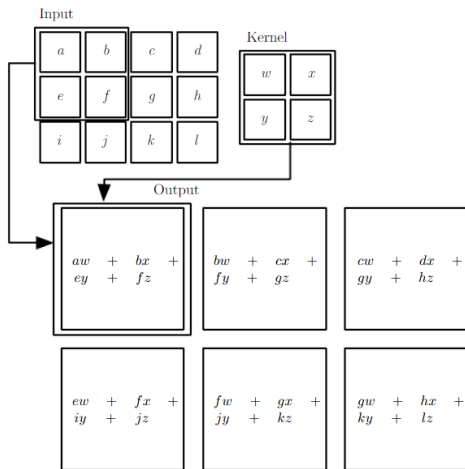- DDR & HBM (modern FPGAs)

# Neural Networks

## ANN Architectures

- Multiclass Perceptron
- Autoencoder
- Convolutional
- Recurrent
- Long short-term memory

# Typical Convolutional Neural Network



Input    Conv + Maxpool    Conv + Maxpool    Conv + Maxpool    Conv + Maxpool    FC    FC    Output

# The Convolution Layer



Input

| | | | |
|---|---|---|---|
| $a$ | $b$ | $c$ | $d$ |
| $e$ | $f$ | $g$ | $h$ |
| $i$ | $j$ | $k$ | $l$ |

Kernel

| | |
|---|---|
| $w$ | $x$ |
| $y$ | $z$ |

Output

| | | |
|---|---|---|
| $aw + bx + ey + fz$ | $bw + cx + fy + gz$ | $cw + dx + gy + hz$ |
| $ew + fx + iy + jz$ | $fw + gx + jy + kz$ | $gw + hx + ky + lz$ |

# The Max-Pooling Layer

## The Fully-Connected Layer



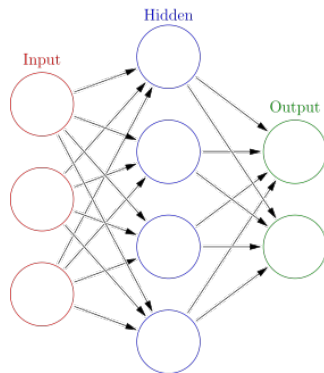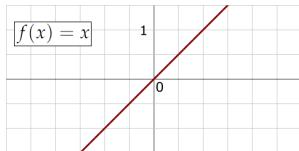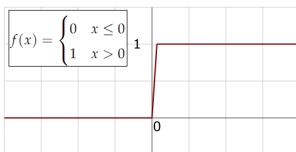$$\text{Output}(i) = \text{Bias}(i) + \sum_{j=1}^{N} \text{Input}(j) * \text{Weight}(i, j), \text{ for } i = 1, 2, ..., M$$

## Activation Functions



$f(x) = x$

Identity

$$f(x) = \begin{cases} 0 & x \leq 0 \\ 1 & x > 0 \end{cases}$$

Binary step

$$f(x) = \sigma(x) = \frac{1}{1+e^{-x}}$$

Sigmoid

$$f(x) = tanh(x) = \frac{e^x - e^{-x}}{e^x - e^{-x}}$$

tanh

$$f(x) = \begin{cases} 0 & x < 0 \\ x & x > 0 \end{cases}$$

ReLU

$$f_i(\mathbf{x}) = \frac{e^{x_i}}{\sum_{j=1}^{J} e^{x_j}}$$

Softmax

# CNN Architectures: LeNet-5



INPUT 32x32 — C1: feature maps 6@28x28 — S2: f. maps 6@14x14 — C3: f. maps 16@10x10 — S4: f. maps 16@5x5 — C5: layer 120 — F6: layer 84 — OUTPUT 10

Convolutions — Subsampling — Convolutions — Subsampling — Full connection — Full connection — Gaussian connections

# CNN Architectures: AlexNet

# CNN Architectures: ZFNet

# CNN Architectures: GoogLeNet / Inception

## CNN Architectures: VGGNet



| | Number of Parameters (millions) | Top-5 Error Rate (%) |
|---|---|---|
| VGG-11 | 133 | 10.4 |
| VGG-11 (LRN) | 133 | 10.5 |
| VGG-13 | 133 | 9.9 |
| VGG-16 (Conv1) | 134 | 9.4 |
| VGG-16 | 138 | 8.8 |
| VGG-19 | 144 | 9.0 |

# CNN Architectures: ResNet

# Related Work

# Software Frameworks

## Software Frameworks: TensorFlow

- By Google
- Most popular
- Lower-level - Much coding - High Configurability
- Python, Javascript, C++, C#, Java, Go & Julia interfaces
- Targeted for production
- Static computation graph - Efficient but less flexible
- CPU, GPU & TPU acceleration
- Server, Mobile & Embedded platforms

## Software Frameworks: PyTorch



- By Facebook
- Based on Torch
- CPU & GPU acceleration
- Dynamically updated graph
- Targeted for prototyping & research
- Contains many pre-trained models

## Software Frameworks: Keras



- Higher-level
- Back-ends: TensorFlow, Theano & CNTK
- Easy huge models - Less configurable
- Targeted for learning & prototyping

## Software Frameworks: Caffe

# Caffe

- By Berkley, University of California
- Written in C++ - Python interface
- CPU & GPU acceleration
- Not all Neural Networks supported
- Caffe2 merged with PyTorch

# Hardware Frameworks

# Hardware Frameworks: Google TPU



Google TPU v3

- By Google since 2015
- Used with TensorFlow
- Accelerate 95% of their AI needs
- Publicly available on Google Compute Engine since 2017
- Initially only inference & 8-bit fixed-point
- 200x compared to Intel Haswell CPU
- 70x compared to NVIDIA K80 GPU
- Version 2 and above add training & floating-point

## Hardware Frameworks: Google TPU v1 Block Diagram



TPU v1 Block Diagram

- Operate on matrices - GPUs operate on vectors
- 2x Matrix Multiplier Units (MXUs)
- MXU based on systolic arrays
- 16k MAC ops/cycle
- Up to 128GB HBM (TPU v3)
- 92 TOPS

# Hardware Frameworks: Google TPU v3 Pods



TPU v3-32
(32 cores, 4x4 slice)

TPU v3-128
(128 cores, 8x8 slice)

TPU v3-512
(512 cores, 16x16 slice)

TPU v3-1024 (1024
cores, 16x32 slice)

- 2048 TPU Cores
- 32TB HBM
- 92 PFLOPS
- Suitable for very large models (weeks/months of training) - no custom operations

## FPGA Frameworks: Xilinx CHaiDNN

- Released in February 2018 - Targeted for CNN inference
- 6-bit & 8-bit fixed-point - variable through layers
- Similar to single-precision floating-point
- Dynamic fixed-point Quantization & Xilinx Quantizer
- 128-1024 Double-Pumped DSPs - Up to 700MHz - URAM supported
- Unsupported layers added via Software
- Fully-Connected & SoftMax layers implemented through Software
- Hardware & Software layers run in parallel
- PetaLinux - Caffe framework
- DietCHai for smaller MPSoCs

## FPGA Frameworks: Xilinx DPU

- Released in February 2019 - Replaces CHaiDNN - Still in development
- Targeted for CNN inference - 8-bit fixed-point
- On-Chip memory utilized as buffer
- All layers are hardware accelerated
- Up to four DPU cores in a single DPU IP
- Double Data Rate / Double-Pumped DSPs
- 512-4096 operations per cycle per core

# FPGA Frameworks: Xilinx DPU v1



- 96x16 DSP Systolic Array

## FPGA Frameworks: Xilinx DPU v2



- Hybrid Computing Array
- Processing Elements - based on fine grained building blocks (multipliers, adders, accumulators)
- Deep Pipeline

## FPGA Frameworks: Xilinx DPU v3



- Multiple Batch Engines
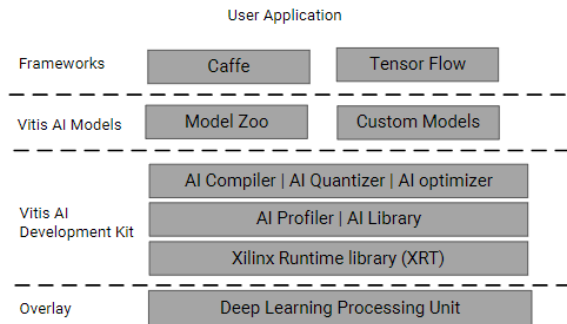- Multiple DPU Cores

# FPGA Frameworks: Xilinx Vitis AI



User Application

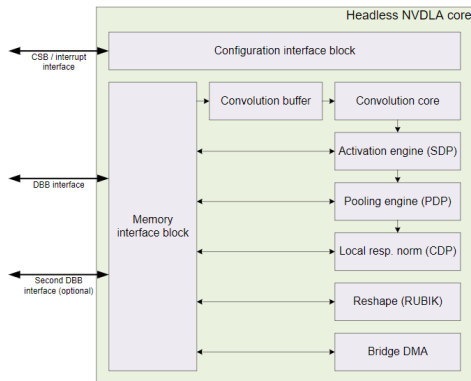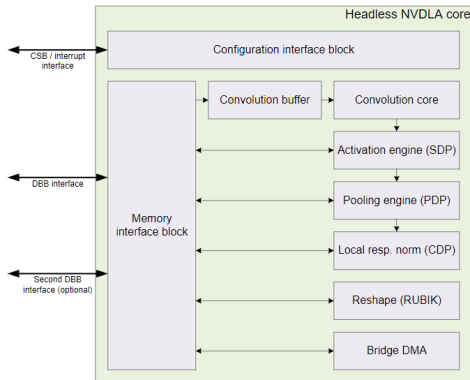| Frameworks | Caffe | Tensor Flow |
| Vitis AI Models | Model Zoo | Custom Models |
| Vitis AI Development Kit | AI Compiler | AI Quantizer | AI optimizer |
| | AI Profiler | AI Library |
| | Xilinx Runtime library (XRT) |
| Overlay | Deep Learning Processing Unit |

- Released in December 2019
- High-level abstraction
- AI inference applications
- Based on Xilinx DPU
- Optimized IPs, tools & libraries
- PetaLinux
- Instruction optimization - Vitis AI Compiler
- Vitis AI Quantizer - 8-bit fixed point parameters

# FPGA Frameworks: NVIDIA Deep Learning Accelerator (NVDLA)



Headless NVDLA core

- Released in Q3 2017
- Free & Open architecture
- Goal to standardize inference DL accelerator development
- Headless implementation: Manager is the main system processor
- Headed implementation: Manager is a companion microcontroller
- Modular & Highly customizable
- Suitable for both FPGAs & ASICs

# FPGA Frameworks: NVIDIA Deep Learning Accelerator (NVDLA)



- Binary and 4-bit integer up to 64-bit floating-point
- Convolution core
- Single Data processor
- Planar Data processor
- Channel Data processor
- Data Reshape Engine
- Memory-to-Memory or Pass-Through

# Theoretical Modeling & Robustness Analysis

## PyTorch, C/C++, MATLAB

PyTorch → pure Python → C/C++ & MATLAB

- Replicate PyTorch functionality
- Evaluation using PyTorch
- PyTorch pre-build & pretrained AlexNet as a reference
- 2500 images of Kaggle cats & dogs database
- Better understanding of underlining algorithms
- Explore Hardware implementation opportunities
- Quantization techniques → reduce memory footprint
- Algorithmic optimizations
- Minimize hardware resources and optimize performance using various tools

## Algorithms

- Convolution Layer
- Max-Pooling Layer
- Fully-Connected Layer
- ReLU
- SoftMax

# Memory Footprint

- Classic hardware architectures $\rightarrow$ Compute bound
- ASICs & FPGAs $\rightarrow$ Memory bound
- Highest benefit: Memory requirements fit into BRAM (order of MBs)
- Otherwise, external memory used $\rightarrow$ latency & IO stalls
- Goal: Minimize memory footprint & bandwidth

## Memory Footprint: AlexNet Parameters (float32)

| Layer | #Parameters | Footprint | Memory (%) |
|-------|-------------|-----------|------------|
| Conv1 | $64 * 3 * 11 * 11 = 23232$ | 92.92KB | 0.04 |
| Conv2 | $192 * 64 * 5 * 5 = 307200$ | 1.22MB | 0.5 |
| Conv3 | $384 * 192 * 3 * 3 = 663552$ | 2.65MB | 1.09 |
| Conv4 | $256 * 384 * 3 * 3 = 884736$ | 3.53MB | 1.45 |
| Conv5 | $256 * 256 * 3 * 3 = 589824$ | 2.35MB | 0.97 |
| FC1 | $9216 * 4096 = 37748736$ | 150.99MB | 61.79 |
| FC2 | $4096 * 4096 = 16777216$ | 67.10MB | 27.46 |
| FC3 | $4096 * 1000 = 4096000$ | 16.38MB | 6.70 |
| **Total** | 61090496 | 244.36MB | 100 |

## Memory Footprint: AlexNet Data Stages (float32)

| Layer | #Data | Footprint | Memory (%) |
|-------|-------|-----------|------------|
| Image | $3 * 224 * 224 = 150528$ | 150.52KB | 6.07 |
| Conv1 | $64 * 55 * 55 = 193600$ | 774.40KB | 31.22 |
| MaxPool1 | $64 * 27 * 27 = 46656$ | 186.62KB | 7.52 |
| Conv2 | $192 * 27 * 27 = 139968$ | 559.87KB | 22.57 |
| MaxPool2 | $192 * 13 * 13 = 32448$ | 129.79KB | 5.23 |
| Conv3 | $384 * 13 * 13 = 64896$ | 259.58KB | 10.46 |
| Conv4 | $256 * 13 * 13 = 43264$ | 173.05KB | 6.98 |
| Conv5 | $256 * 13 * 13 = 43264$ | 173.05KB | 6.98 |
| MaxPool3 | 9216 | 36.86KB | 1.49 |
| FC1 | 4096 | 16.38KB | 0.66 |
| FC2 | 4096 | 16.38KB | 0.66 |
| FC3 | 1000 | 4KB | 0.16 |
| **Total** | 682856 | 2.48MB | 100 |

## Memory Footprint Reduction

- Data type bit-width shortening (float64-float16)
- Simpler data types (fixed-point/integers)
- Binary
- Quantization
- Quantization aware training
- Compression
- K-Means clustering
- Second Level Codebook

Diploma Thesis

Introduction  Neural Networks  Related Work  **Theoretical Modeling**  Architecture Design  FPGA Implementation  Results  Conclusion

## Memory Footprint Reduction

# Trading accuracy to performance.

## Memory Footprint Reduction: Evaluation

- Baseline: PyTorch pre-trained pre-build AlexNet model
- Inferencing 2500 pre-transformed Kaggle cats & dogs images
- Top-1 error rate
- MATLAB implementation used
- PyTorch & C/C++ do not support half-floating point

## Memory Footprint Reduction: Floating Point

Convert 32-bit floats to their closest representation.

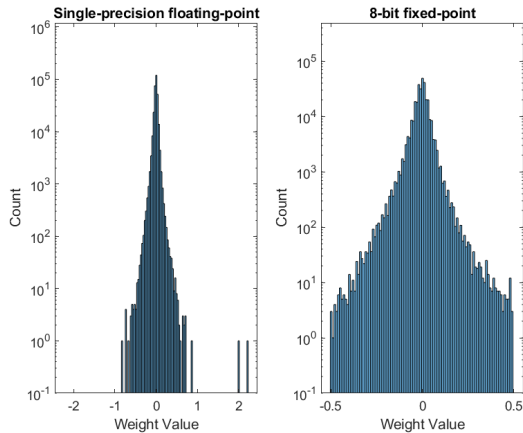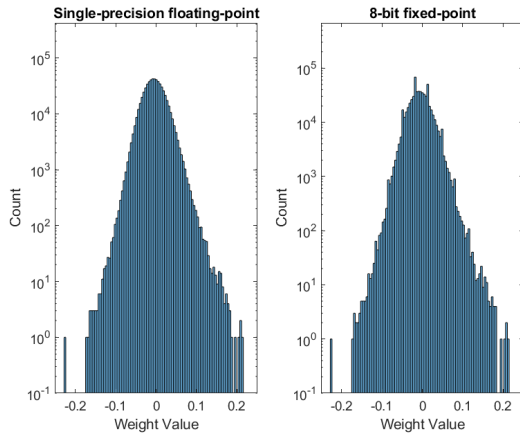| Tool | Data type | Top-1 Error rate (%) | Avg. inference time (sec) |
|------|-----------|---------------------|---------------------------|
| PyTorch | float64 | 0 | 0.091 |
| PyTorch | float32 | 0 | 0.034 |
| MATLAB | float64 | 0 | 6.624 |
| MATLAB | float32 | 0 | 8.162 |
| MATLAB | float16 | 0.36 | 147.480 |

## Memory Footprint Reduction: Fixed Point Parameters

- Convert 32-bit float number sets to fixed-point
- Select best radix-point position to most accurately represent the number set
- Use same scale factor on whole set
- Every layer has its own scale factor

$$\text{Position} = \text{argmin}_{i=0}^{W}[\frac{\sum_{j=1}^{size(S)} |S_j - FixPtConvert(S_j, W, i)|}{size(S)}]$$

## Memory Footprint Reduction: Fixed Point Parameters

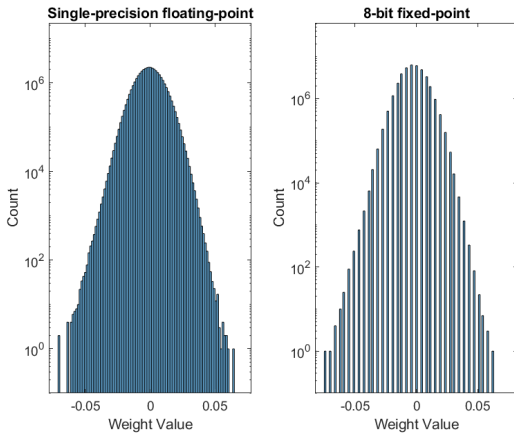| Tool | Data type | Top-1 Error rate (%) | Avg. inference time (sec) |
|------|-----------|----------------------|---------------------------|
| MATLAB | fixed64 | 0 | 7.318 |
| MATLAB | fixed32 | 0 | 7.692 |
| MATLAB | fixed16 | 22 | 6.650 |
| MATLAB | fixed14 | 28.44 | 6.813 |
| MATLAB | fixed12 | 36.24 | 6.797 |
| MATLAB | fixed10 | 77.07 | 6.929 |
| MATLAB | fixed8 | 100 | 6.312 |

## Memory Footprint Reduction: Fixed Point Parameters



Histogram limits significantly altered

# Memory Footprint Reduction: Fixed Point Parameters



Significant spiking

## Memory Footprint Reduction: Fixed Point Parameters



Severe subsampling

## Memory Footprint Reduction: Fixed Point Parameters

Use Mean Squared Error (MSE)

$$\text{Position} = \text{argmin}_{i=0}^{32}\left[\frac{\sum_{j=1}^{size(S)} |S_j - \text{FixPtConvert}(S_j, W, i)|^2}{size(S)}\right]$$

Use Mean Quarted Error (MQE)

$$\text{Position} = \text{argmin}_{i=0}^{32}\left[\frac{\sum_{j=1}^{size(S)} |S_j - \text{FixPtConvert}(S_j, W, i)|^4}{size(S)}\right]$$
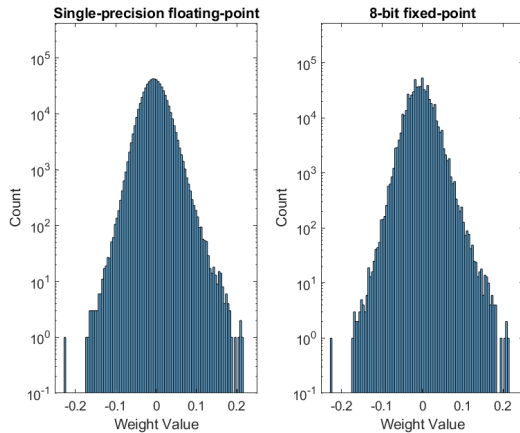
## Memory Footprint Reduction: Fixed Point MQE

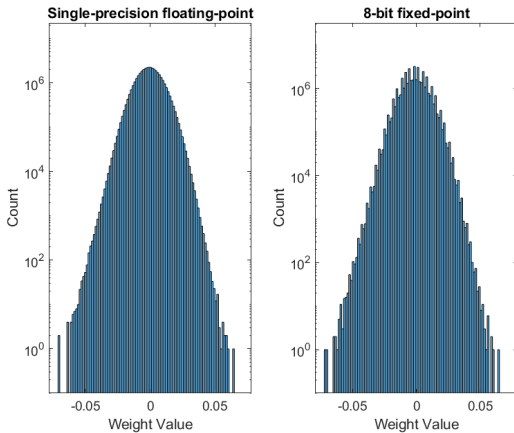| Tool | Data type | Top-1 Error rate (%) |
|------|-----------|----------------------|
| MATLAB | fixed64 | 0 |
| MATLAB | fixed32 | 0 |
| MATLAB | fixed16 | 4.42 |
| MATLAB | fixed14 | 17.59 |
| MATLAB | fixed12 | 48.11 |
| MATLAB | fixed10 | 86.91 |
| MATLAB | fixed8 | 99.3 |

## Memory Footprint Reduction: Fixed Point Parameters



Histogram limits identical

## Memory Footprint Reduction: Fixed Point Parameters
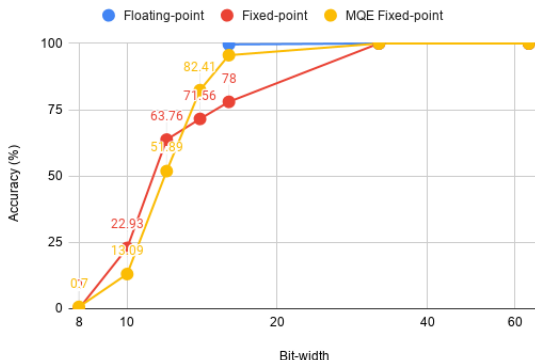


Slight spiking

## Memory Footprint Reduction: Fixed Point Parameters



No subsampling, but spiking

# Memory Footprint Reduction: All data types tested



- Accuracy degradation is expected
- Model dependent
- Training dependent

## Memory Footprint Reduction: Fixed Point Activations

- Floating-point arithmetic scales automatically
- Fixed-point activations may overflow
- Quantize every layer outputs
- Keep the upper N most significant bits to retain accuracy
- Finding on runtime the uppermost on is computationally intensive $\rightarrow$ any fixed-point benefits get obsolete

## Memory Footprint Reduction: Fixed Point Activations

Calculate optimal activation scale factor per layer.

$$\text{Theoretical}_{\text{bitWidth}} = \text{input}_{\text{bitWidth}} + \text{weight}_{\text{bitWidth}} + \lceil \log_2 \#\text{Additions} \rceil$$

## Memory Footprint Reduction: Fixed Point Activations

| Layer | Theoretical bit-width | Practical bit-width |
|-------|-----------------------|---------------------|
| Input | 8 | 8 |
| Conv1 | 25 | 17 |
| Conv2 | 27 | 14 |
| Conv3 | 27 | 15 |
| Conv4 | 28 | 15 |
| Conv5 | 28 | 17 |
| FC1 | 30 | 17 |
| FC2 | 28 | 17 |
| FC3 | 28 | 17 |

- Theoretical worse case scenario significantly differs from practical
- Inference 2000 images, find maximum absolute valued activation per layer
- Maximum theoretical bit-width is 30-bits $\rightarrow$ all activations fit in 32-bit integers before quantization

## Memory Footprint Reduction: Fixed Point Activations

| Layer | Weights | Bias | Output |
|-------|---------|------|--------|
| Input | -7 | - | - |
| Conv1 | -7 | -5 | -2 |
| Conv2 | -5 | -7 | 0 |
| Conv3 | -7 | -7 | 3 |
| Conv4 | -8 | -6 | 5 |
| Conv5 | -9 | -5 | 10 |
| FC1 | -10 | -10 | 15 |
| FC2 | -10 | -9 | 19 |
| FC3 | -9 | -9 | 23 |

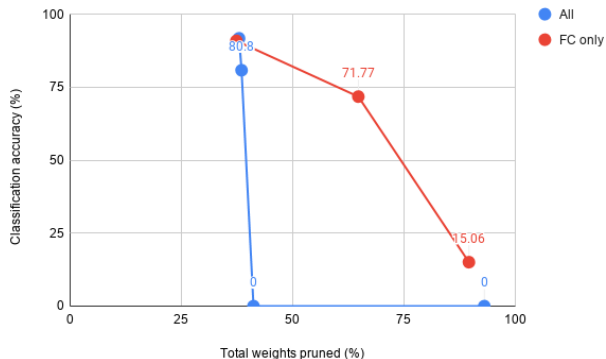Optimal scale factor per layer.

## Weight Pruning

- Biological brain phase from birth until mid-20s
- Network compression
- Weak weights get pruned $\rightarrow w\epsilon[-f, f] = 0$, f: pruning factor
- Calculations can get skipped
- Higher memory & power efficiency & inference performance
- Accuracy-Performance tradeoff
- Weight pruning amount varies per network
- Global pruning factor: Not a good idea!
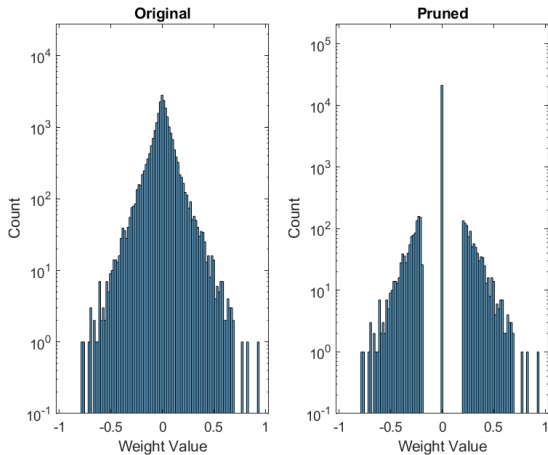
## Weight Pruning

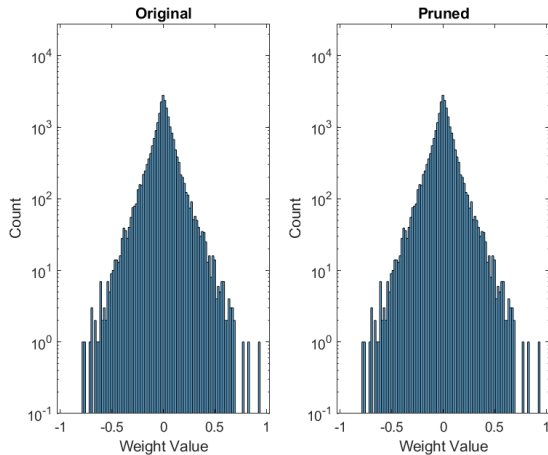| Layer | Test 1 | Test 2 | Test 3 | Test 4 | Test 5 | Test 6 | Test 7 |
|-------|--------|--------|--------|--------|--------|--------|--------|
| Conv1 (%) | 7.15 | 13.66 | 91.3 | 91.3 | 0 | 0 | 0 |
| Conv2 (%) | 13.82 | 26.9 | 95.83 | 95.83 | 0 | 0 | 0 |
| Conv3 (%) | 13.54 | 26.63 | 98.62 | 98.62 | 0 | 0 | 0 |
| Conv4 (%) | 15.32 | 29.99 | 93.14 | 93.14 | 0 | 0 | 0 |
| Conv5 (%) | 15.55 | 30.53 | 94.02 | 94.02 | 0 | 0 | 0 |
| FC1 (%) | 41.23 | 41.23 | 41.23 | 94.48 | 41.23 | 71.89 | 96.61 |
| FC2 (%) | 36.69 | 36.69 | 36.69 | 90.61 | 36.69 | 62.52 | 90.61 |
| FC3 (%) | 27.27 | 27.27 | 27.27 | 89.68 | 27.27 | 47.74 | 75.56 |
| **Total (%)** | 37.97 | 38.54 | 41.22 | 93.11 | 37.38 | 64.79 | 89.65 |
| **Accuracy (%)** | 91.74 | 80.8 | 0 | 0 | 90.87 | 71.77 | 15.06 |

# Weight Pruning



- Less Pruning → Higher Accuracy
- Convolution layers more prone to error
- Pruning also denoises - Low valued weights act as noise

# Weight Pruning



- Aggressive pruning
- High concentration of zeroes
- High compression factor
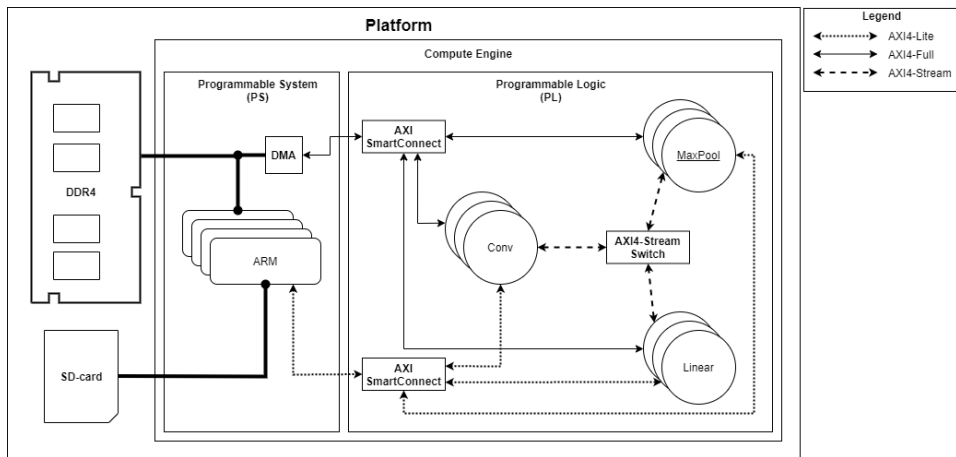- Severe absence of near-to-zero valued weights

# Weight Pruning



- Fine-tuned pruning
- Normal concentration of zeroes
- No discontinuation

# Architecture Design

## The platform

- Targeted for FPGAs
- Flexible & Versatile → easy transfer
- Scalable → multi-FPGA platforms
- Expandable → Easy adding of new layer types & accelerators
- Capable of running various CNN models
- Easy experimentation & development
- Minor to no code changes

# Platform Block Diagram

## Platform: Non-Volatile Memory

- Storage Medium
- Network model configurations
- Weights & Biases
- Class labels
- Input data, e.g. Images
- SD card - M.2 SSD (QFDB)
- External storage devices via Ethernet & JTAG

## Platform: Volatile Memory

- Main system memory: DDR
- DDR loaded using PS part
- No global BRAM module
- Accelerators: BRAM caches
- Accelerators responsible to load their BRAM & locate data on DDR
- BRAM caches are private to their accelerator

## Platform: Compute Engine

- Both PS & PL part utilized
- Bulk of computation on PL part through hardware accelerators
- Sophisticated work on PS part → Initialization, Data loading, Input data preprocessing, accelerator configuration & scheduling
- Network layers both software & hardware
- Convolution, Max-Pooling & Fully-Connected accelerators
- Accelerator driver only knows how to configure it
- Driver is reusable

## Platform: I/O

- Memory-Mapped I/O (MMIO)
- Streaming (AXI4-stream)
- BRAM MMIO - Not used

| Port bit-width | MMIO avg. cycles | Streaming avg. cycles |
|---|---|---|
| 32-bit | 62700922 | 65611580 |
| 128-bit | 15761270 | 16201797 |

40MB data - 40KB bursts

## Platform: Software

- Accelerator drivers - Abstract form, every accelerator implements same functions
- Scheduler
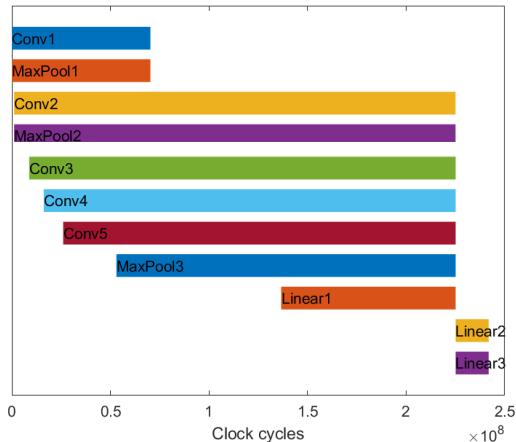- Application Logic
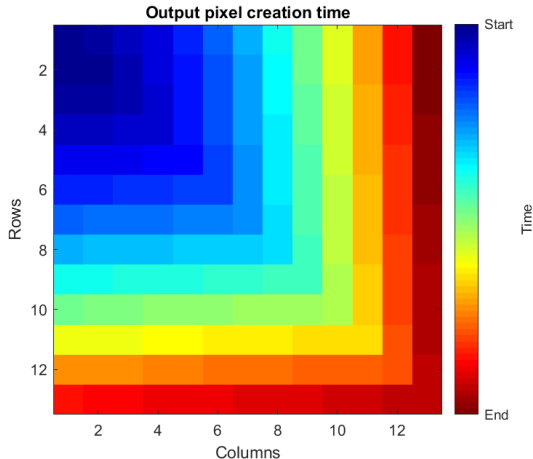- User Interface

## Platform: Software Flowchart

## Platform: Serial Scheduler



- Simple
- Best suited for debugging & validating
- About 90% of total inference time is consumed by convolution layers
- Possibility for deadlocks

# Platform: Layer Pipelining Scheduler



- A layer gets its input as soon as its previous generates a single output
- Accelerator instances needed as many as there are in the model
- Almost 3x speedup
- Decreases latency & increases throughput
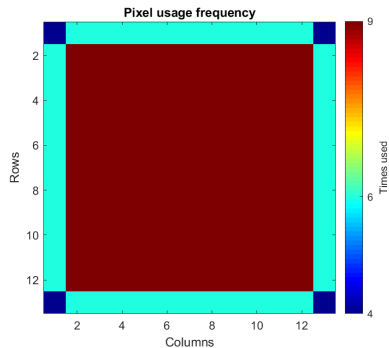- Accelerators need to support pipelining
- Relatively complex

## Output Pixel Creation Time



Output pixel creation time

- Outputs generated in specific order to become useful inputs
- Convolution layers use cubes of inputs
- Max-Pooling layers use squares of inputs

## Pixel Usage Frequency

# Larger BRAM caches required

## Platform: Multi-Inference Scheduler

- Multiple accelerator instances
- Multiple inferences in parallel
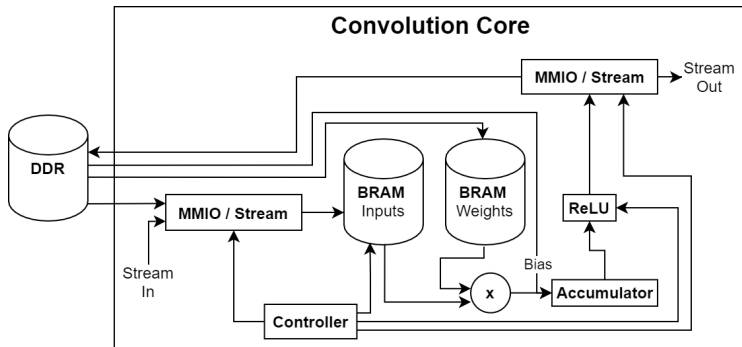- Increases batch size & throughput
- Possibility for deadlocks

## Platform: Image-Pipelining Scheduler

- Combination of Layer Pipelining & Multi-Inference
- Every layer handles a different input image
- Decreases latency & increases throughput
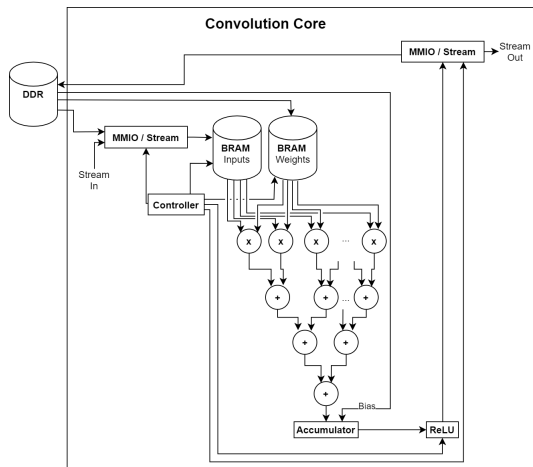- Possibility for deadlocks

# Accelerator Architectures

- Two versions: Simple serial & High performance
- Many have been tested
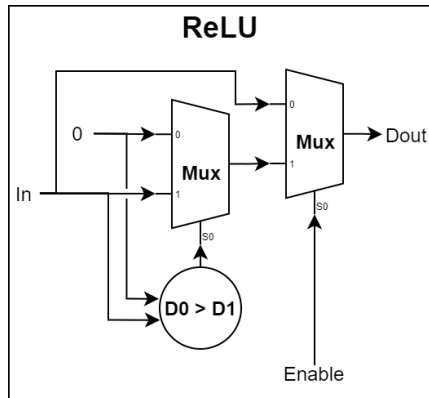- Convolution layer
- Max-Pooling layer
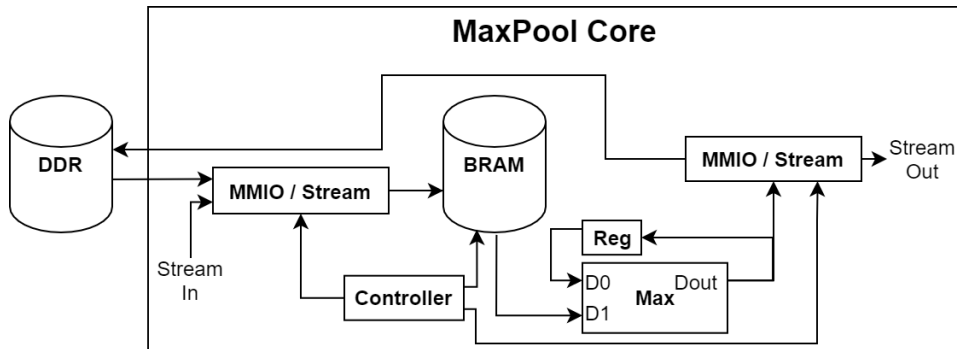- Fully-Connected layer

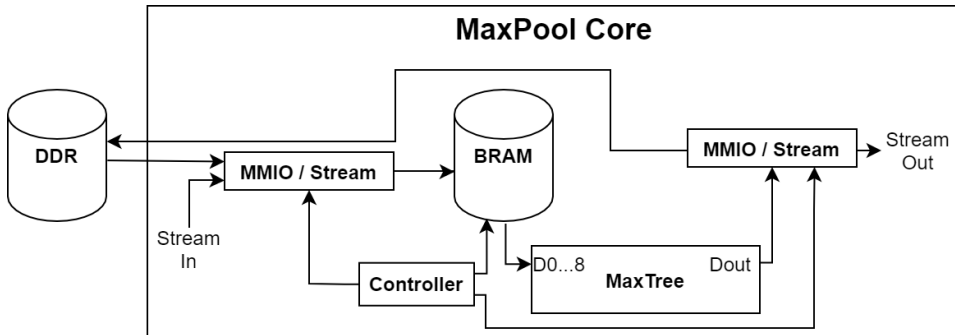## Convolution Accelerator

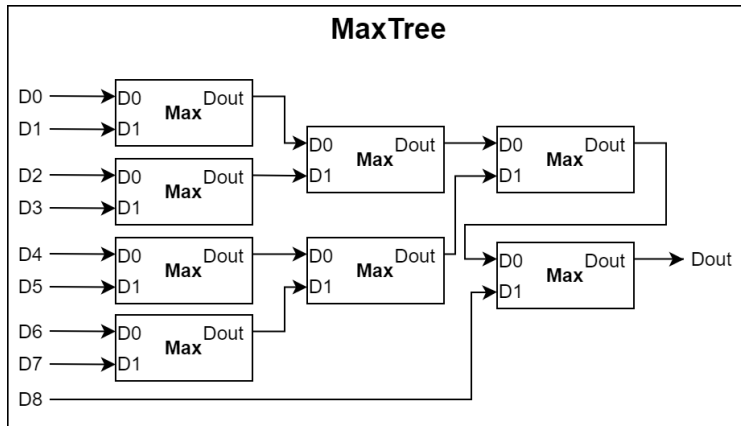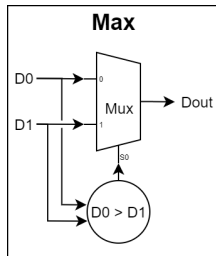## Convolution Accelerator

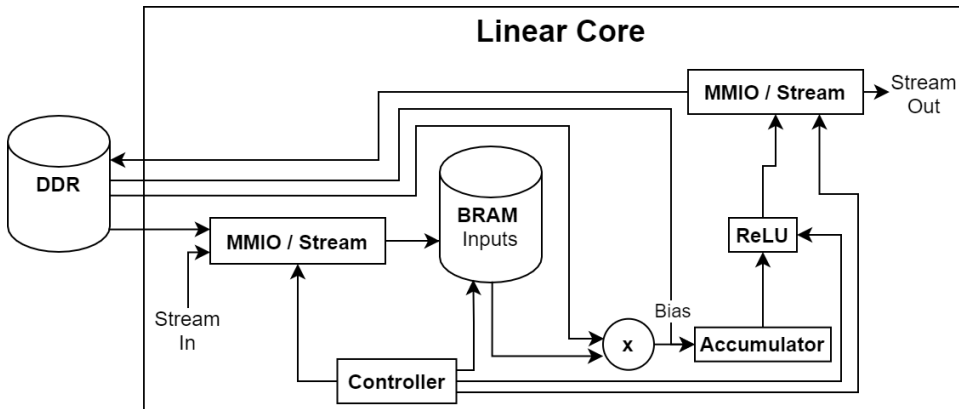## ReLU component

## Max-Pooling Accelerator
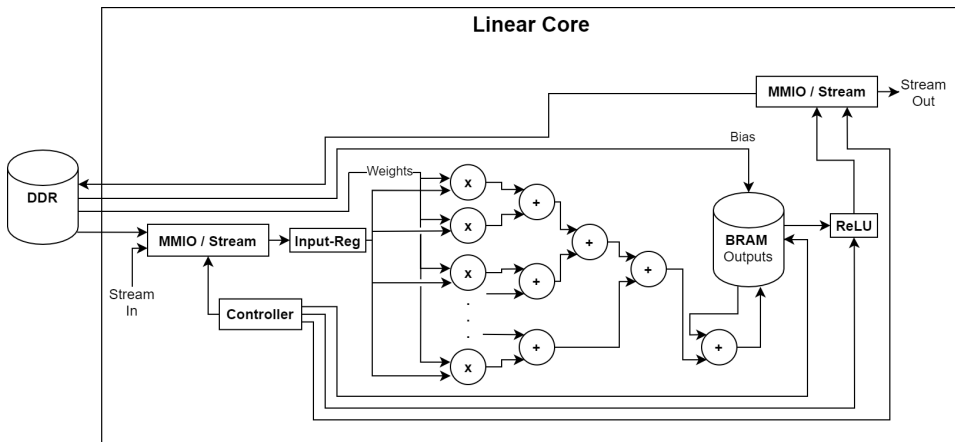
## Max-Pooling Accelerator

## Max & Max-Tree components

## Fully-Connected Accelerator

## Fully-Connected Accelerator

# FPGA Implementation

# Xilinx ZCU102 Evaluation Kit

# Tools Used: Xilinx Vivado HLS



- Now Vitis HLS
- High-level design using C/C++, SystemC, OpenCL
- Generates VHDL & Verilog HDL designs
- Directives
- C/C++ testbench
- C/RTL Cosimulation
- Synthesis Report

## Tools Used: Xilinx Vivado IDE



- VHDL & Verilog
- IP Integrator Tool
- Vivado HLS RTL designs
- Synthesis, Implementation & Download RTL designs
- RTL Simulators & Integrated Logic Analyzer IPs

# Tools Used: Xilinx SDK/Vitis IDE
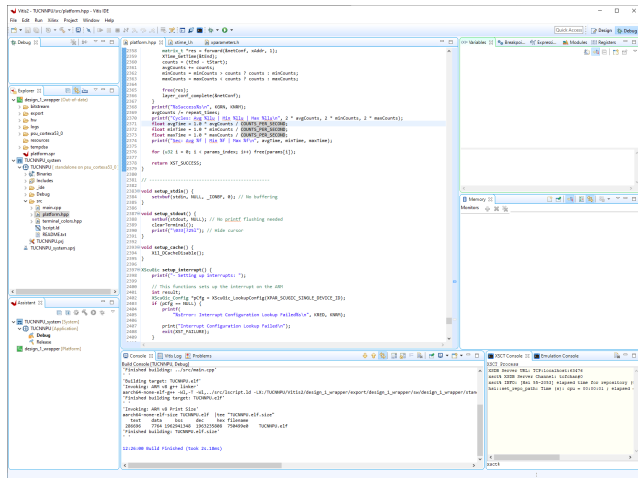


- Vitis IDE integrates SDK, SDAccel, SDSoC tools
- C/C++ IDE
- Application development for PS part
- PetaLinux & FreeRTOS
- Download bitstreams
- Debugging tools

Diploma Thesis

Introduction    Neural Networks    Related Work    Theoretical Modeling    Architecture Design    FPGA Implementation    **Results**    Conclusion

# Results

## Compared Platforms: CPU

### Intel i7 4710MQ

| | |
|---|---|
| **Cores / Threads** | 4/8 |
| **Max Turbo Frequency** | 3.5GHz |
| **TDP** | 47W |
| **Max Memory Bandwidth** | 25.6GB/s |
| **Lithography** | 22nm |

## Compared Platforms: GPU

### NVIDIA RTX-2060 Super 8GB

| | |
|---|---|
| **CUDA Cores** | 2176 |
| **Tensor Cores** | 32 |
| **GPU Memory** | 8GB GDDR6 |
| **Boost Clock** | 1650 MHz |
| **Memory Interface** | 256-bit |
| **Memory Bandwidth** | 448GB/s |
| **Power Consumption** | 175W |

## Compared Platforms: FPGA

### Xilinx CHaiDNN

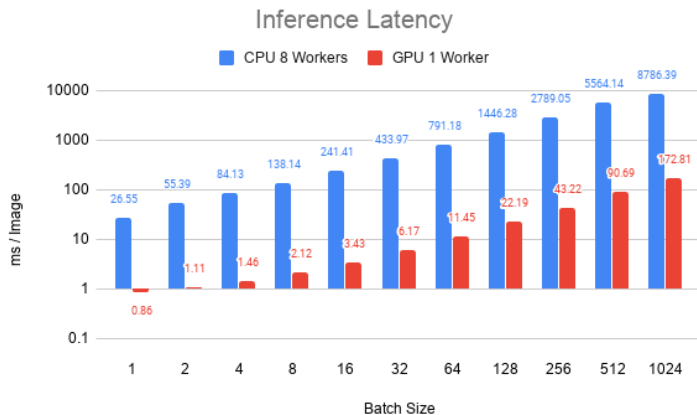| | |
|---|---|
| **PL/DSP Clock Frequency** | 250/500 MHz |
| **LUT Usage** | 59.1% |
| **FF Usage** | 27.66% |
| **BRAM Usage** | 74.12% |
| **DSP Usage** | 53.65% |

## Compared Platforms: FPGA

### Proposed Platform

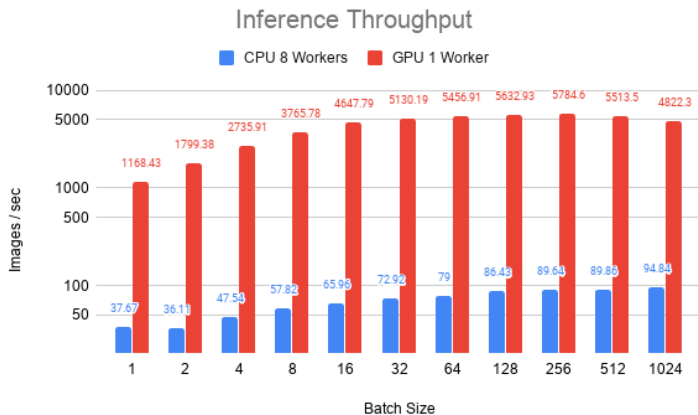| | |
|---|---|
| **Clock Frequency (MHz)** | 300MHz |
| **LUT Usage** | 7.34% |
| **LUTRAM Usage** | 2.05% |
| **FF Usage** | 4.03% |
| **BRAM Usage** | 7.51% |
| **DSP Usage** | 1.9% |

## CPU & GPU Performance

- Inference 2500 images
- Use all worker & batch-size combinations
- PyTorch pre-built pre-trained AlexNet
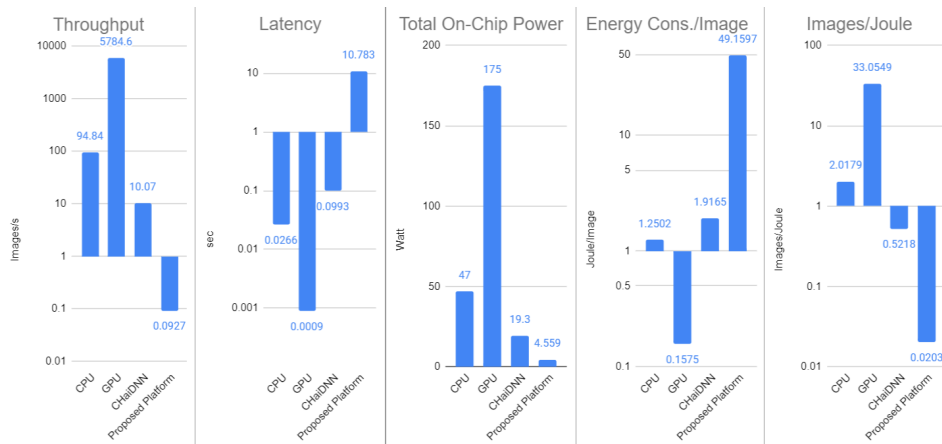
# CPU & GPU Performance: Latency



Inference Latency

# CPU & GPU Performance: Throughput



Inference Throughput

■ CPU 8 Workers   ■ GPU 1 Worker

## Final Performance

|  | CPU | GPU | CHaiDNN | Proposed Platform |
|---|---|---|---|---|
| **Clock Frequency (MHz)** | 3500 | 1650 | 250/500 | 300 |
| **Throughput (Images/s)** | 94.84 | 5784.6 | 10.07 | 0.0927 |
| **Throughput Speedup** | 1x | 60.9933x | 0.1062x | 0.001x |
| **Latency (s)** | 0.0266 | 0.0009 | 0.0993 | 10.783 |
| **Latency Speedup** | 1x | 29.5556x | 0.2679x | 0.0025x |
| **Total On-Chip Power (Watt)** | 47 | 175 | 19.3 | 4.559 |
| **Power Efficiency** | 1x | 0.2686x | 2.4352x | 10.3093x |
| **Energy Cons./Image (Joule)** | 1.2502 | 0.1575 | 1.9165 | 49.1597 |
| **Energy Efficiency** | 1x | 7.9378x | 0.6523x | 0.0254x |
| **Images/Joule** | 2.0179 | 33.0549 | 0.5218 | 0.0203 |

# Final Performance

# Conclusions & Future Work

## Conclusions

- Neural Networks need hardware acceleration
- Proposed platform provides an easy and structured methodology for scalable & expandable accelerator implementation
- Memory reduction is a necessity
- Further development $\rightarrow$ higher performance

## Future Work

- Quantization: better classification accuracy, K-Means clustering, Lloyd's, Pair and Quad compression, and Second Level Codebook
- Integrating the pooling layer into the convolution layer
- Pruning enabled accelerators
- Systolic arrays as their main compute engine
- Multiple accelerator instances

## Future Work

- Layer-Pipelining
- Bigger FPGA devices & multiple interconnected FPGAs (FORTH QFDB & CRDB).
- Monte Carlo Dropout for increased confidence of classification results
- Designs using VHDL & Verilog for resource & performance optimization
- CPU-FPGA partitioning

# Thank You!
Any Questions?