

Embedded Hardware Systems Design Practical

Task description - Implement a Conv2D operator and then utilize it to implement a Convolutional layer. The implementation should be tested on an Ultra96 board using a test dataset. Furthermore, precision scaling would be considered to explore accuracy-performance trade-offs. Both HDL and HLS can be used for the implementation.

VIVADO HLS - For this task I used custom IP, which performs convolution operation and has parameters like Image and kernel matrix size, stride and padding using Vivado HLS software. I choose HLS as it is easy and convenient to compile, simulate, debug and optimise algorithm implement in C/C++. The parameters are fixed but can be changed in the header file according to the requirements and again export the IP.

```
header.h
1 #include <stdio.h>
2 #include <hls_stream.h>
3 #include <ap_int.h> // For last signal (1-bit) of streaming protocol
4
5
6 #define image_size 3 // Size of Input image
7 #define kernel_size 3 // Kernel matrix size
8 #define stride 1 // Step size of kernel movement (Assumed equal in both dimension)
9 #define padding 1 // Additional pixel added to input image (Assumed equal in both dimension)
10
11
12 // Define the data type of matrix inputs and outputs
13 typedef double Mat_Dtype;
14
15
16 // Define the axis data structure
17 struct axis_data {
18     Mat_Dtype data;
19     ap_uint<1> last;
20 };
21
22
23 void Conv(hls::stream<axis_data> &data_in, hls::stream<axis_data> &data_out);
24
25
```

- The data is read and sent by IP using AXI input stream and AXI output stream channels create using below pragmas.
 - #pragma HLS INTERFACE ap_ctrl_none port=return
 - #pragma HLS INTERFACE axis register both port=data_in
 - #pragma HLS INTERFACE axis register both port=data_out
- Image and kernel matrix are created and fed in values using stream input. Then padding is applied to the input image matrix. Then the convolution operation is carried out as below and then output image is sent via stream output channel.

```
// Read data for Input Image Matrix
L1_Image_Matrix: for(i=0; i < image_size; i++){
    L2_Image_Matrix: for(j=0; j < image_size; j++){
        local_stream = data_in.read();
        input_mat[i][j] = local_stream.data;}}

// Read data for Kernel Matrix
L1_Kernal_Matrix: for(i=0; i < kernel_size; i++){
    L2_Kernal_Matrix: for(j=0; j < kernel_size; j++){
        local_stream = data_in.read();
        kernal_mat[i][j] = local_stream.data;}}

// Convolution operation
C1: for (i = 0; i < output_size; ++i) {
    C2: for (j = 0; j < output_size; ++j) {
        output_image[i][j] = 0;
        C3: for (m = 0; m < kernal_size; ++m) {
            C4: for (n = 0; n < kernal_size; ++n) {
                output_image[i][j] += padded_image[i * stride + m][j * stride + n] * kernal_mat[m][n]; } } } }

// Output stream the result
output_loop_1: for(i=0; i < output_size; i++){
    output_loop_2: for(j=0; j < output_size; j++){
        local_stream.data = output_image[i][j];

        // Last signal and the strobe signal
        if((i==output_size-1) && (j==output_size-1))
            local_stream.last = 1;
        else
            local_stream.last = 0;
        data_out.write(local_stream);}}
```

Un-optimized code

- 2 different types of IP were implemented and their performance were compared. One was without any optimisations using directives and other was fully optimised. For optimisation, I used following pragmas.
- **#pragma HLS LOOP_FLATTEN** - Allows nested loops to be flattened into a single loop hierarchy with improved latency. This saves clock cycles, potentially allowing for greater optimization of the loop body logic.
- **#pragma HLS LOOP_MERGE** - It merges consecutive loops into a single loop to reduce overall latency, increase sharing, and improve logic optimization. Reduces the number of clock cycles required in the RTL to transition between the loop-body implementations.
- **#pragma HLS PIPELINE II=5** – A pipelined function or loop can process new inputs every 5 clock cycles, where 5 is the II (initiation interval) of the loop or function. It allows the concurrent execution of operations.

```
// Read data for Input Image Matrix
L1_Image_Matrix: for(i=0; i < image_size; i++){
  L2_Image_Matrix: for(j=0; j < image_size; j++){
    #pragma HLS LOOP_FLATTEN
    local_stream = data_in.read();
    input_mat[i][j] = local_stream.data;}}

// Read data for Kernal Matrix
L1_Kernal_Matrix: for(i=0; i < kernal_size; i++){
  #pragma HLS LOOP_MERGE
  L2_Kernal_Matrix: for(j=0; j < kernal_size; j++){
    local_stream = data_in.read();
    kernal_mat[i][j] = local_stream.data;}}

// Convolution operation
C1: for (i = 0; i < output_size; ++i) {
  C2: for (j = 0; j < output_size; ++j) {
    #pragma HLS PIPELINE II=5
    output_image[i][j] = 0;
    C3: for (m = 0; m < kernal_size; ++m) {
      C4: for (n = 0; n < kernal_size; ++n) {
        output_image[i][j] += padded_image[i * stride + m][j * stride + n] * kernal_mat[m][n]; } } } }

// Output stream the result
output_loop_1: for(i=0; i < output_size; i++){
  #pragma HLS LOOP_MERGE
  output_loop_2: for(j=0; j < output_size; j++){
    local_stream.data = output_image[i][j];

    // Last signal and the strobe signal
    if((i==output_size-1) && (j==output_size-1))
      local_stream.last = 1;
    else
      local_stream.last = 0;
    data_out.write(local_stream);}}
```

Optimized code

- Analysing both design after synthesis

Table -1 Comparison between both design

Not – Optimised design	Optimised design																												
<ul style="list-style-type: none">• Latency (clock cycles)<ul style="list-style-type: none">◦ Summary<table><tr><th colspan="2">Latency</th><th colspan="2">Interval</th><th rowspan="2">Type</th></tr><tr><th>min</th><th>max</th><th>min</th><th>max</th></tr><tr><td>1172</td><td>1172</td><td>1172</td><td>1172</td><td>none</td></tr></table>	Latency		Interval		Type	min	max	min	max	1172	1172	1172	1172	none	<ul style="list-style-type: none">• Latency (clock cycles)<ul style="list-style-type: none">◦ Summary<table><tr><th colspan="2">Latency</th><th colspan="2">Interval</th><th rowspan="2">Type</th></tr><tr><th>min</th><th>max</th><th>min</th><th>max</th></tr><tr><td>217</td><td>217</td><td>217</td><td>217</td><td>none</td></tr></table>	Latency		Interval		Type	min	max	min	max	217	217	217	217	none
Latency		Interval		Type																									
min	max	min	max																										
1172	1172	1172	1172	none																									
Latency		Interval		Type																									
min	max	min	max																										
217	217	217	217	none																									

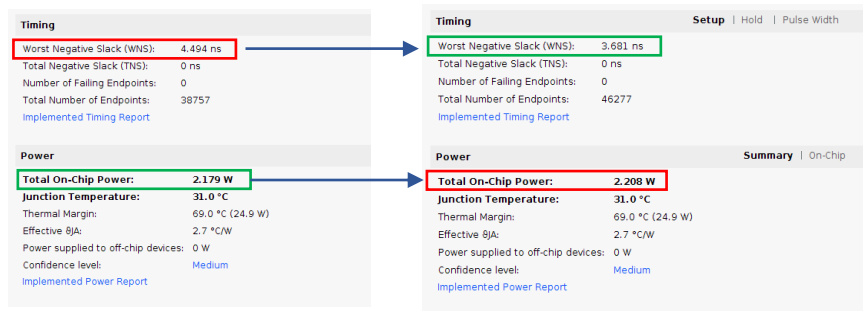


Figure-2 Timing and power report un-optimized and optimized design

- Uploading the .bit and .hwh files in Jupiter notebook and running the python script.
- The .py scripts first send the image and kernel data via send channel and receive output in output buffer. Then convolution function implemented in python using numpy library calculates the output.
- Both results are compared for equality. The end result are as follows.

```
print("Input Matrix:")
print(A_matrix)
print("\nKernel Matrix:")
print(B_matrix)
print("\nConvolution result by python script:")
print(result)

Input Matrix:
[[1. 2. 3.]
 [4. 5. 6.]
 [7. 8. 9.]]

Kernel Matrix:
[[-1. -2. -1.]
 [ 0.  0.  0.]
 [ 1.  2.  1.]]

Convolution result by python script:
[[ 13.  20.  17.]
 [ 18.  24.  18.]
 [-13. -20. -17.]]

In [9]: # Checking if both output matrix are equal
if np.array_equal(Output_Matrix, result):
    print("The matrices are equal.")
else:
    print("Both matrix doesn't match")

The matrices are equal.

In [10]: print('FPGA run time: ', fpga_run_time)
print('ARM PS run time: ', ps_run_time)

FPGA run time: 29.145351886749268
ARM PS run time: 0.0034723281808351562
```

Result of Unoptimized bit file

```
result = convolution(A_matrix, B_matrix, stride, padding)

end = time.time()
ps_run_time = end - start

print("Input Matrix:")
print(A_matrix)
print("\nKernel Matrix:")
print(B_matrix)
print("\nConvolution result by python script:")
print(result)

Input Matrix:
[[1. 2. 3.]
 [4. 5. 6.]
 [7. 8. 9.]]

Kernel Matrix:
[[-1. -2. -1.]
 [ 0.  0.  0.]
 [ 1.  2.  1.]]

Convolution result by python script:
[[ 13.  20.  17.]
 [ 18.  24.  18.]
 [-13. -20. -17.]]

In [20]: # Checking if both output matrix are equal
if np.array_equal(Output_Matrix, result):
    print("The matrices are equal.")
else:
    print("Both matrix doesn't match")

The matrices are equal.

In [21]: print('FPGA run time: ', fpga_run_time)
print('ARM PS run time: ', ps_run_time)

FPGA run time: 35.36079216063418
ARM PS run time: 0.003140687942504883
```

Result of Optimized bit file

Observation and Conclusion

- Pragma's were added to reduce the critical path and latency. It significantly **reduced the latency** by 81% as seen in table-1.
- We can also observe from table-1 that **no. of loop iteration also significantly decreased** by using Loop flatten and merge pragma.
- But then **resource utilization increased** in the optimised case that means overall design area increased too.
- Total on-chip **power increased slightly** as well. Which was expected as no. of resource increased, they would utilize more power. Refer figure-2.
- So, there is always trade-off between timing, power and area. All cannot be optimised at the same time. In this project timing was prioritised over others.

Submitted by,

Kishan Dadhania

Matriculation no. - 5112720