# Final Project:

# CNN Accelerator for MNIST Classification

20201082 H. Kim, 20211184 H. Yang

## I. Used Quantization methods
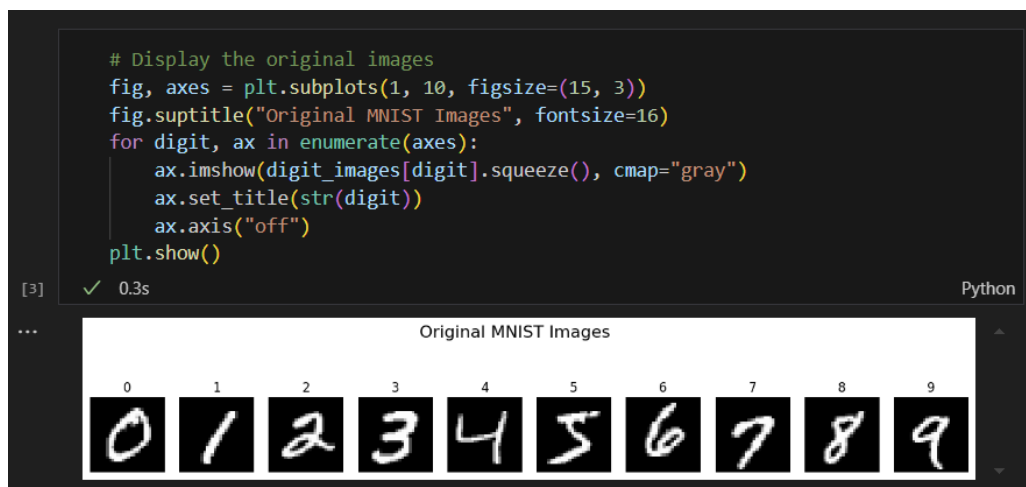
### 1. Binary Neural Network (BNN)

In this project, binary neural network (BNN) was used for better hardware utilization. To do so, we binarized weights and activation to -1 and 1 with `binary_quantization` function as shown below, based on the value of threshold.

```python
11  def binary_quantization(input_array, threshold=threshold, binary_range=(-1.0, 1.0)):
12      binary_min, binary_max = binary_range
13
14      if isinstance(input_array, torch.Tensor):
15          input_array = input_array.numpy()
16
17      if isinstance(threshold, torch.Tensor):
18          threshold = threshold.item()
19
20      quantized_array = np.where(input_array >= threshold, binary_max, binary_min)
21
22      return torch.tensor(quantized_array, dtype=torch.float32)
```
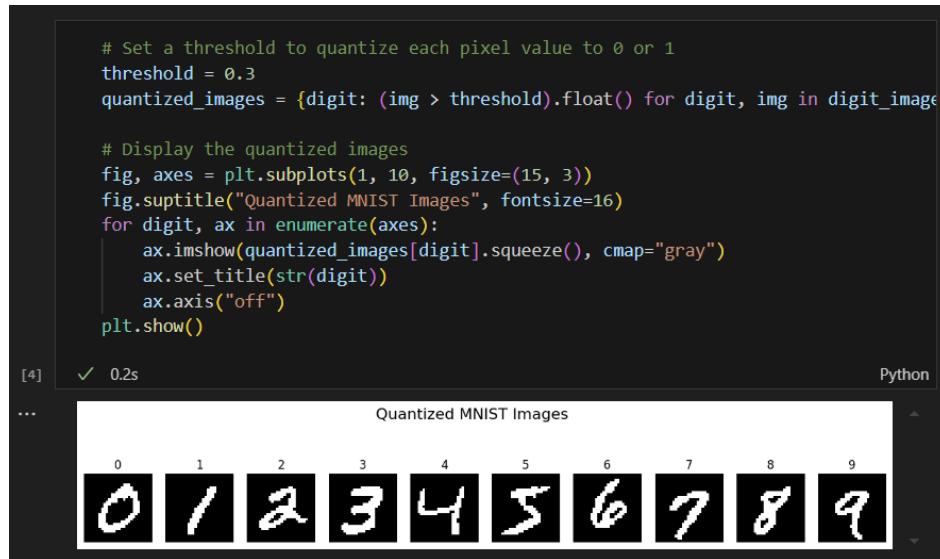
*Figure 1 Binary Quantization*

### 2. Input Binarization

Here, the value of threshold is set by visualizing each label.

```python
# Display the original images
fig, axes = plt.subplots(1, 10, figsize=(15, 3))
fig.suptitle("Original MNIST Images", fontsize=16)
for digit, ax in enumerate(axes):
    ax.imshow(digit_images[digit].squeeze(), cmap="gray")
    ax.set_title(str(digit))
    ax.axis("off")
plt.show()
```



*Figure 2 Original MNIST Data*

With simulation, the best threshold for our model was 0.2.

```python
# Set a threshold to quantize each pixel value to 0 or 1
threshold = 0.3
quantized_images = {digit: (img > threshold).float() for digit, img in digit_image

# Display the quantized images
fig, axes = plt.subplots(1, 10, figsize=(15, 3))
fig.suptitle("Quantized MNIST Images", fontsize=16)
for digit, ax in enumerate(axes):
    ax.imshow(quantized_images[digit].squeeze(), cmap="gray")
    ax.set_title(str(digit))
    ax.axis("off")
plt.show()
```

[4]   ✓  0.2s                                                                    Python

...



*Figure 3 Tested Quantized MNIST Images from threshold 0.1 to 0.9, and concluded that 0.2 shows the clearest image.*

## 3. Weight Binarization / Activation Binarization

```python
26    # Binary Activation (Applied STE for sign function)
27    class BinaryActivation(torch.autograd.Function):
28        @staticmethod
29        def forward(ctx, input):
30            output = torch.sign(input)
31            ctx.save_for_backward(input)
32            return output
33
34        @staticmethod
35        def backward(ctx, grad_output):
36            input, = ctx.saved_tensors
37            grad_input = grad_output * (1 - torch.tanh(input) ** 2)
38            #grad_input = grad_output * (torch.abs(input) <= 1).float()
39            #grad_input = grad_output * (1 - input / (1 + torch.abs(input)) ** 2)
40            return grad_input
41
42    class BinaryWeights(torch.autograd.Function):
43        @staticmethod
44        def forward(ctx, weights):
45            return torch.sign(weights)   # Forward -> binarize weights
46
47        @staticmethod
48        def backward(ctx, grad_output):
49            return grad_output   # Backward -> apply STE
```

*Figure 4 Activation Binarization*

In case of sign function, the back propagation is not applied thus we used a STE method which is:

$$\text{grad\_input} \approx 1 - \tanh^2 x$$

It is derived from the derivative properties of the hyperbolic tangent (tanh)

function. The tanh function is defined as:

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

Its range is $(-1, 1)$, with values approaching $1$ as $x \to \infty$ and $-1$ as $x \to -\infty$. The square of $\tanh(x)$ spans the range $[0, 1)$ and describes how closely the value of $\tanh(x)$ approaches its upper bound. The derivative of $\tanh(x)$ is:

$$\frac{d}{dx}\tanh(x) = 1 - \tanh^2(x)$$

Here, the equation $\text{grad\_input} \approx 1 - \tanh^2(x)$ works because it is derivative of approximated sign function.

However, there are a few drawbacks to this method.

1. Precision Near Zero
   Around $x = 0$, the approximation may deviate significantly from the true sign function due to the nature of $\tanh(x)$, which does not change abruptly at $x = 0$.

2. Approximation Error
   For extreme values of $x$, $1 - \tanh^2(x)$ does not converge exactly to $\pm 1$, leading to slight inaccuracies compared to the ideal sign function.

Thus, the proposed method for sign function results in a reduction of final accuracy. This could be resolved by avoiding the values near zero to be determined incorrectly, with other methods.

```
68      def _forward_features(self, x):
69          x = binary_quantization(x)
70
71          binary_conv1_weight = BinaryWeights.apply(self.conv1.weight)
72          x = F.conv2d(x, binary_conv1_weight, stride=1)
73          #x = F.relu(x)
74          x = BinaryActivation.apply(x)
75
76          binary_conv2_weight = BinaryWeights.apply(self.conv2.weight)
77          x = F.conv2d(x, binary_conv2_weight, stride=1)
78          #x = F.relu(x)
79          x = BinaryActivation.apply(x)
80
81          x = F.avg_pool2d(x, 2)
82          x = BinaryActivation.apply(x)
83          return x
```

*Figure 5 In* `BinaryNN` *function,* `BinaryActivation` *functions were applied between each layer.*

The training step was done with binarization functions between each layer.

```
 95    # Weight Initialization
 96    def initialize_weights(model):
 97        for m in model.modules():
 98            if isinstance(m, nn.Conv2d) or isinstance(m, nn.Linear):
 99                nn.init.xavier_uniform_(m.weight)
100
101    # Dataset load / Transform definition
102    transform = transforms.Compose([
103        transforms.ToTensor(),
104        transforms.Normalize((0.0,), (1.0,))
105    ])
```

For better accuracy, image normalization and appropriate weight initialization function were used.

```
116    train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
117    test_loader = DataLoader(test_dataset, batch_size=32, shuffle=False)
118
119    # model, loss function, and optimizer
120    model = BinaryNN()
121    initialize_weights(model)
122    criterion = nn.CrossEntropyLoss()
123    #optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.9)
124    optimizer = optim.Adam(model.parameters(), lr=0.001)
```

For BNN, smaller batch size leads to better training thus we've set the batch size to 32. For gradient descent method, we used Adam.

```
Epoch [1/10], Loss: 7.8605, Accuracy: 80.74%
Epoch [2/10], Loss: 5.9129, Accuracy: 85.65%
Epoch [3/10], Loss: 5.7389, Accuracy: 86.14%
Epoch [4/10], Loss: 5.2453, Accuracy: 87.34%
Epoch [5/10], Loss: 4.7340, Accuracy: 88.55%
Epoch [6/10], Loss: 4.0968, Accuracy: 89.48%
Epoch [7/10], Loss: 4.0990, Accuracy: 89.52%
Epoch [8/10], Loss: 3.8908, Accuracy: 89.69%
Epoch [9/10], Loss: 3.5470, Accuracy: 90.52%
Epoch [10/10], Loss: 3.0475, Accuracy: 91.33%
```

*Figure  6 Training result of BNN*

The training results are shown above.

## 4. Quantization Performance (Inference Accuracy)

The inference accuracy of quantized BNN was 93.62%

```
 88%|
atch 620: Accuracy = 93.72%
 90%|
atch 640: Accuracy = 93.75%
 94%|
atch 660: Accuracy = 93.77%
 96%|
atch 680: Accuracy = 93.69%
100%|
Final Accuracy: 93.62%
```
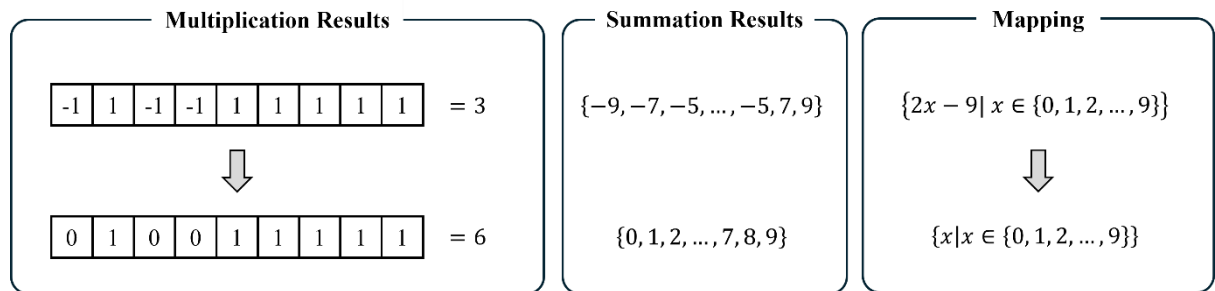
while FP32 showed 98.26% accuracy.

```
 86%|                                              | 600/700 [05:23<00:52,  1
 89%|                                              | 620/700 [05:34<00:41,  1
 91%|                                              | 640/700 [05:44<00:32,  1
 94%|                                              | 660/700 [05:55<00:21,  1
 97%|                                              | 680/700 [06:06<00:12,  1
100%|                                              | 700/700 [06:17<00:00,  1
Final Accuracy: 98.26%
```

## II. Hardware Architecture

### 1. XNOR and Popcount

The weights and activation, which was set as -1 and 1 in software, is assigned to 0 and 1 in hardware. In hardware, we can implement multiplication of weights and activations in software with XNOR operation, and addition with popcount.

| A | B | Result |
|---|---|--------|
| -1 | -1 | 1 |
| -1 | 1 | -1 |
| 1 | -1 | -1 |
| 1 | 1 | 1 |

| A | B | Result | XNOR |
|---|---|--------|------|
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 |

*Figure 7 Results of XNOR is the equivalent with Multiplication of BNN.*

**Multiplication Results**

| -1 | 1 | -1 | -1 | 1 | 1 | 1 | 1 | 1 | = 3

| 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | = 6

**Summation Results**

$\{-9, -7, -5, \dots, -5, 7, 9\}$

$\{0, 1, 2, \dots, 7, 8, 9\}$

**Mapping**

$\{2x - 9 | x \in \{0, 1, 2, \dots, 9\}\}$

$\{x | x \in \{0, 1, 2, \dots, 9\}\}$

**Figure 8 Popcount can replace Summation Results in convolution when BNN.**

In weight binarization step, we save trained weights in `.npy` file in two different ways: -1/1 for software computation and 0/1 for hardware computation.

```
∨ PYTHON
  ≡ binary_model_uint8_-11.npy
  ≡ binary_model_uint8_01.npy
```

The one with weight values of -1 and 1 was used on inference step in CPU (for

verifying pretrained model), while the one with weight values of 0 and 1 was used to hardware design.

```python
 9    ############################# weight binarize ####################
10
11    # Load model
12    torchdict = torch.load('model.pth')
13
14    # binarized weight dictionary
15    numpydict = {}
16
17    # weight binarization and convert to uint8
18    def binarize_and_convert_to_uint8(tensor):
19        binary_tensor = torch.sign(tensor)  # -1 & 1
20        binary_tensor = (binary_tensor + 1) // 2  # mapping to 0 & 1
21        return binary_tensor.to(dtype=torch.int8)  # uint8
22
23    # binarized conv1 weight
24    conv1_weight = torch.tensor(torchdict['conv1.weight'])
25    binary_conv1_weight = binarize_and_convert_to_uint8(conv1_weight)
26    numpydict['conv1w'] = binary_conv1_weight.numpy()
27
28    # binarized conv2 weight
29    conv2_weight = torch.tensor(torchdict['conv2.weight'])
30    binary_conv2_weight = binarize_and_convert_to_uint8(conv2_weight)
31    numpydict['conv2w'] = binary_conv2_weight.numpy()
32
33    # binarized fc weight
34    fc_weight = torch.tensor(torchdict['fc.weight'])
35    binary_fc_weight = binarize_and_convert_to_uint8(fc_weight)
36    numpydict['fc3w'] = binary_fc_weight.numpy()
37
38    # binarized weights saved as .npy
39    np.save('binary_model_uint8_01.npy', numpydict, allow_pickle=True)
```

*Figure  9 Creation of numpy file based on model file.*

```python
47    # Feed-forward function
48 ∨ def feed_forward(X0):
49        X0 = binary_quantization(X0)
50        X0 = torch.tensor(X0.reshape(-1, 1, 28, 28), dtype=torch.float32)  # Shape: [batch_size, 1, 28, 28]
51
52        # Conv1
53        X1 = F.conv2d(X0, conv1w)  # Shape: [batch_size, 16, 26, 26]
54        X1 = binary_quantization(X1, threshold=0.0)
55
56        # Conv2
57        X2 = F.conv2d(X1, conv2w)  # Shape: [batch_size, 16, 24, 24]
58        X2 = binary_quantization(X2, threshold=0.0)
59
60        # Avg Pooling
61        A2 = avg_pool2d(X2, kernel_size=2, stride=2)  # Shape: [batch_size, 16, 12, 12]
62        A2 = binary_quantization(A2, threshold=0.0)
63        A2 = A2.view(A2.size(0), -1)  # Flatten to [batch_size, 2304]
64
65        # Fully connected layer
66        X3 = A2 @ fc3w.T  # Shape: [batch_size, 10]
67        return X3
68
```

With the weights saved as numpy files, the verification of model is done by evaluating accuracy when going through the same algorithm with hardware as forward propagation.

```
 79%|
atch 560: Accuracy = 93.97%
 82%|
atch 580: Accuracy = 93.78%
 85%|
atch 600: Accuracy = 93.62%
 88%|
atch 620: Accuracy = 93.72%
 91%|
atch 640: Accuracy = 93.75%
 94%|
atch 660: Accuracy = 93.77%
 97%|
atch 680: Accuracy = 93.69%
100%|
Final Accuracy: 93.62%
```

Inferencing gives us 93.62% accuracy as well as training step.

```cpp
#pragma HLS UNROLL factor=16
uint16_t xnor_result_each_channel =
    ~(conv2_kernels[c][nk] ^ input_layer2[t][c][row][col]);
popcount_result_all_channel += popcount9(xnor_result_each_channel, popcount_lut_9);
```

```cpp
int16_t sum = 0;
for (int i = 0; i < 288; i++) {
    uint8_t xnor_result = ~(packed_output[t][i] ^ fc_weights[out_idx][i]);
    sum += popcount8(xnor_result, popcount_lut_8);
}
```

```cpp
34    const int8_t popcount_lut_8[256] = {-8, -6, -6, -4, -6, -4, -4, -2, -6, -4, -4, -2, -4, -2,
35    ;
36
37
38    const int8_t popcount_lut_9[512] = {-16, -14, -14, -12, -14, -12, -12, -10, -14, -12, -12,
39    ;
40
```

The popcount is done with LUT which reduces latency.

In our hardware implementation, we utilized `uint8_t` and `uint16_t` data types, which adhere to the C++ standard. Consequently, a custom packing method was developed to support XNOR and popcount operations. Although the `ap_uint<>` format offered a more efficient alternative, time constraints prevented us from revising the architecture to incorporate it.

**2. Bit Packing**



*Figure  1  0 Data Dimension transformation after Bit Packing*

The 3x3 kernel-based convolution operation was implemented using the IM2COL approach, as illustrated in the figure. This method transforms the input data into a column-major matrix, enabling efficient computation. In this case, for a 28x28 MNIST input, the sliding window approach is replaced with IM2COL to optimize for parallel computation, significantly improving hardware performance for operations involving small kernels.

By flattening the overlapping regions into a structured matrix, the convolution process becomes better suited for matrix multiplication, facilitating parallelism and hardware acceleration.



*Figure   1   1 Filling Remaining Slots of uint16 without effecting popcount result*

As a result, each 3x3 kernel region contains a total of 9 bits of data. For the image data, the upper 7 bits are padded with zeros, while for the weights, the upper 7 bits are padded with ones. This ensures that the XNOR operation produces a base value of -7. Consequently, the binary activation function in hardware is implemented with -7 as the threshold. This approach is applied specifically to `conv1w` and `conv2w`, whereas the `fc_weight` fully utilizes all 8 bits for packing.

The weights, packed according to the defined rules, are hardcoded into a header file to maintain weight stationarity during computation. This ensures efficient access and reuse of weights throughout the operation.

## III. Parallelization

### 1. Input Stream

The MNIST images are received directly using an input stream without any preprocessing. An input stream refers to a sequential flow of data, where each element (e.g., pixel values) is transmitted in a pipeline-like manner. This method is particularly efficient for hardware accelerators as it minimizes memory latency and allows real-time data ingestion. The input stream delivers grayscale MNIST images of size 28x28 as 8-bit unsigned integers (`uint8_t`).

The `load_input` function loads the data from the input stream into a 3D array `input_image` for further processing. The data from the stream is accessed using the `pop_stream` function, which extracts pixel values sequentially. The indices are computed based on the batch index (t), row, and column to correctly map the streamed data into the 3D array.

```
11   void load_input(AXI_VAL_uint8* input_stream,
12                    uint8_t input_image[BATCH_SIZE][28][28]) {
13
14       for (int t = 0; t < BATCH_SIZE; t++) {
15           // t is each image
16           for (int row = 0; row < 28; row++) {
17               for (int col = 0; col < 28; col++) {
18                   int idx = t * 28 * 28 + row * 28 + col;
19                   input_image[t][row][col] = pop_stream<uint8_t>(input_stream[idx]);
20               }
21           }
22       }
23   }
24
```

*Figure  1  2 Load input with Input Stream*

- `AXI_VAL_uint8`: Represents the input stream data type for AXI interfaces.
- `pop_stream`: Extracts values from the stream one by one.

The indices ensure that the pixels are placed in the correct position in the `input_image` array for each batch, row, and column.


### 2. IM2COL

Once the MNIST images are loaded, they are preprocessed using the IM2COL technique to prepare for convolution. IM2COL stands for "Image-

to-Column," which is a method that rearranges image data into a structured matrix to align with matrix multiplication. It converts overlapping regions of an image, based on the kernel size (3x3 in this case), into contiguous columns. This transformation enables parallel computation of convolution operations, particularly when used with matrix multiplication hardware accelerators.

In the `pre_processing` function, the 3x3 regions of each image are packed into a 9-bit representation and stored as a 16-bit value (`uint16_t`). Each bit in the packed value corresponds to a single pixel in the 3x3 window, shifted into the correct position:

- Bit 8 to Bit 0:
  Represent the 9 pixels of the 3x3 region, extracted using bitwise AND (&) and shifted using <<.

- This packed format reduces memory usage and accelerates subsequent XNOR operations in binary convolution layers.

The packed values are stored in a 3D array `input_layer1`, where each 3x3 window is represented as a single 16-bit value.

```
25   void pre_processing(uint8_t input_image[BATCH_SIZE][28][28],
26                       uint16_t input_layer1[BATCH_SIZE][26][26]) {
27       for (int t = 0; t < BATCH_SIZE; t++) {
28           for (int row = 0; row < 26; row++) {
29               for (int col = 0; col < 26; col++) {
30                   uint16_t packedValue = 0;
31
32                   packedValue |= (input_image[t][row][col] & 1) << 8;
33                   packedValue |= (input_image[t][row][col + 1] & 1) << 7;
34                   packedValue |= (input_image[t][row][col + 2] & 1) << 6;
35                   packedValue |= (input_image[t][row + 1][col] & 1) << 5;
36                   packedValue |= (input_image[t][row + 1][col + 1] & 1) << 4;
37                   packedValue |= (input_image[t][row + 1][col + 2] & 1) << 3;
38                   packedValue |= (input_image[t][row + 2][col] & 1) << 2;
39                   packedValue |= (input_image[t][row + 2][col + 1] & 1) << 1;
40                   packedValue |= (input_image[t][row + 2][col + 2] & 1);
41
42                   input_layer1[t][row][col] = packedValue;
43               }
44           }
45       }
46   }
47
```

*Figure 1 3 Preprocessing of Input Image (IM2COL)*

i.  `input_image`:
    3D array holding the original MNIST images of size 28x28.

ii. `input_layer1`:
    3D array storing the packed 3x3 regions of size 26x26 (as the 3x3 kernel sliding window reduces the effective size).

iii. Bit Packing:
The least significant bit (LSB) of each pixel is extracted and shifted to its corresponding position in the 9-bit structure.

iv. Parallelization:
The code is "unrolled" to explicitly handle each pixel in the 3x3 window, enhancing readability and ensuring efficient parallel implementation.

## 3. Convolution Layer 1

In the first convolutional layer, computations are performed using XNOR and popcount operations, which are highly efficient for binary data processing. These operations accelerate the convolution process compared to conventional arithmetic operations. The `#pragma HLS PIPELINE` directive is applied to the innermost loop with an Initiation Interval (II) of 1, ensuring that operations are pipelined effectively.

Pipelining allows multiple operations to overlap, improving hardware throughput by processing new iterations before previous ones are fully completed. This is particularly important in the convolution layer, where the same computation (XNOR and popcount) is repeated for a large number of inputs and kernels. The loop is "flattened" and explicitly written to enhance hardware synthesis and maximize parallel execution.

In the following function, `popcount9` computes the number of matching bits from the XNOR result, and the activation function is applied to binarize the output values.

```
48    void layer1_conv2d(uint16_t input_layer1[BATCH_SIZE][26][26],
49                        uint8_t output_layer1[BATCH_SIZE][16][26][26]) {
50
51        for (int t = 0; t < BATCH_SIZE; t++) {
52            for (int nk = 0; nk < 16; nk++) {
53                for (int row = 0; row < 26; row++) {
54                    for (int col = 0; col < 26; col++) {
55    #pragma HLS PIPELINE II=1
56                        int16_t popcount_result_all_channel = 0;
57                        uint16_t xnor_result_each_channel =
58                            ~(conv1_kernels[nk] ^ input_layer1[t][row][col]);
59
60                        popcount_result_all_channel = popcount9(xnor_result_each_channel, popcount_lut_9);
61
62                        output_layer1[t][nk][row][col] = (popcount_result_all_channel >= -7) ? 1 : 0;
63                        // without bit packing
64                        // activation function line (instead of relu)
65                    }
66                }
67            }
68        }
69    }
70
```

*Figure  1 4 conv1 Layer computation*

After passing through the first convolutional layer, the output values are packed into a `uint16_t` format using the IM2COL approach. This transformation reorganizes the binary output data into a compact structure,

aligning it for the next convolutional layer.

IM2COL simplifies convolution into matrix multiplications, enabling efficient reuse of the same data and improving hardware throughput. The loop is unrolled and explicitly written to allow parallel computation for each 3x3 sliding window. Similar to the previous function, `#pragma HLS PIPELINE` is applied to ensure that the bit-packing operations are pipelined for maximum efficiency.

```
71  void layer1_packing(uint8_t output_layer1[BATCH_SIZE][16][26][26],
72                      uint16_t input_layer2[BATCH_SIZE][16][24][24]) {
73
74      for (int t = 0; t < BATCH_SIZE; t++) {
75          for (int nk = 0; nk < 16; nk++) {
76              for (int row = 0; row < 24; row++) {
77                  for (int col = 0; col < 24; col++) {
78  #pragma HLS PIPELINE II=1
79                      uint16_t packedValue = 0;
80
81                      packedValue |= (output_layer1[t][nk][row][col] & 1) << 8;
82                      packedValue |= (output_layer1[t][nk][row][col + 1] & 1) << 7;
83                      packedValue |= (output_layer1[t][nk][row][col + 2] & 1) << 6;
84                      packedValue |= (output_layer1[t][nk][row + 1][col] & 1) << 5;
85                      packedValue |= (output_layer1[t][nk][row + 1][col + 1] & 1) << 4;
86                      packedValue |= (output_layer1[t][nk][row + 1][col + 2] & 1) << 3;
87                      packedValue |= (output_layer1[t][nk][row + 2][col] & 1) << 2;
88                      packedValue |= (output_layer1[t][nk][row + 2][col + 1] & 1) << 1;
89                      packedValue |= (output_layer1[t][nk][row + 2][col + 2] & 1);
90
91                      input_layer2[t][nk][row][col] = packedValue;
92                  }
93              }
94          }
95      }
96  }
97
```

*Figure 1 5 Output of conv1 Layer is packed into* `uint16_t`.

## 4. Convolution Layer 2

The second convolutional layer is the most computationally expensive part of the network, as it processes all channels with a high number of operations. To address this bottleneck, significant optimization is applied using loop unrolling and pipelining techniques to maximize parallel execution and minimize latency.

- Loop Unrolling:
  The inner loop that iterates over the channels (`c = 0 to 15`) is unrolled with a factor of 16 using the `#pragma HLS UNROLL factor=16` directive. This aggressive level of unrolling allows simultaneous computation across all 16 channels, significantly reducing the number of clock cycles required. Compared to the implementation without unrolling, this optimization resulted in a 5x speedup.

- Pipelining:
  The second part of the layer performs average pooling over a 2x2

region. To further optimize latency, `#pragma HLS PIPELINE II=1` is applied to ensure that new iterations of the pooling loop begin every clock cycle, maintaining a steady flow of data.

After average pooling, the output values are flattened into a 1D array to prepare for the fully connected layer. This flattening process ensures the data is stored sequentially, ready for further processing.

## 5. Packing Flatten Output

After the average pooling and flattening stages, the output data is packed into a more compact format using 8-bit containers (`uint8_t`). This process involves filling each byte starting from the most significant bit (MSB) with consecutive binary values from the flattened data. The goal is to reduce the memory footprint and prepare the output for efficient transfer or further computation.

The function `pack_flatten_output` iterates through the flattened data and shifts the binary values into their respective bit positions within an 8-bit `current_byte`. Once all 8 bits of the `current_byte` are filled, the byte is stored in the output array `packed_output`.

Given that the packing process is relatively lightweight in terms of computation time, no further optimizations, such as pipelining or unrolling, were applied.

```
142    void pack_flatten_output(uint8_t flatten[BATCH_SIZE][16 * 12 * 12],
143                             uint8_t packed_output[BATCH_SIZE][2 * 12 * 12]) {
144        for (int t = 0; t < BATCH_SIZE; t++) {
145            int packed_idx = 0;
146            uint8_t current_byte = 0;
147
148            for (int i = 0; i < 16 * 12 * 12; i++) {
149                current_byte |= (flatten[t][i] & 1) << (7 - (i % 8));
150
151                if ((i + 1) % 8 == 0) {
152                    packed_output[t][packed_idx++] = current_byte;
153                    current_byte = 0;
154                }
155            }
156        }
157    }
158
```

## 6. Fully Connected Layer

The fully connected (FC) layer is implemented to perform the final classification stage of the network. This layer computes the weighted sum of the binary-packed input data using XNOR and popcount operations, which are efficient for binary neural networks. Although this design considers optimization, additional resource constraints, such as memory usage, were prioritized, leading to a straightforward implementation.

Compared to earlier stages (convolution and pooling), the fully connected layer requires relatively less computation time and is not a critical performance bottleneck. Therefore, further optimizations like pipelining or unrolling were not applied.

```
159  void layer3_fc(uint8_t packed_output[BATCH_SIZE][2 * 12 * 12],
160                 int16_t fc_output[BATCH_SIZE][10]){
161      for (int t = 0; t < BATCH_SIZE; t++) {
162          for (int out_idx = 0; out_idx < 10; out_idx++) {
163              int16_t sum = 0;
164              for (int i = 0; i < 288; i++) {
165                  uint8_t xnor_result = ~(packed_output[t][i] ^ fc_weights[out_idx][i]);
166                  sum += popcount8(xnor_result, popcount_lut_8);
167              }
168              fc_output[t][out_idx] = sum;
169          }
170      }
171  }
172
```

## 7. Output Stream

The final output of the fully connected layer is sent out using an output stream interface, which efficiently transmits data in a sequential manner. This approach aligns with AXI-stream protocols often used in hardware systems, facilitating real-time and resource-efficient data communication.

The `store_output` function iterates through the output data (`fc_output`), which contains the classification results for each batch and outputs the values into an AXI-compatible stream (`output_stream`).

The `push_stream` function wraps the data and appends an end-of-stream (EOS) flag for the last value in the stream, signaling the completion of data transfer.

```
173  void store_output(AXI_VAL_int16* output_stream,
174                    int16_t fc_output[BATCH_SIZE][10]) {
175
176      for (int t = 0; t < BATCH_SIZE; t++) {
177          for (int k = 0; k < 10; k++) {
178              int idx = t * 10 + k;
179              bool is_last = (t == BATCH_SIZE - 1) && (k == 9);
180              output_stream[idx] = push_stream<int16_t>(fc_output[t][k], is_last);
181          }
182      }
183  }
184
```
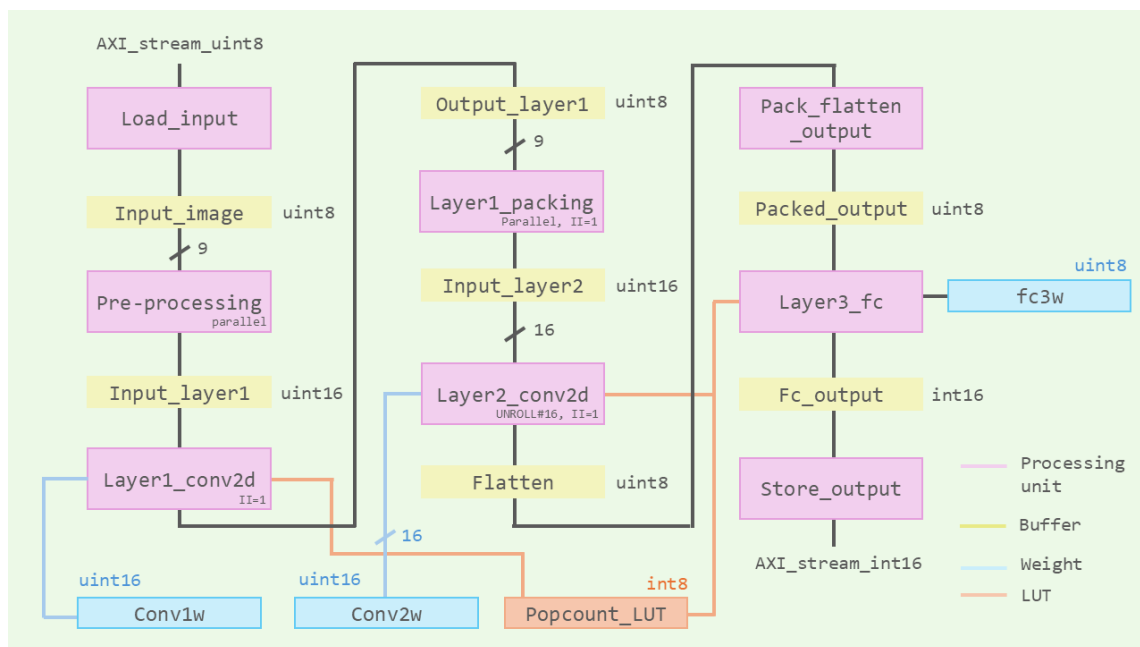
## 8. Overall

Overall, the structure of the top module is shown below, where array partitioning of buffers is applied to facilitate parallel access during computation.

```
211        // Process each tile
212
213        // 1. Load Input Tile
214        load_input(input_stream, input_image);
215
216        // 2. Pre-processing input_image
217        pre_processing(input_image, input_layer1);
218
219        // 3. Layer 1: Convolution
220        layer1_conv2d(input_layer1, output_layer1);
221
222        // 4. Layer 1: Packing
223        layer1_packing(output_layer1, input_layer2);
224
225        // 5. Layer 2: Convolution + Avg Pooling
226        layer2_conv2d(input_layer2, flatten);
227
228        // 6. Pack Flattened Output
229        pack_flatten_output(flatten, packed_output);
230
231        // 7. Fully Connected Layer
232        layer3_fc(packed_output, fc_output);
233
234        // 8. Store Output Tile
235        store_output(output_stream, fc_output);
```

## 9. Hardware Dataflow and Structure



## IV. Hardware Synthesis

## 1. C-simulation

The testbench is set as below. (See the number patterns on the right.)

```
C+ CNN_tb_pre.cpp > ...
 1   #include <iostream>
 2   #include <vector>
 3   #include <algorithm>
 4   #include "dma_template.h"
 5
 6   #define INPUT_TILE_SIZE 784 // 28 * 28
 7   #define OUTPUT_TILE_SIZE 10
 8
 9   //typedef ap_axiu<16, 4, 5, 5> AXI_VAL_int16;
10
11   // Function declaration for the top function
12   void top_function(AXI_VAL_uint8* input_stream, AXI_VAL_int16* output_stream);
13
14   uint8_t array[10][28][28] = {
15     {
16       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
17       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
18       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
19       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
20       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
21       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0},
22       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0},
23       {0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
24       {0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
25       {0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
26       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
27       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
28       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
29       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
30       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
31       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
32       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
33       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
34       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0},
35       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
36       {0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
37       {0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
38       {0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
39       {0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
```



```
Vivado HLS 2020.1 - cnn_ip (C:\Users\krist\cnn\cnn_ip)
File   Edit   Project   Solution   Window   Help

Explorer 🔲        Synthesis(solution1)(top_function_csynth.rpt)   cnn_hw_top_csim.log 🔲
v 🗁 cnn_ip          108 Max Index per 10 element batch:
  > 🔗 Includes       109 Batch 0: Max Index = 5, Value = 202
  v 🔗 Source         110 Batch 1: Max Index = 0, Value = 230
      📄 CNN_hw.cpp   111 Batch 2: Max Index = 4, Value = 116
      📄 dma_template.h  112 Batch 3: Max Index = 1, Value = 140
  v 🔗 Test Bench     113 Batch 4: Max Index = 9, Value = 182
      📄 CNN_tb.cpp   114 Batch 5: Max Index = 2, Value = 208
  v 🗁 solution1      115 Batch 6: Max Index = 1, Value = 174
    v ✳ constraints  116 Batch 7: Max Index = 3, Value = 222
        📄 directives.tcl  117 Batch 8: Max Index = 1, Value = 154
        📄 script.tcl  118 Batch 9: Max Index = 4, Value = 176
    v 🗁 csim          119 INFO: [SIM 1] CSIM done with 0 errors.
      > 🗁 build       120 INFO: [SIM 3] ************* CSIM finish *************
                      121
```

The results show that Max Index matches the number shown in each input.

## 2. C synthesis

After synthesizing the ip, the Synthesis Report shows:

**Synthesis Report for 'top_function'**

**General Information**

| | |
|---|---|
| Date: | Tue Dec 17 20:03:14 2024 |
| Version: | 2020.1 (Build 2897737 on Wed May 27 20:21:37 MDT 2020) |
| Project: | cnn_ip |
| Solution: | solution1 |
| Product family: | zynq |
| Target device: | xc7z020-clg400-1 |

**Performance Estimates**

**Timing**

**Summary**

| Clock | Target | Estimated | Uncertainty |
|---|---|---|---|
| ap_clk | 10.00 ns | 9.634 ns | 1.25 ns |

**Latency**

**Summary**

| Latency (cycles) | | Latency (absolute) | | Interval (cycles) | | |
|---|---|---|---|---|---|---|
| min | max | min | max | min | max | Type |
| 1290965 | 1290965 | 12.910 ms | 12.910 ms | 1290965 | 1290965 | none |

**Detail**

- Instance
- Loop

**Utilization Estimates**

**Summary**

| Name | BRAM_18K | DSP48E | FF | LUT | URAM |
|---|---|---|---|---|---|
| DSP | - | - | - | - | - |
| Expression | - | - | 0 | 207 | - |
| FIFO | - | - | - | - | - |
| Instance | 20 | 1 | 1989 | 6916 | 0 |
| Memory | 211 | - | 0 | 0 | 0 |
| Multiplexer | - | - | - | 8338 | - |
| Register | - | - | 157 | - | - |
| Total | 231 | 1 | 2146 | 15461 | 0 |
| Available | 280 | 220 | 106400 | 53200 | 0 |
| Utilization (%) | 82 | ~0 | 2 | 29 | 0 |

**Detail**

- Instance
- DSP48E
- Memory
- FIFO
- Expression
- Multiplexer
- Register

**Interface**

Export the report(.html) using the Export Wizard

Open Analysis Perspective    Analysis Perspective

- Clock Frequency: 10.00ns (estimated: 9.634ns)
- Resource Utilization:

**Utilization Estimates**

**Summary**

| Name | BRAM_18K | DSP48E | FF | LUT | URAM |
|---|---|---|---|---|---|
| DSP | - | - | - | - | - |
| Expression | - | - | 0 | 207 | - |
| FIFO | - | - | - | - | - |
| Instance | 20 | 1 | 1989 | 6916 | 0 |
| Memory | 211 | - | 0 | 0 | 0 |
| Multiplexer | - | - | - | 8338 | - |
| Register | - | - | 157 | - | - |
| Total | 231 | 1 | 2146 | 15461 | 0 |
| Available | 280 | 220 | 106400 | 53200 | 0 |
| Utilization (%) | 82 | ~0 | 2 | 29 | 0 |

- # of Cycles: 1,292,965 (12.910 ms)
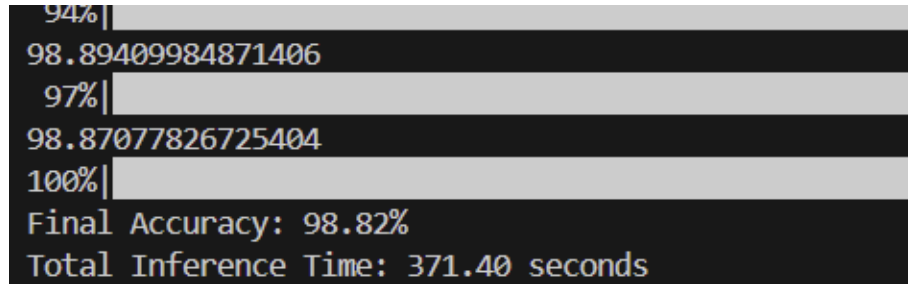
## 3. Vivado Top diagram





## 4. Execution on PYNQ



Running on PYNQ, the time took for execution of all 70,000 images was 97.88 seconds



*Figure 1 6 Execution time of TPU on PYNQ*

while FP32 CPU took 371.40 seconds.



*Figure  1  7 Accuracy and Runtime of FP32 CPU*



Runtime breakdown:

- Load Input (8,421): 0.65%
- Pre-processing (41,101): 3.18%
- Conv1 (108,961): 8.44%
- Conv1 Packing (460,804): 35.69%
- Conv2 (515,681): 39.95%
- Fully Connected (86,621): 6.71%
- Store Output (221): 0.02%

## V. Discussion

In our project, the hardware implementation achieved an accuracy of 93.89%, which is slightly lower than the 98.82% accuracy obtained on the FP32 CPU implementation. However, the execution time for our hardware was 97.88 seconds, significantly reducing the FP32 inference time of 371.40 seconds by approximately 73.65%. On the PYNQ board, the FP32 inference took even longer, requiring about 3.5 hours, which underscores the efficiency of our hardware implementation.

*Figure 1 8 Final Accuracy of BNN on PYNQ board*

Despite these improvements, there are still areas for further optimization. As mentioned earlier, the sign function was approximated using the STE-based tanh function, which can introduce errors for values close to zero due to the approximation's inherent nonlinearity. This limitation may slightly affect the final accuracy of the network.

Additionally, we used the C++ standard data types `uint8` and `uint16` for activation and weight storage. This approach was suboptimal for the 9-bit binary activations and weights, as the design used 16 bits for storage despite requiring only 9 bits. This resulted in inefficient resource utilization. By replacing the standard types with `ap_uint<>`, we expect to achieve better utilization of hardware resources, reducing unnecessary overhead and improving efficiency.

Overall, while our implementation demonstrates substantial improvements in inference time compared to FP32, future work addressing data type efficiency and approximation accuracy can further enhance the performance and resource utilization of the design.

**VI. Role & contributions of each team member**

20201082 H. Kim

- Design and implementation of hardware architecture, including HLS coding, functional debugging, and performance optimization
- Development of the initial BNN software pipeline, including training and inference framework setup
- Project management through schedule planning, role distribution, and overall task coordination

20211184 H. Yang

- BNN model optimization through parameter tuning (threshold adjustments) to achieve high-accuracy weight extraction
- Development of auxiliary scripts for HW implementation, including LUT value formatting and bit-packing automation
- Refining final report into the completed version, creating project presentation (PPT) and drafting the presentation script