

# ECE 284

## Project Report

Fall 2024



*k-furthest-neighbors*

Anushka Chaudary      **A69029714**

Arjit Verma      **A69030364**

Chainika Shah      **A69030115**

Hariram Ramakrishnan      **A69030260**

Shail Vaidya      **A69030307**

## 1. Introduction

This project aims to implement a reconfigurable AI accelerator utilizing 2D systolic array architecture, enhanced with features such as sparsity-aware clock gating and Max Pooling, and optimized for compute time. The design is validated using a VGGNet model trained with quantization-aware training with the CIFAR-10 dataset.

## 2. VGG16 quantization-aware training

The **27th layer of the VGG16 model** was modified to have 8 input and 8 output channels in its convolutional layer, aligning it with the requirements for mapping onto an 8x8 2D systolic array. Additionally, we removed the batch normalization layer that typically follows the convolution, simplifying the structure to just "conv -> relu." By implementing these changes and training the model, we successfully achieved an accuracy of **92%**.

```
print(model.features[26])
print(model.features[27])
print(model.features[28])

QuantConv2d(
  256, 8, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False
  (weight_quant): weight_quantize_fn()
)
BatchNorm2d(8, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
ReLU(inplace=True)
QuantConv2d(
  8, 8, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False
  (weight_quant): weight_quantize_fn()
)
ReLU(inplace=True)
```

*Fig 1: VGG16 network with modified layers*

## 3. “Vanilla” Version

### 3.1. RTL Design

In addition to the mac\_array, L0 and OFIFO, we developed the RTL code for the SFU module to perform the accumulation and ReLU operations, and instantiated all the required modules at the corelet hierarchy. The corelet was then instantiated in the core hierarchy, along with the **pmem** SRAM, of size 512 words x 128 bits, to store the partial sums, and the **xmem** SRAM, of size 256 words x 32 bits, to store the input weight and activations. The FIFO depth for the L0 and OFIFO were also reduced to 16 to save area.

### 3.2. Testbench

The testbench, acting as the control unit, performs the tasks of loading weight and activation data from the TXT files into the SRAMs, and then loading the L0 with first the weights, followed by the activations. Consecutive loops were used as much as possible to maximize efficiency.

Once the computations are complete, the testbench loads the output partial sums to the OFIFO, and then to the PSUM memory. It also implements the cherry-picking algorithm in order to use the SFU and calculates the final output feature value and compares it to the expected output, delivering appropriate messages for the comparison. Here, in addition to the final output, the PSUMs are also compared to the expected PSUM, with results displayed accordingly.

### 3.3. FPGA Mapping

Table 1 shows the results of the FPGA mapping. We are able to see the area benefit gained through FIFO depth reduction, resulting in only 4132 registers needed for the entire 'Vanilla' version of the design.

## 4. Reconfigurable Version

### 4.1. Reconfigurable PE & Instruction Map

The PE was enhanced to enable output stationary mode computation. This required the addition of an extra instruction bit, with the new instruction mapping shown in Table 2. The MSB indicates if the present computation refers to weight or output stationary mode mapping, and the PE is ready to receive weight and activation inputs and modify the calculation to use the internal registers accordingly.

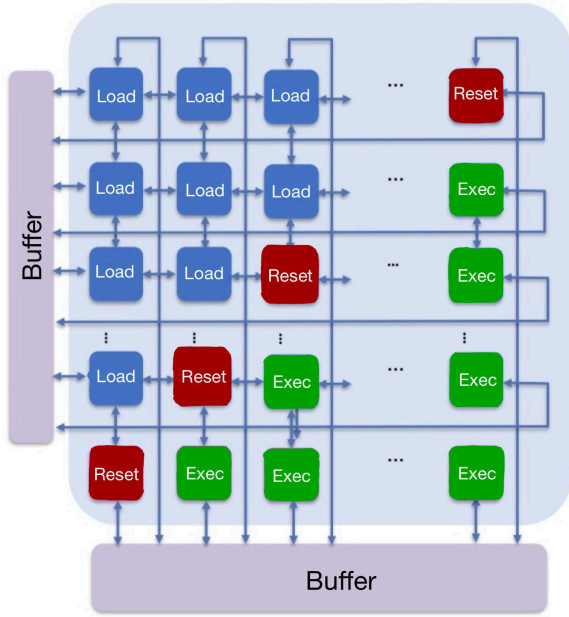
Cyclone IV FPGA Mapping of Vanilla Version	
Fmax	122.2 MHz
Switching Activity	20%
PVT	ss-1.2V-100C
Thermal Dynamic Power	18.76mW
Total Logic Elements	6886
Total Registers	4132
Total Operations	8x8x2 = 128
GOPs/s	16
TOPs/W	0.834

Table 1: FPGA Mapping Results

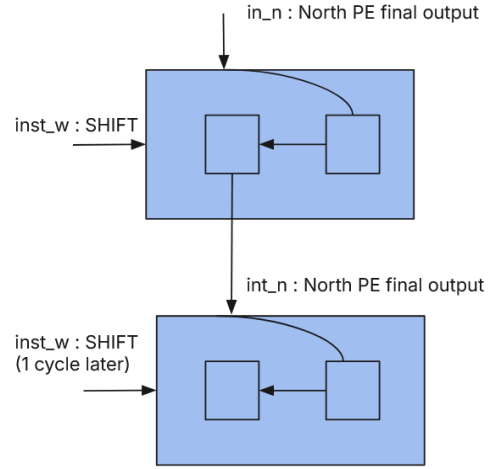
instr[2:0]	Description
000	IDLE
001	W_LOAD
010	W_EXEC
011	NOT_USED
100	RESET
101	O_SHIFT
110	O_EXEC
111	IDLE

Table 2: Reconfigurable PE Instruction Map





*Fig 3: Execution, Soft Reset & Kernel loading pipelined across PEs*



*Fig 4: SHIFT instruction: Input stored in current PE; Output sent to South PE*

### 5.1.2 Pipelined Output Shifting

To shift out the calculated output pixel values in the most efficient manner, each PE is designed to receive the SHIFT instruction as soon as the computation is completed, independent of the state of the neighboring PEs. As shown in Fig. 4, the output pixel value from the north is captured, while the current output pixel is sent south. This means the PE can begin shifting in the immediate next cycle after computation, without any IDLE cycle.

### 5.1.3. Testbench Optimization

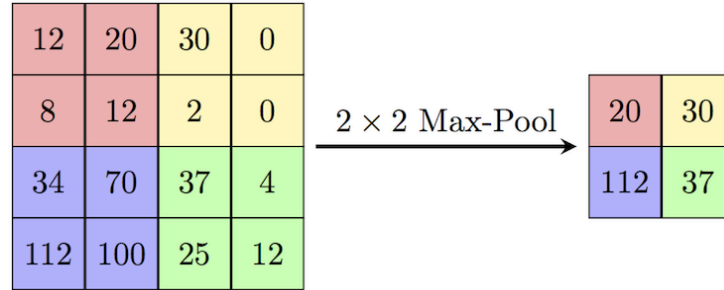
The testbench functions as the control unit in our project, and we have designed it to sequence instructions with no idle cycles, making full use of the pipelining implemented in the PEs. This leads to an overall cycle count of **586 cycles** in weight stationary and **47 cycles** in output stationary per output feature calculation.

## 5.2. Sparsity-Aware Clock Gating

The design has sparsity aware-clock gating to save dynamic power. This is achieved inside the PEs which accept a weight-zero and activation-zero flag, preventing the switching of the partial sum registers if either of these flags are high. The PE rows were enhanced to dynamically

disable the computation if the incoming weight or activation is zero, which prevents unnecessary toggling and dynamic power consumption.

### 5.3. Max Pooling Layer Implementation



*Fig 5: 2x2 MaxPool Operation*

The SFU is capable of implementing a MaxPool2d layer. When *max\_pool\_en* signal is asserted, the SFU computes the result of the MaxPool operation along with ReLU, with any kernel size and stride. This is verified for the MaxPool2d(2,2) present in VGG16, by enabling the MaxPool mode. The max pooling operation happens over 4 cycles of cherry-picked input data. The result is provided by the SFU as soon as the *max\_pool\_en* is pulled low.

## 6. Conclusion

We have presented all the details pertaining to the developed 2-D Systolic Array based AI accelerator design. The project completes all the requirements for the “vanilla” version, enabling weight-stationary mode computation and verification, as well as the reconfigurable version. This enables output-stationary mode computation which has also been verified with weight and activation data from VGG16 neural network. In addition, we have presented enhancements that are twofold - leading to overall compute time reduction through reduced cycle count, and additional features to help save power and implement other layers used in neural networks.