# Semester Project Report

# Real-Time Image Classification using Convolutional Neural Networks on Embedded Platforms using SoC / FPGA Convolutional Layer Accelerator

**Student:**                        Patryk Oleniuk
**Laboratory:**                     Processor Architecture Lab (**LAP**)
**Supervisors:**                    Mikhail Asiatici, Rene Beauchat
**Tools to train the CNN:**         Python, Tensorflow, Jupyter Notebook
**Algorithm Analysis:**             Python, C/C++
**Embedded System:**                Xilinx Zynq ZC-7006 on PetaLinux / BareMetal
**Tools for Embedded System:**      Sigasi, Vivado 2017, Xilinx C/C++ SDK 2017

**Abstract:**
Cameras are used everywhere as flexible sensors for numerous applications. For mobility and privacy reasons, the required image processing should be local on embedded computer platforms with performance requirements and energy constraints. Dedicated acceleration of Convolutional Neural Networks (CNN) can achieve these targets with enough flexibility to perform multiple vision tasks. A challenging problem is the design of efficient accelerators for that purpose, since some CNNs require a lot of repetitive image processing. In this project, I show and implement solutions to quantize the trained CNN as well as perform and optimize computations on an Embedded System. The design flow is evaluated by implementing the previously trained CNN to recognize facial emotions from face image implemented in python on a PC. The project explains the process of porting the CNN algorithm from python to C/C++ and then executing it on a ZYNQ FPGA board running with PetaLinux. The bottleneck of the CNN is the convolutional layers and that is why different solutions for that accelerator are analysed and some of them implemented on VHDL, connected to the embedded processor and their performance is tested.

# Contents

# List of figures

## List of tables

## List of Sample Codes

## List of Terms and Abbreviations:

**CNN**              – Convolutional Neural Network.
**Tensorflow**    – python toolbox for Machine Learning.
**Convolutional Kernel** – feature map, image or images representing a pattern to detect.
**FPGA**            – Field Programmable Gate Array
**SoC**              – System on Chip
**FFs**              – flop flops
**LUTs**            – Look Up Tables
**BRAMs**          – internal Block RAM (inside the FPGA)
**LSb**              – Least Significant bit
**MSb**              – Most significant bit
**ZynQ**            – ARM microcontroller used in ZC-06 board
**GPIO**            – General Purpose Input Output
**GPO**              – General Purpose Output
**GPI**              – General Purpose Input
**L1- L6**          – Indicators of Convolutional Layer Numbers of the CNN

## 1. Introduction

Cameras are used everywhere as flexible and cheap sensors for numerous applications. For mobility and privacy reasons, the required image processing should be local on embedded computer platforms with performance requirements and energy constraints that would allow to perform certain tasks, e.g. image classification. [1]

Currently the image classification tasks are performed very often by data structures called Convolutional Neural Networks (CNNs). The CNNs show very promising classification accuracy especially for image classification, but also for many other fields. The CNN properties and obtained accuracy are based not only of the chosen internal architecture, but solely on the training data and its amount, i.e. it requires a lot of trained, labelled data to train the model properly. The CNN used in this example was trained before and is designed to classify 6 emotions (happy, sad, surprised, angry, normal, scared) from 48x48 greyscale image containing frontal face. Executing single classification for this CNN on a standard PC on python takes few seconds. It was designed like any other typical CNN hence all the methodology and codes below <u>can be used directly</u> or with minor modifications to <u>port and accelerate any other CNN to classify different objects</u>. The CNN training, or the training acceleration is NOT the topic of this document.

The project shows the analysis of the computations behind a typical CNN classification (forward propagation) and the approach to port it into an Embedded Platform, so that the classification algorithm can be executed there. The process includes:

1. Migration of the trained model from python / tensorflow (the CNN is usually trained in some high-level language like python or Java) to C/C++.
2. Profiling of the C/C++ code to identify bottlenecks of the classification algorithm.
3. Quantisation of the coefficient to decrease computational power but keep good classification precision.
4. Case studies on 2 sample CNN accelerator designs using FPGA present in the Embedded Platform.
5. Implementation of the accelerator on the Embedded System.
6. Acceleration results.

A challenging problem is the design of efficient accelerators of CNN forward propagation, i.e. the "new image" classification task based on the previously trained CNN. The document describes a dedicated accelerator of the CNN that can achieve fast timing targets with enough flexibility to perform multiple vision tasks simultaneously, since some CNNs require a lot of repetitive image processing. In the project, I show and implement efficient solutions to quantize the trained CNN as well as perform and optimize computations on an Embedded System. The design flow is evaluated by implementing the CNN and executing it on a ZYNQ FPGA board running with PetaLinux. The bottleneck of the CNN are the convolutional layers and that is why different solutions for that accelerator are analysed and some of them implemented on VHDL, connected to the embedded processor and the performance of the best accelerator solution is tested.

## 2. The Trained CNN Model

The purpose of the CNN used in this example is to classify facial images into 6 emotion categories. It uses a lot of images' convolutions to find specific features in the images as well as non-linear operations to perform the classification task based on the detected features.

The input of the CNN is 48x48 greyscale image (array of pixels), each pixel traditionally represented by 8 bits (0-255). The output of the CNN is the array with floating-point classification probabilities. The overview of the system is shown below.
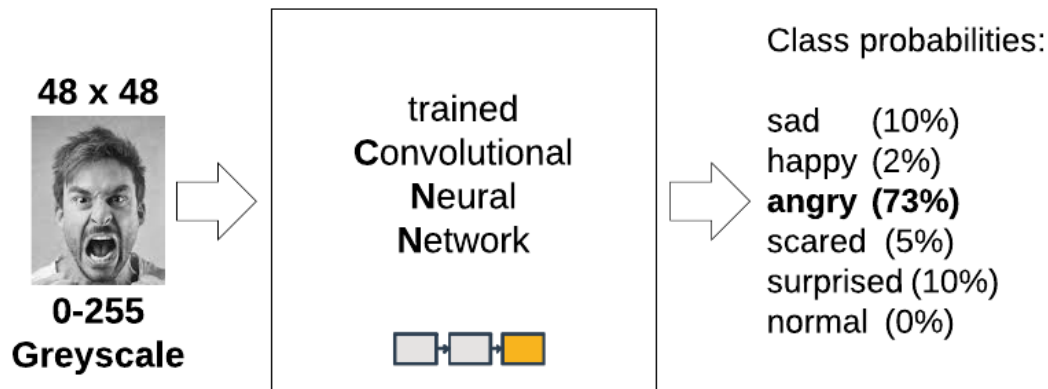


**Figure 1. General view of CNN classification.**

The CNN computational model contains multiple layers of very similar operations. There are multiple convolutional layers, which detect multiple features by performing the 2D image convolution, and fully connected layer at the end which maps the image pixels to the classified probability by simple matrix-matrix multiplication. The important thing to notice is that during the training period, each feature map and fully-connected matrix changes in order to minimise the error with the training data. The training of the CNN was described in [ADD REFERENCE] and this document is <u>not</u> to describe and/or accelerate the training of the CNN.

**2.1 Computational Graph – Details**

The CNN introduces 7 similar stages of image processing. Each of the stage is called a layer. In the CNN image processing domain, we could distinguish 2 main types of layers:
- **Convolutional** – the input image is convoluted in 2D with a trained feature image. It is usually much smaller than the image, e.g. input – 48x48 face image, feature map: 8x8 eye → this feature map will detect the position of the certain type of eye in the image). Each layer can have certain number of features (e.g. 22 8x8 feature maps). There can be also more input images (e.g. 22 images from the previous layer). <u>Then usually for each resulting image there is a "transfer" feature map so that each of the resulting images has a feature map that related to any of the input image.</u> As An example, for 22 input images and 22 output images, there are 22*22=484 feature maps. The resulting images for each of the input images are summed together to create the resulting image.
- **Fully connected Layers** – the input image is transformed into a matrix and multiplied to obtain another form of the data, e.g. transforming 22 6x6 images into 6x classification probabilities in the last layer of the aforementioned CNN).

To obtain better classification and training results, different image operations can be performed in between layers, i.e.:

- **ReLU** – Nonlinear Operator which "flattens" the negative part of the image (if the image is defined as signed). The function is applied to all the pixels in the image. The Function of the ReLU is shown below.
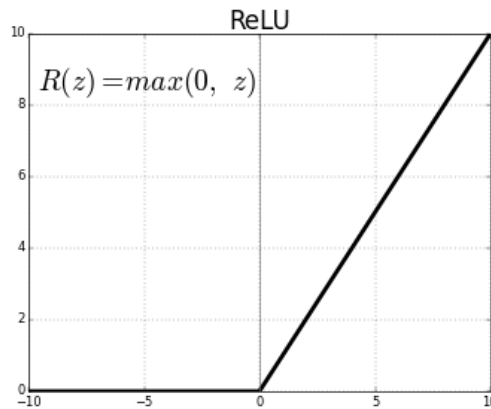


**Figure 2. ReLU rectifier function.**

- **Max Pooling** – It decreases the size of the image by choosing the pixels with the biggest value from the pool size. As an example, 48x48 image with max_pooling2x2 will result in 24x24 image with the biggest pixels in each 2x2 pool. The Max Pooling methodology is shown in the figure below.
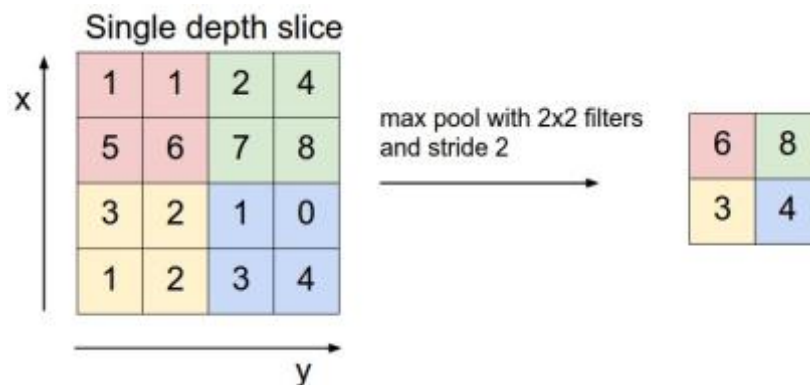


**Figure 3. Max Pooling 2x2 explanation.**

In the used CNN there are 6 convolutional layers using 8x8 or 4x4 feature maps, each followed by a ReLU and sometimes by max-pooling. Each layer recognises 22 feature maps. The last 7th layer is a fully connected layer which transforms 22 6x6 images into the classification constants. The detailed overview of the algorithm is shown below.
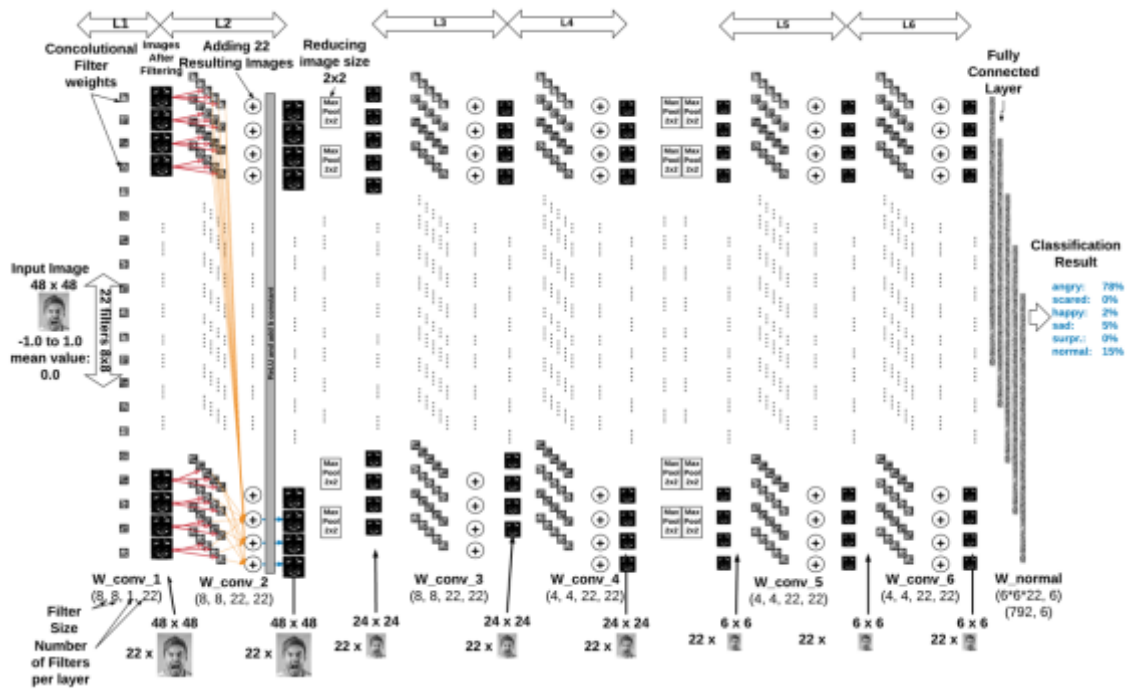
**Figure 4. Simplified Computational Structure of the CNN.**

Comments:

- Note the names and sizes of the Feature Maps (**W_Conv_1 (8,8,1,22)** ). The notation is consistent with the python/tensorflow implementation.

The summary of the 6 convolutional Layers parameters is shown on a table below:

**Table 1. Summary of each layers' properties.**

|                        | L1 | L2  | L3  | L4  | L5  | L6  | L7       |
|------------------------|----|-----|-----|-----|-----|-----|----------|
| n 2D conv to compute   | 22 | 484 | 484 | 484 | 484 | 484 | 6x6x22   |
| kernel length          | 8  | 8   | 8   | 4   | 4   | 4   | →6       |
| n img before layer (in)| 1  | 22  | 22  | 22  | 22  | 22  | matrix   |
| n img after layer(out) | 22 | 22  | 22  | 22  | 22  | 22  | multipli |
| img length             | 48 | 48  | 24  | 24  | 6   | 6   | cation   |

About the features and training:

To train all the feature maps automatically, a labelled dataset of more than 35 000 images set was used. One can observe, that some of the features have a real physical meaning and could represent part of a human face. An example of trained features and the images after applying them is shown below.
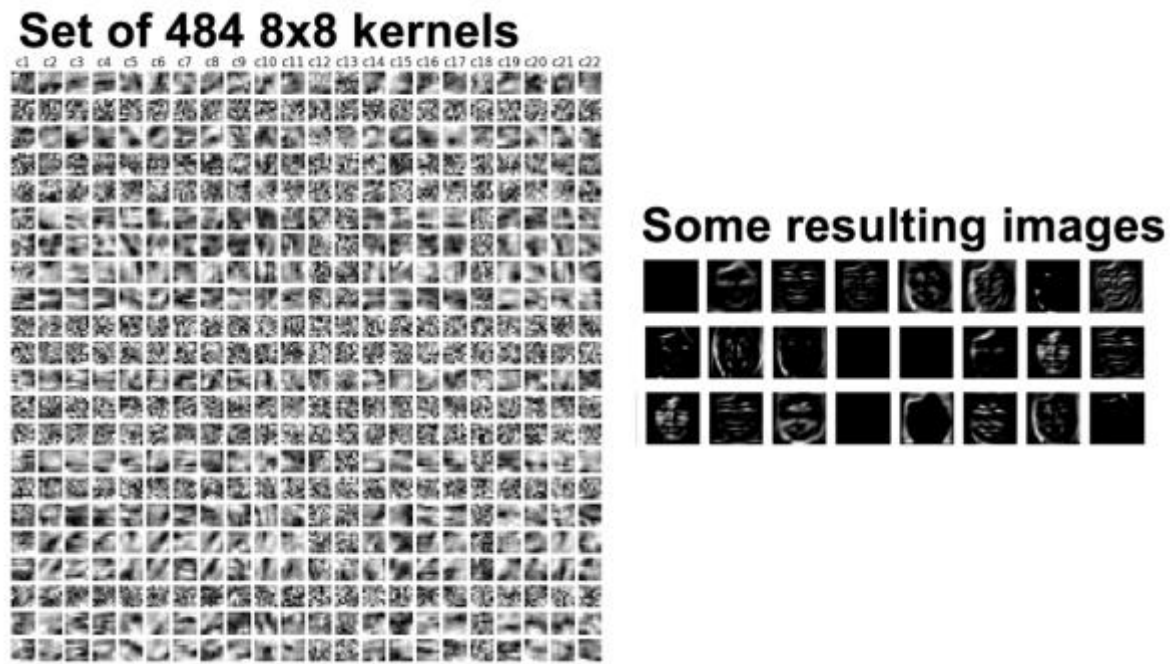
**Figure 5. Example of feature maps and resulting images.**

Both feature maps and the resulting images can be positive or negative, floating point. After performing a set of convolutions (22*22 in that case) on the corresponding input images (22 in that case, by row) with different kernels, columns of the resulting images(22 in that case) are added (pixel by pixel) to obtain again the number of resulting images after the layer (22 in that case). After the convolution, some the constant values are added to the (22 in that case) images and ReLU operation is performed on each pixel.

The code below is showing a single layer operation based on second layer of the CNN.

```python
## convolved2   = np.zeros((22,48,48)) # 22- number of inter-layer imgs, image 48x48
## convolved2_p = np.zeros((22,24,24)) # 22- number of inter-layer imgs, image 24x24
(maxpool2x2)

## LAYER 2 ##
for i in range(22):
    for j in range(22):  #accumulating convolutions for each row
        #single convolution operation (2 for loops inside)
        convolved2[i,:,:] = convolved2[i,:,:] + conv2d_single(convolved1[j,:,:], Wc2[:,:,j,i])
    convolved2_p[i,:,:] = maxpool2x2(convolved2[i,:,:])    #maxpooling 2x2
    convolved2_p[i,:,:] = ReLU( convolved2_p[i,:,:] + b_c2[0][i] )    #ReLU

## LAYER 3 operates on convolved2_p
```

**Code 1. Pseudo-code showing operations in layer 2.**

To fully understand the algorithm performed by a tensorflow library, a python and then C/C++ implementation has been written. To allow the reader to understand the algorithm structure and operations performed, both codes are attached to the document.

Important notices: Each layer contains a series of for-loops which then will be referred to in the accelerator section. The codes written in C++ and python are single threaded, therefore cannot use any other processor core, even if present.

8

**2.2 Profiling on a PC and Embedded Platform**

The python and C/C++ implementations of the algorithms were used to identify the classification time for each of them and understand the acceleration needs, if any. The algorithm execution times, averaged by 10x execution, are shown below.
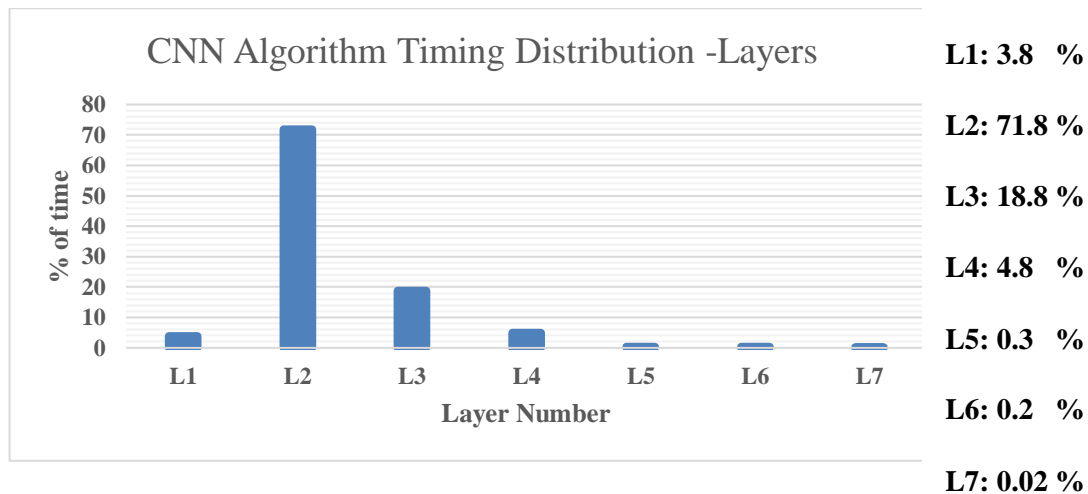
**Table 2. Timing results per single 48x48 face image**

| | | |
|---|---|---|
| 1. Macbook PC (3.4GHz Intel) in Python | –python / tensorflow | 3.55 s |
| 2. Macbook PC (3.4GHz Intel) in Python | –for loops (single thr)– | 126 s |
| 3. Macbook PC (3.4GHz Intel) in C/C++ | –for loops (single thr)– | 0.6 s |
| 4. Zynq PetaLinux(900MHz) in C/C++ | –for loops (single thr)– | 11.5 s |

Observations:
- The execution in python/tensorflow on a PC is very fast as for the high-level language, since the library is well optimized. Hence, even the python code (1.) is high-level, it's very close to the C++ (3.).
- Execution 3. Is the fastest, because it's executed on the fastest machine and the code is most low-level (C++).
- We could observe how much longer is the execution on the Embedded platform. For most of the processing tasks, the classification time of more than few seconds is unacceptable. Especially for real-time emotion recognition from the camera, which the CNN was designed to, is unacceptable.

To identify the bottlenecks of the algorithm, the profiling was executed on the C-code. The results for each layer are shown below.
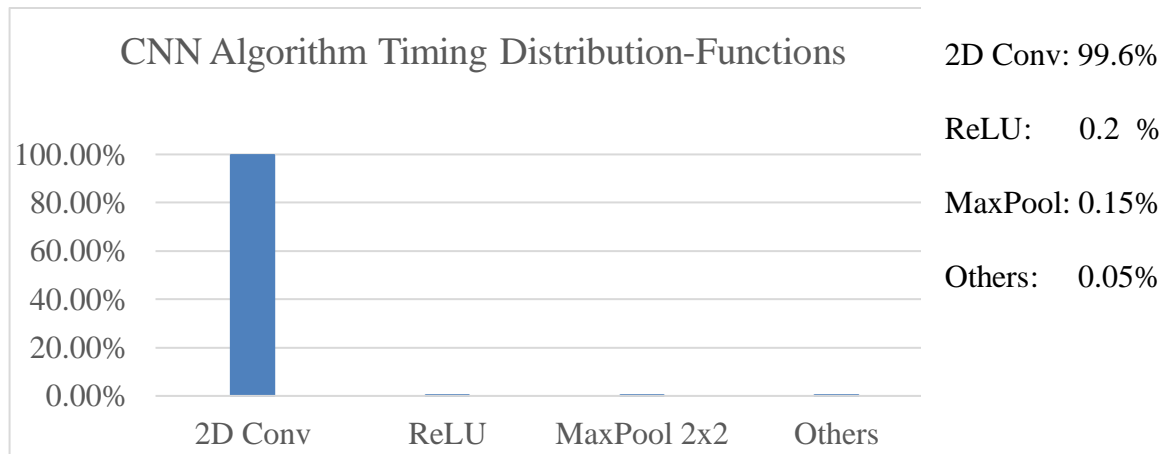


CNN Algorithm Timing Distribution -Layers

L1: 3.8 %

L2: 71.8 %

L3: 18.8 %

L4: 4.8 %

L5: 0.3 %

L6: 0.2 %

L7: 0.02 %

CNN Algorithm Timing Distribution-Functions

2D Conv: 99.6%

ReLU:    0.2 %

MaxPool: 0.15%

Others:    0.05%

**Figure 6. CNN Algorithm Timing Results.**

Observations:
- The most computationally expensive is Layer 2, since it has 22 48x48 input images (it has the highest parameters from each of the layers).
- By analyzing the detailed profiling info, the main computational task could be specified – the multiple 2D convolutions.
- The ReLU or max-pooling did not show to take more than 0.5% of the whole processing time.

Conclusions:
- Acceleration of the following functions should visibly shorten the classification time:
  o 2D image convolution, because it's the most computationally  heavy task executed with the most of time.
  o ReLU, because it's very easy to implement it in the FPGA as a multiplexer (if x<0  assign 0, else assign x).
- Due to the acceleration, there will be some delays to exchange information between the processor and the accelerator. Since we need to share a lot of data, it'd be valuable to design a solution that would also minimise this communication and therefore minimise the accelerator⟷processor time.
- Another effort of the data transfer between the float ⟷ fixed point, which would need to be executed every transfer to/from the accelerator.

## 2.3 Quantization

The usual requirement for signal processing using FPGA is quantisation, since they cannot easily handle the floating-point operations. For the purpose of the CNN acceleration, all of the operations to be accelerated have to be performed in the fixed-point domain. To clarify the needs for the CNN, the summary of the filter coefficients, as well as some example middle-results (the images in between the layers) content are shown below.

```
The filter coefficients (constant):
Max Wc1:   0.581492    Avg Wc1:   0.00039835 Min Wc1:   -0.942267
Max Wc2:   0.919815    Avg Wc1:  -0.0357445  Min Wc2:   -1.39372
Max Wc3:   0.817923    Avg Wc2:  -0.014101   Min Wc3:   -1.02589
Max Wc4:   0.604326    Avg Wc3:  -0.0391957  Min Wc4:   -0.990101
Max Wc5:   0.626381    Avg Wc4:  -0.0127036  Min Wc5:   -1.15402
Max Wc6:   0.624091    Avg Wc5:  -0.0218916  Min Wc6:   -0.839514
Max Wn6:   0.68207     Avg Wn6:  -0.0214262  Min Wn6:   -1.14417


The adder coefficients b (constant):
Max bc1:   0.0133187   Avg bc1:  -0.0378035  Min bc1:   -0.13341
Max bc2:   0.223921    Avg bc2:   0.0253297  Min bc2:   -0.243458
Max bc3:   0.307745    Avg bc3:   0.0571535  Min bc3:   -0.153931
Max bc4:   0.345763    Avg bc4:  -0.000638   Min bc4:   -0.35225
Max bc5:   0.297876    Avg bc5:   0.0265728  Min bc5:   -0.291342
Max bc6:   0.0384562   Avg bc6:   0.00243238 Min bc6:   -0.0225041


The middle results in between the layers h
(calc for 64 test images) :
Max h1:   0.424823    Avg h1:   0.0094548   Min h1:   0.0
Max h2:   2.96478     Avg h2:   0.0251095   Min h2:   0.0
Max h3:   7.43573     Avg h3:   0.267952    Min h3:   0.0
Max h4:   18.9776     Avg h4:   0.128376    Min h4:   0.0
Max h5:   23.1608     Avg h5:   0.653597    Min h5:   0.0
Max h6:   15.8007     Avg h6:   0.350184    Min h6:   0.0
```

**Table 3. Min, avg and max values of the coefficients and example inter-layer results.**

The aforementioned results were generated using python notebook and are based on the tensorflow CNN model. For each execution, the filter and b coefficients are constant and have positive or negative values. The inter-layer results are different for each image. To take a reasonable example to see what is the maximum and average value of the images generated after each layer, the values were calculated from 64 example images. The minimal value of the inter-layer images are 0, because of the ReLU function which zeros the negative part of the image after the resulting image is calculated.

The maximum value of the coefficients are usually very close to 1 and none of the values are bigger than 2.0. However, we can see that, especially for the layers 3-15, the maximum values are quite high and exceed 2.0 value. That is why, the fixed point value used in the computation should include some more bits for the high integer values. Another consideration is using different fixed point scheme for filter coefficients and the result calculation (which is considered to be implemented in the future). For simplicity, the same fixed-point scheme for all operations will be applied.

To understand how the quantisation will affect the classification results, the following experiment in C++ on float coefficients was performed:

1. Prepare the set of 64 test images.
2. Classify each of the image and obtain the classification probabilities for each class.
3. Quantise the coefficients with the certain, signed fixed-point scheme.
4. Perform the CNN algorithm on the same 64 test images, with the quantised coefficients from 3. and "flatten" each of the resulting imaged after each layer, according to maximum value that the signed point can store.
5. Compare the 64*6 class probabilities with results from 2. and calculate the mean and max error.*

\* Note, that all the probability values was taken into consideration when calculating the error. As an example, if the original classification probabilities were [ 0.1 0.1 0.5 0.2 0.0 0.1] and the result after quantisation was [0.0 0.0 0.55 0.45 0.0 0.0 ], the mean error is 0.108 (0.65/6) and the maximum error os 0.25 (0.45-0.2).

The quantisation sweep was performed for the following fixed-point schemes:
- 2.5 to 2.10 - max values: (-2,2)
- 3.5 to 3.10 - max values: (-4,4)
- 4.5 to 4.10 - max values: (-8,8)
- 5.5 to 5.10 - max values: (-16,16)

The fixed point fractional values from .5 to .10 represent steps from 0.03125 to 0.00097. The result of the experiment is shown on the figure and table below.



**Figure 7. Plots showing relative error vs. different fixed-point quantisation.**

CNN Quantisation - Mean Error / Max Error

INTEGER BITS

| FRACTIONAL BITS | 2. | 3. | 4. | 5. |
|---|---|---|---|---|
| .5 | 0.18 / 0.99 | 0.17 / 0.99 | 0.18 / 0.99 | 0.16 / 0.99 |
| .6 | 0.15 / 0.99 | 0.15 / 0.99 | 0.13 / 0.99 | 0.11 / 0.99 |
| .7 | 0.14 / 0.99 | 0.11 / 0.99 | 0.07 / 0.96 | 0.07 / 0.95 |
| .8 | 0.13 / 0.99 | 0.1 / 0.99 | 0.03 / 0.52 | 0.02 / 0.49 |
| .9 | 0.12 / 0.98 | 0.09 / 0.98 | 0.02 / 0.31 | 0.01 / 0.31 |
| .10 | 0.12 / 0.98 | 0.09 / 0.93 | 0.01 / 0.19 | 0.01 / 0.20 |

**Table 4. Detailed results from the CNN quantisation.**

We can observe, that the quantisation error drops dramatically when increasing the number of bits. The best trade-off between the number of bits and reasonable error is fixed point 4.8.

That is why all the future accelerators will be done for 12 bit fixed-point (4.8) per pixel.

### 3. Resource Available on the Xilinx SoC platform

To accelerate the algorithm, Xilinx ZC7006 platform has been chosen. It has multiple available resources, including a ZynQ processor – to run the algorithm in software - and FPGA which could accelerate it. The summary of its properties is shown below:

| Name | Number | Comments |
|---|---|---|
| DSPs | 900 | |
| BRAMs | 545 | Each BRAM at least 8kB |
| FFs | 437200 | |
| LUTs | 218600 | |

**Table 5. Properties of the ZC706 platform.**

The most important property of the platform is the number of DSPs, since for convolution a lot of multiplications-addition operations are done. Therefore, the accelerator design, especially for Solutiuon 2, will be fitted to the available DSPs. The availability of internal BRAMs is also important, since they can be used to store parts of image data for faster access.

### 4. Image Convolution Acceleration

When accelerating a signal processing algorithm, first there's a need to understand the algorithm behind. Algorithm for single layer convolutions, with multiple for loops, is shown below.

```
## convolved2   = np.zeros((22,48,48)) # 22- number of inter-layer imgs, image 48x48
## convolved2_p = np.zeros((22,24,24)) # 22- number of inter-layer imgs, image 24x24
(maxpool2x2)

## Wc2 = np.zeros((22,22,8,8)) # set of 484 kernels

## LAYER 2 ##
for i in range(22): -- 1st for loop
   for j in range(22):  -- 2nd for loop
      # img-kernel convolution operation
      for k in range(48):  -- 3rd for loop
         for l in range(48): --4th for loop
            # convolution result for single pixel with kernel
            for ker_w in range(8): --5th for loop
               for ker_h in range(8): --6th for loop
                  convolved2[i,k,l] += convolved1[j,k+ker_w,l+ker_h] *Wc2[i,j,ker_w,ker_h])
```

**Code 2. Pseudo-code showing multiple loops of the CNN convolution.**

Note: Some of the details, like 0-padding have been omitted to simplify the code above.

While analysing this code, there are multiple ways of acceleration. The main purpose in the FPGA acceleration is operation parallelization. As we can see, most of the data processed above is not dependent on each other, therefore can be easily parallelized. The operation used here is multiplication – accumulation. That is why probably the best solution will be to use DSPs present on FPGA platform. Since the same image pixel and same kernel pixels are used a lot, there's a huge

The main challange is to minimise pixel and kernel data transfer → re-use them inside the FPGA. For algorithm parallelization, 2 main solutions have been analysed and implemented.

**4.1 Solution 1**

The first idea to accelerate the kernel was to accelerate loop 5[th] and 6[th] from Code 2. Therefore, a pixel shifting logic has been designed and implemented. The simplified idea is shown in figure below.



14

**Figure 8. Block diagram showing the main idea of solution 1.**

Characteristics:
- 1 output pixel per clock cycle (after filling the buffer at the beggining).
- 64 Multiplications and additions done at a single clock cycle.

## 4.2 Evaluation of solution 1

Drawbacks:
- the pixels shifting logic is very complicated and takes a lot of resources
- the design takes 70% of the FPGA LUT/FF resources and most of it is for pixel shifting logic.
- The maximum clock frequency has been set as ~70MHz
- With the aforementioned characteristics, the estimated classification rate is every 0.5s.

The solution of the pixel shifting logic has been implemented and synthesized. After seeing the results of the synthesis and estimating the classification rate, a search for better solutions has been started, since the acceleration was not satisfactory.

## 4.3 Solution 2

The idea of this solution is inspired from the paper "Optimizing Loop Operation and Dataflow in FPGA Acceleration of Deep Convolutional Neural Networks [2]. Unlike the previous solution, the registers here do not take a lot of FPGA resources and it's using a different approach to the parallelization. The solution is using an additional level of image storage: internal FPGA BRAMs, which have 2 ports and are accessible both from the processor (by some automated interface) and from the accelerator. The paper presents very efficient pixel management technique and different level of parallelization.

**Figure 9. Slide taken from [2] showing the way of pixel re-utilization and parallelization. This example is based on 9-pixel kernel (3x3) and 9-pixel block (3x3).***

*Note: After 9 clock cycles, the convolution results of 3x3 block of pixels is available in the DSPs. Moreover, the same pixel shifter can be used to use it with another sets of DSPs and another sets of Kernels to compute another block results in parallel. To fully understand the solution, read the referenced paper.

When analyzing the paper, we can observe that the pixel shifting logic is complicated, but the operations on the shifter pixels are very generic and once understood and properly implemented for one block, can be re-used for all the pixel blocks.

Based on the method presented in the paper, similar simulations of pixel re-utilization for kernel size of 8x8 (64 length) has been created in MS Excel (poleniuk_CNN_method2.xlsx file, tab: registers). The results have shown a similar way to shift the pixels and allowed to think about a design of a pixel shifting logic.

Choice of the block size and levels of parallelizations
The decision about the exact level of parallelization was determined on basis of the number of DSPs. In particular, 900 DSPs are available. The CNN algorithm has image sizes of 48x48, 24x24, 6x6 → so it makes sense to make the single block size (that is computed as single convolution result after N clock cycles) as 6x6 and since the CNN has 22 images as input/output, so it makes sense to use 6 * 6 * 22 = 792 DSPs. This means, 88% of the available DSPs are utilised.

With the aforementioned level of parallelization, the algorithm together with parallelized operations is presented below.

```
## Wc2 = np.zeros((22,22,8,8)) # set of 484 kernels

## LAYER 2 ##
for i in range(22): --- 1st for loop
   for j in range(22): -- 2nd for loop
      for k in range(48): -- 3rd for loop     # partly parallelised (blocks by 6)
         for L in range(48): --4th for loop   # partly parallelised (blocks by 6)
            for ker_w in range(8): --5th for loop
               for ker_h in range(8): --6th for loop
                  convolved2[i,k,L] += convolved1[j,k+ker_w,L+ker_h] *Wc2[i,j,ker_w,ker_h])
convolved2[:,:,:] += b; # adding constant → also done in FPGA
ReLU(convolved2[:,:,:] ) # ReLU rectivier → also done in FPGA
```

**Code 3. Example code showing L2 computations. Parallelization are highlighted in red.**

### 4.3.1 Implementation Explanation

The aforementioned method has been implemented in VHDL for kernel size 8x8 and image sizes 48x48, 24x24 and 6x6.

#### 4.3.1.1   Pixel Shifting Logic

Based on the Excel file simulations of the 8x8 kernel and image pixels, the design of single pixel shifting block has been designed.



**Figure 10. Block Diagram of the single block pixel shifter.**

At first clock cycle, the initial pixels are loaded (either from the previous blocks or from the img BRAM) and the pixels start to shift, shifting in the NEXT_PIX which also comes from BRAM or previous block. At the end, all the loaded pixels appear in the output registers INIT_SAVE_REG(6 pixels to initialize the next block) an SAVE_REG(7 pixels saved and then multiplexed as NEXT_PIX for the next block) to be utilized.

The test bench showing pixel shifter is shown below.



**Figure 11. Test Bench showing basic operation of the pixel shifter.**

Each pixel shifter is responsible for 6 pixels and therefore we need 6 of them to shift 6x6 block image. The pixels are loaded from BRAM at the first cycle, and then only last pixel shifter takes the image from different image BRAMs and all others are taking the values from the previous block.

The block Diagram of all the pixel shifting logic is shown below.



**Figure 12. Diagram Showing details of pixels shifting logic.***

Each pixel shifter block needs to store $7 + 6 = 13$ pixels. In the pixel logic there are 6 pixel shifter, so all the logic need only $13 * 6 = 78$ registers storing 12-bit pixels.

**\*Notes:**
-Between single pixel shifter, there are combinatorial interfaces (grey) that are set of multiplexers that are interfacing the data between pixel shifters.
- each image shifter is taking it's data from it's own BRAM(first 8 cycles) or the previous pixel shifter(other cycles).
- The last (right) pixel shifter takes its data from different image BRAMs every clock cycle (it's a data source for the system).
- the image is split in between 6 image BRAMs without any repetitions. See Excel file for image BRAM structure.
- the control inputs of the pixel shifters are managed by the state machine logic.

The test bench showing the operation of the pixel shifting logic is shown below.



**Figure 13. Test Bench showing pixel shifting logic.**

### 4.3.1.2   Pixel Multiplication / addition

For the purpose of parallel addition and multiplication, a block of 792 DSP has been designed. Each cycle it's supposed to multiply the 36 input pixels by 22 different kernel pixels and accumulate it into separate DSP registers. It has a synchronous zero (zero the accumulation result) and accumulate enable signals. The block diagram of the design is shown below.



**Figure 14. Block diagram of the pixel accumulator block (792 DSPs).**

The accumulator design handles maximum and minimum value saturation and after each cycle updates the result at the output. When accumulate enable signal is turned off, the output result is still available.

### 4.3.1.3   ReLU Rectifier in Hardware

Before saving the pixel values to BRAM, the rectifier function has to be applied. The function is simply a hardware realization of $y = max(x,0)$ in signed 12 bits. If the MSb of the input pixel is 1(the value of pixel is negative then), all the results is zeroed. If not, the results is the same as input.

### 4.3.1.4   Zero Padding

To solve zero padding problem when performing the kernel operation, 3 pixel boundary is added in the image. That is why, the image size in the BRAM is 54 x 54 pixels. The values in the boundary should be kept at 0. The example image (with numbered pixels) is shown below, on basis on "IMG BRAM" tab in the excel file "poleniuk_cnn_method2.xslx".

**Figure 15. An 48x48image with zero-padding (54x54), nr of pixels used for further considerations.**

#### 4.3.1.5   BRAM Image Structure

The BRAM image structure is designed in a way, that during the first 8 cycles of the single block computation each pixel shifter can take the data only from single BRAM image (apart from the last one) and so that there is no repetitive data in BRAMs. This means, that there are 6 BRAMS and each BRAM cell stores 6 pixels (72 bits). Pixel bits are concatenated with each other.

**Figure 16. Single image BRAM structure (in 6 image BRAMS).**

Each consecutive 9 BRAM cells contain a row of image. Therefore, rows 0-8 of BRAM 1 contain 1st line(1-54) of the image, and then rows 9-18 contain 6th line of the image etc. For BRAM2, rows 0-8 contain 2nd line of the image and 9-18 contain 7th line of the image. Each image takes 81 rows of image BRAMs 1-6.

It's needed to store 22 full images; therefore, it takes 81* 22 = 1782 addresses. Each image is stored just after the previous one. Since both read and resulting images are stored in the same BRAM, the BRAM is divided into 2 sectors of the 22-image size and for different layers images are read and write in the swapped sectors. If the image has different size(24x24, 6x6), it takes smaller space (but the sector and size per image remains the same). The table showing R/W sectors for each layers are shown below.

21

| Layer | Read Addresses | Write addresses |
|-------|----------------|-----------------|
| L1 | 0-1781 | 1782-3563 |
| L2 | 1782-3563 | 0-1781 |
| L3 | 0-1781 | 1782-3563 |
| L4 | 1782-3563 | 0-1781 |
| L5 | 0-1781 | 1782-3563 |
| L6 | 1782-3563 | 0-1781 |

**Table 6. Addressing of different image sectors.**

### 4.3.1.6   BRAM addressing

Each BRAM has 2 ports. Each port works the same way. One port is connected to the processor, where the software can read and write images. The other port is connected to the accelerator, which reads the input image and writes the result.

Even though the BRAM needs to store 72 bits per cell, its size been configured to 128 bits to be able to handle it from the processor side easily (the available cell sizes are only powers of 2). The 72 image pixels are aligned to LSb and MSb's are padded with 0's. Since the address of the BRAM is byte aligned, keeping the previous notation, the real address that the image should be should be written or read is shifted left (padded 0's) 4 bits. As an example, cell address 81 (0x51) will become 1296 (0x510).

### 4.3.1.7   BRAM Image Writing

After the result of single block (6x6 pixels)  is ready, the image needs to be written back to the BRAM memory into the proper space, i.e. not from the beginning address, where the zero-padding zero begins, but where the real image begins. The image below shows the sample image result block (on the right) and the image BRAMs (on the left).



**Figure 17. Bram structure when writing: Unaligned pixels.**

Please note, that the results need to be written partly, into 2 cells in the BRAM. That is why, saving state machine first needs to read what is in the BRAM, save it, write the new values, and then do the same for the second address.

### 4.3.1.8   BRAM output registers for Read Access

To allow to obtain the timing between BRAM and the pixel shifting design, a register has been added in the output of the BRAM. Hence, after putting the correct address, the result is present in BRAM output after 2 clock cycles. The writing works as in the documentation, after 1 clock cycle.

### 4.3.1.9   Kernel BRAMs

The kernel BRAM structure is shown in the excel file "poleniuk_cnn_method2.xslx" in the tab "KERNEL BRAM". Since the kernel shifting logic had only been designed for 8x8 kernels, the 4x4 kernels for layers L3-L6 has been padded by 0's and used as 8x8.

Each kernel is 8x8 = 64 pixels and 22 kernel results are parallelized. Hence, for each kernel single pixel needs to be accessible in parallel. For instance, $1^{st}$ pixel of all 22 kernels for L1 needs to be available in parallel.

The kernel pixels are also 12 bits, and it's need 22 of them in one cell to be read in parallel. Therefore, $22 * 22 = 264$. The BRAM cell size needs to be the power of 2, so the BRAM cell size has been chosen as 512 bits. The kernel values are aligned to LSb and all the MSb are padded with 0's. The kernel BRAM structure is shown in the excel file "poleniuk_cnn_method2.xslx" in tab "KERNEL BRAM".

$1^{st}$ layer uses only 22 kernels but L2-L6 uses 484 kernels each. That is why, the number of BRAM kernel cells is $64 * 1 + 5 * 64 * 22 = 7104$. The 6 b constants for each layer, are stored at the end of the BRAM size. That is why, the final number of cells is $7104 + 6 = 7110$.

### 4.3.1.10  Accessing BRAMs from ZynQ processor

**It is possible to access each 32-bit part of the BRAM cell with the byte-aligned addressing. The address space is shown below:**

| Cell | Slave Interface | Base Name | Offset Address | Range | | High Address |
|---|---|---|---|---|---|---|
| ∨ 🖈 processing_system7_0 | | | | | | |
| ∨ ▦ Data (32 address bits : 0x40000000 [ 1G ]) | | | | | | |
| ▭ axi_bram_ctrl_0 | S_AXI | Mem0 | 0x4000_0000 | 64K | ▾ | 0x4000_FFFF |
| ▭ axi_bram_ctrl_1 | S_AXI | Mem0 | 0x4200_0000 | 64K | ▾ | 0x4200_FFFF |
| ▭ axi_bram_ctrl_2 | S_AXI | Mem0 | 0x4400_0000 | 64K | ▾ | 0x4400_FFFF |
| ▭ axi_bram_ctrl_3 | S_AXI | Mem0 | 0x4600_0000 | 64K | ▾ | 0x4600_FFFF |
| ▭ axi_bram_ctrl_4 | S_AXI | Mem0 | 0x4800_0000 | 64K | ▾ | 0x4800_FFFF |
| ▭ axi_bram_ctrl_5 | S_AXI | Mem0 | 0x4A00_0000 | 64K | ▾ | 0x4A00_FFFF |
| ▭ axi_bram_ctrl_6 | S_AXI | Mem0 | 0x4C00_0000 | 512K | ▾ | 0x4C07_FFFF |
| ▭ axi_gpio_0 | S_AXI | Reg | 0x4120_0000 | 64K | ▾ | 0x4120_FFFF |
| ▭ axi_gpio_1 | S_AXI | Reg | 0x4121_0000 | 64K | ▾ | 0x4121_FFFF |
| ▭ axi_gpio_2 | S_AXI | Reg | 0x4122_0000 | 64K | ▾ | 0x4122_FFFF |

**Table 7. Kernel and BRAM address space in AXI bus.**

### 4.3.1.11 State Machine

The block diagram of the full accelerator system and state machine is shown below.

**Figure 18. Simplified Block Diagram of Solution 2.**

The main purpose of the state machine is to manage the addresses of the BRAMs, based on layer number (which determines image size, sector to R/W, number of input images etc.) and control the pixel shifting logic, DSPs and writing the result. The simplified state machine diagram is shown below.

**Figure 19. Simplified State machine Diagram of Solution 2.**

The explanation for each of the states and different counters are shown below:
counters:
**Nr_pix** – counts from 0 to 63, when the DSPs are accumulating the kernel result.
**Nr_img_saved** – from 0 to 21, when saving the images needed to iterate through the DSP results.
**Block_nr (block_row and block_col)** – the results are computed in 6x6 blocks and each image has 8x8 blocks (for 48x48 image), 4x4 blocks (for 24x24 image) or 1x1 block(for 6x6 image).

States:
**Wait_start** – waiting for start signal. Connected to GPIO.

**CALC_LAYER** – wait for the result, change the kernel pixels and shift the pixel values.

**Other kernels to add to the same block ?** – depending on the layer, after computing single results for 22 kernels we need to iterate to another set / column of kernels(L2-L5) or not (L1).

**Add_b_const** – puts the kernel address to point to the b constant for the current layer and make them accumulate it in the DSP (the pixels at this point are multiplexed to be 1.0, and therefore the DSP_value += pixel_value * b_const(from kernel BRAM), pixel value = 1.0 so DSP_value += b_const.

**Save_6x6_block** – reads what in the current cell , writes the partial value from the DSPs, reads another cell and writes another partial result from the DSPs.

**Finished** – when all the image blocks(8x8 for L1-L2, 4x4 for L3-L4, 1x1 for L5-L6) are saved, the done flag is put to 1.

The test bench for the whole accelerator design, as an example for L1, is shown below:



**Figure 20. Test bench of the accelerator logic.**

#### 4.3.1.12 Block Design with ZynQ processor

**To be able to load the kernel coefficients and images, each BRAM ports 1 has been connected to the BRAM controller (connected to AXI bus and ZynQ) and BRAM ports 2 to the accelerator. The accelerator start, reset and layer_nr signals are managed by connecting it into a GPO, as well as the done signal to GPI.**



**Figure 21. System block diagram with BRAMs and ZynQ processor.**

### 4.4 Evaluation of Solution 2

The aforementioned solution had been synthesized in Vivado suite with standard implementation parameters. The details of the implementation are shown below:

**FPGA Frequency : 110MHz**

**Figure 22. Place & Route details of Solution 2.**

Observations:
- Most of the design is pixel accumulators-DSPs (yellow) and the pixel shifting logic takes only small part of it.
- 792 DSPs are utilised.
- 896kB of BRAM space is used.

To check the timing of the solution, a ZYNQ software has been used to trigger the start of the accelerator for different layers. The timing results are shown below.

| | Software timings[s] | Accelerated timings [s] | FPGA accel. (xNR) | posiible when maxpool on FPGA (xNR) |
|---|---|---|---|---|
| **python tensorflow** | 3 | | 149 | 48051 |
| **Mac C/C++** | 0.6 | | 30 | 1508 |
| **ZynQ Petalinux** | 11.5 | 0.02 | 571 | 28212 |

| | before | | now | future |
|---|---|---|---|---|
| **Classifications per second** | 0.1 to 2 | | 50 | ~500 |

**Figure 23. Timing Results After Acceleration of Solution 2. \***

*Only pixel results after L1 has been validated, the mean error (comparison to the non-quantized C software CNN execution) after converting to float is smaller than 3%.

## 5. Results and Conclusions

### 5.1 Acceleration ratios

One can clearly observe on Figure 23 that for the current acceleration results the 2D convolution is not a bottleneck any more. Thanks to multiple DSP utilization and smart pixel shifting logic the amount of time to compute convolutions and image accumulations results was reduced massively. Unfortunately, there is a max-pooling algorithm to be performed between L2 and L3 and between L4 and L5. This algorithm is not implemented in hardware, therefore all the images need to be transferred from BRAM to the operational memory, the algorithm is executed and the results is transferred back to BRAMs. The transfer between L2 and L3 is much bigger (and therefore much longer) because the image size there is 48x48 – after maxpoooling 24x24, transfer of 22 images. The transfer between L4 and L5 is much smaller – BRAM -> mem : 22 images, 24x24 , mem-> BRAM 22 img 6x6 each.

If we reduce the bottleneck of MaxPooling (implement it in hardware), the acceleration results can be obtained much faster, i.e. at least every ~3ms.
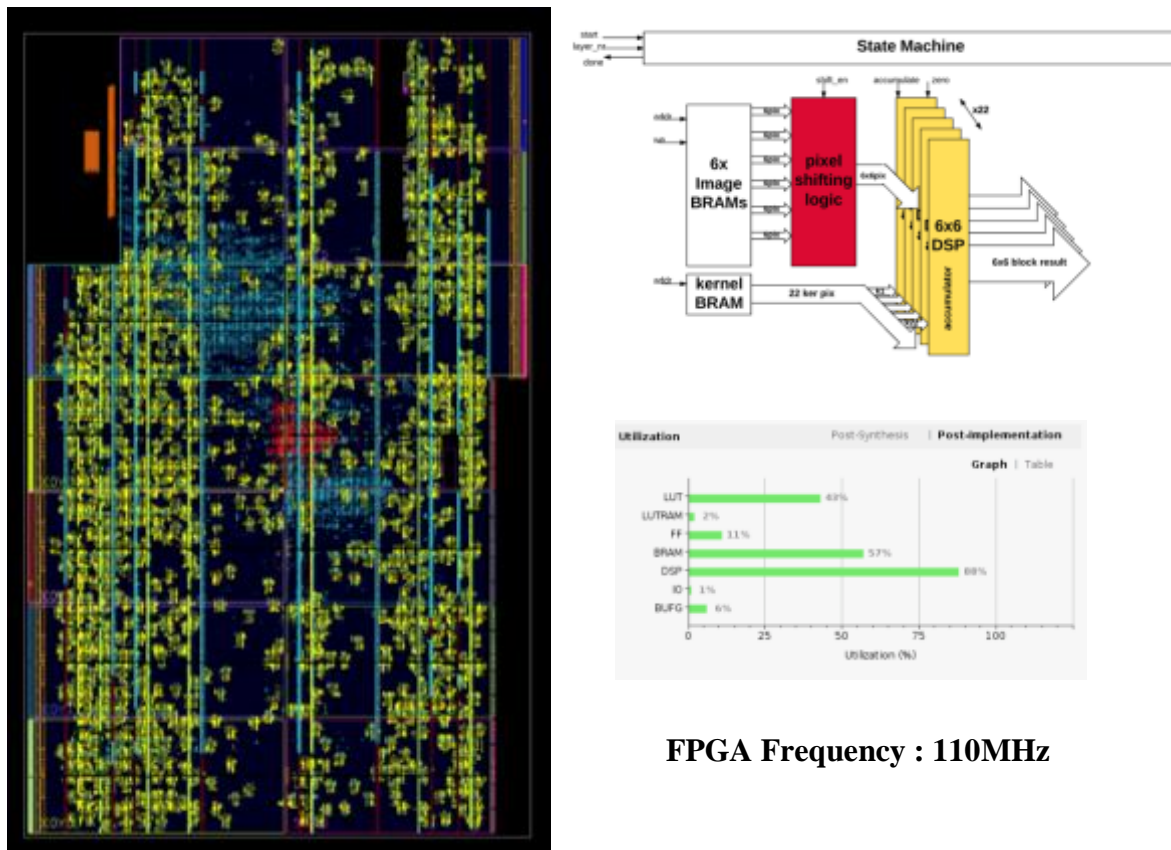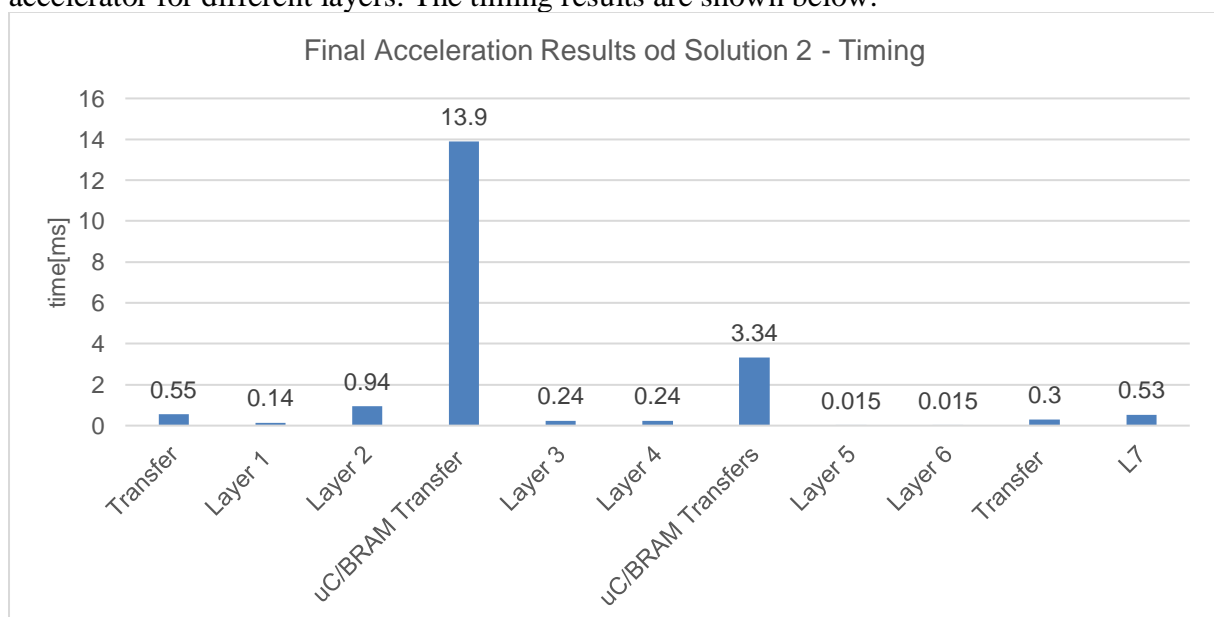
## 6. Challenges and Problems

### 6.1.1 Not Solved Problems

#### 6.1.1.1   Pixel results of L2 .. L6 are not fully validated

Pixel values for L1 has been fully validated – the values in the memory after accelerating L1 with a test images are consistent with the siftare-only version of the CNN algorithm within 3% of error which is cause by computational limitatiopns of the DSPs. However, the pixel values for L2-L6 are not matching with the software execution of CNN algorithm, therefore the result of the classification using the accelerated design is wrong. Because of the limited timeframe of the semester project, the solution to this problem could not be found.

**Pixel saving in L2-L6**
Pixels in the memory are saved in the correct place, with respecting the zero-padding etc.

**Wrong Pixel Values Debugging in L2-L6**

Probably there is a mistake in the  addressing state machine when iterating through the input images

- ❖ Very complicated design, debugging FPGAs
- ❖ Didn't manage to fix it within the given time frame,
- ❖ Would need about 2-4 more weeks of hardware debugging with test patterns.

### 6.1.2 Solved Challenges:

### 6.1.2.1    Analyzed the computational graph of the CNN.

The computations of the CNN has been fully understood in the means of single * and +. The difficulty was mainly to understand the high-level functions used in python / tensorflow and to map it into the C code executable efficiently on the embedded platform.

### 6.1.2.2    Profiled the CNN algorithm

After implementing fully functional classification algorithm in C, it was possible to precisely profile which operation was the bottlenech of the classification. That is why it was discovered that the 2D image convolution take more than 99% of the classification time.

### 6.1.2.3    Made the quantization analysis

Generated coefficients with different quantization levels and executed the algorithm with "flattening functions" to simulate the DSP behaviors. Then compared the generated classification errors in comparison to the "not quantized/flattened" coefficient solutions.

### 6.1.2.4    Reviewed different accelerator architectures.

During the project, multiple designs of convolution acceleration has been reviewed and the solutions appropriate to the platform, the CNN and the level of complexity for the semester project has been chosen.

### 6.1.2.5    Implemented and reviewed simple Solution 1

When implementing the simplified solution 1 in VHDL it was observed, that the design cannot efficiently accelerate the CNN very much because of the size of internal registers logic and prone utilization of DSPs. That initiated search for better acceleration architectures.

### 6.1.2.6    Implemented and reviewed Solution 2

The state-of-art solution 2 has been implemented for kernel size 8x8 and all the layers, which allowed to accelerate the CNN classification algorithm very efficiently by utilizing BRAMs and shifting the pixels in a very generic manner, therefore it could be implemented using much less logic cells of the FPGA than Solution 1 and so also more accumulator logic could fit into it.

### 6.1.2.7    Python Software to Generate the Constants in *.h file

After the CNN is trained in python, there is a need to transfer the trained CNN kernels and other coefficients (like B constants) and test faces into C code arrays. The python code has been prepared that could automatically convert tensor flow coefficients into C multi-dimensional arrays.

### 6.1.2.8   C libraries to transfer to / from 12-bit fixed point

To transfer the data to/from the BRAMs, functions has been written to transfer 12 bits in the proposed fixed point format from/ to double floating point.

### 6.1.2.9   C Libraries to control the accelerator

The functions has been written to reset, stop and check done flag of the accelerator (pins connected to different GPIOs).

### 6.1.2.10  C Libraries to transfer images to/from the BRAMs

The functions has been written to take the double image (array) and put it in the BRAM (align the 12 bit pixels etc. ). The same the other way, the function to read a single image (return array of double) from the BRAM address has been written.

### 6.1.2.11  Summary of Code Written:

- ~800 lines of python for:
  - analysis
  - quantisation
- ~3000 lines of C/C++
  - BRAM control
  - accelerator control
  - image manipulation
  - profiling
- ~ 2500 lines of  generic  VHDL:
  - implementation of Solution 1 and test benches
  - implementation of Solution 2 and test benches

### 6.1.2.12  Pixel values for L1 fully validated

The values in the memory after accelerating L1 with a test images are consistent with the software-only version of the CNN algorithm within 3% of error (in comparison which is probably caused by computational limitations of the DSPs.

## 7. Next Steps

### 7.1 Make it fully work with the addressing.

The L2-L6 pixel results still have to be debugged and the result fully validated.

### 7.1 Implement Max Pooling in hardware

Currently, the main classification speed bottleneck is tr image transfers to and from the software to do the maxpooling. Since the maxpooling algorithm is not very complicated, it should be implemented in hardware and then the classification result can be reduced into ~2ms (the previous timings subtracted by internal transfer times for maxpooling).

### 7.2 Make 4x4 kernels pixel shifting logic

Currently the pixel shjifting logic is working only for 8x8 kernels and 4x4 kernels are padded by 0's and treated as 8x8 kernels. The number of computations required for computing 4x4 kernels is much smaller (single block result in 16 not 64 cycles) and therefore the obtained acceleration could be up to 8 times faster for L4-L6.

### 7.3 Adding non-equal fixed-point scheme for filter coeff. and the inter-layer images.

Since the kernel coefficients values are from -2.0 to 2.0, their fixed point scheme could be adjusted and more bits for fractional points could be used. For the images, to obtain better classification accuracy and not "flatten" the result in some cases, it'd be better to have more bits in the integer part of the image. Further simulations needs to be dine for that purpose.

### 7.4 Saving state machine can take less cycles

Currently the saving state machine used 6 clock cycles to save single 6x6 block result (Read - 2clk and Write -1 1clk, for 2 addresses). This can be reduced to only 2 cycles by using byte-aligned write enable signal in BRAMs.

## 8. References

1. Maurice Peemen, et al *Memory-Centric Accelerator Design for Convolutional Neural Networks,* Department of Electrical Engineering, Eindhoven University of Technology
2. Yufei Ma et al., Optimizing Loop Operation and Dataflow in FPGA Acceleration of Deep Convolutional Neural Networks, Arizona University.
3. DianNao: A Small-Footprint High-Throughput Accelerator for Ubiquitous Machine-Learning