

## 1.0 Introduction

The following document serves as the lab report for assignment 2 of ECE 1756 (Fall 2021). This report is written by Sheng Zhao (1003273913).

Section 2.0 describes the FPGA implementation of the convolution function, its operation and performance while Section 3.0 describes the same metrics of the CPU and GPU implementation. A comparison of the different implementations will be done in Section 4.0.

## 2.0 FPGA Implementation

To enable a comparison between FPGA, CPU and GPU, a convolution engine was implemented on an Intel Arria 10 GX device (10AX115N2F45I1SG). The following section describes the details of this implementation as well as performance metrics.

### 2.1 Convolution and Input Data Details

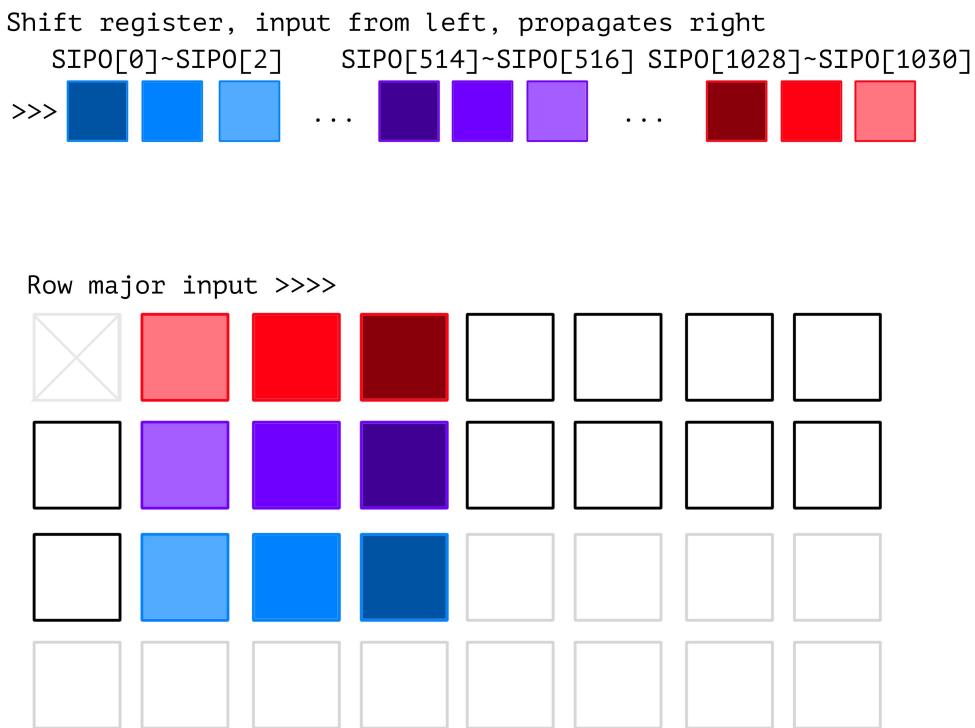


Figure 1: Data input and shift register

The goal of the convolution engine implementation on the FPGA is to compute the convolutions on a 512 pixel wide, any height image, along with a provided 3x3 filter. As shown in the bottom half of Fig. 1, image data is given in a row major format, sending in pixel by pixel starting with the top left. Convolution is then to be performed on every 3x3 set of pixels.

Given the row major input format, in order to compute convolutions on pixels that span across multiple rows, pixel data must be held for a period of time. To achieve this, a Serial-in, Parallel-out (SIPO) shift register is implemented. Given that the length of the image is fixed at 514 pixels (512 wide image with 0 padding), the minimum required length of the SIPO is 1031. Input pixels given in row major format will enter from the left (top half of Fig. 1), with the data propagating to the next register every clock cycle.

To better illustrate the pixels involved the convolution, consider an example as shown in Fig. 1. All the pixels required for the computation is highlighted in color, each with a different shade to better illustrate their position in the SIPO as the order is important in the correct calculation of the convolution.

## 2.2 FPGA Datapath

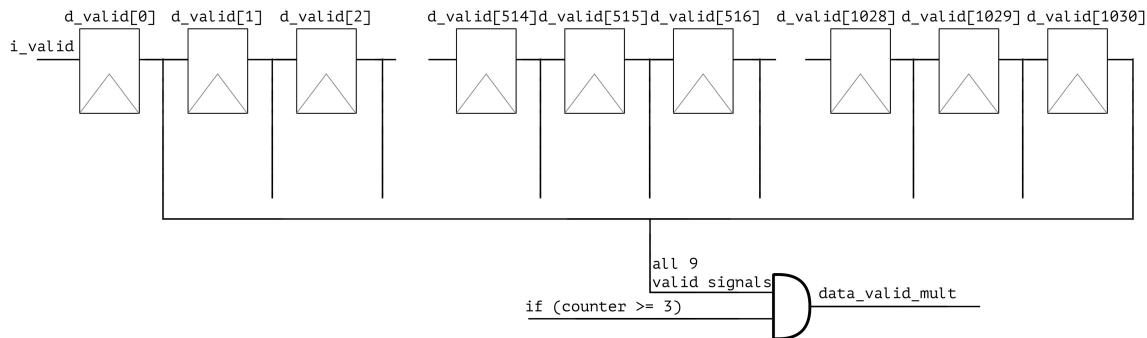


Figure 2: Data valid datapath

With the insight from Section 2.1, the datapaths for the FPGA implementation is created. A streaming dataflow model is chosen for the implementation and computations are pipelined as much as possible to enable the highest operating frequency possible.

The datapath for the computation is presented in Fig. 3. The datapath has 9 multipliers performing the 9 multiplications in parallel while the outputs are summed in an adder tree. Data is pipelined in between each stage of multiplier and adders.

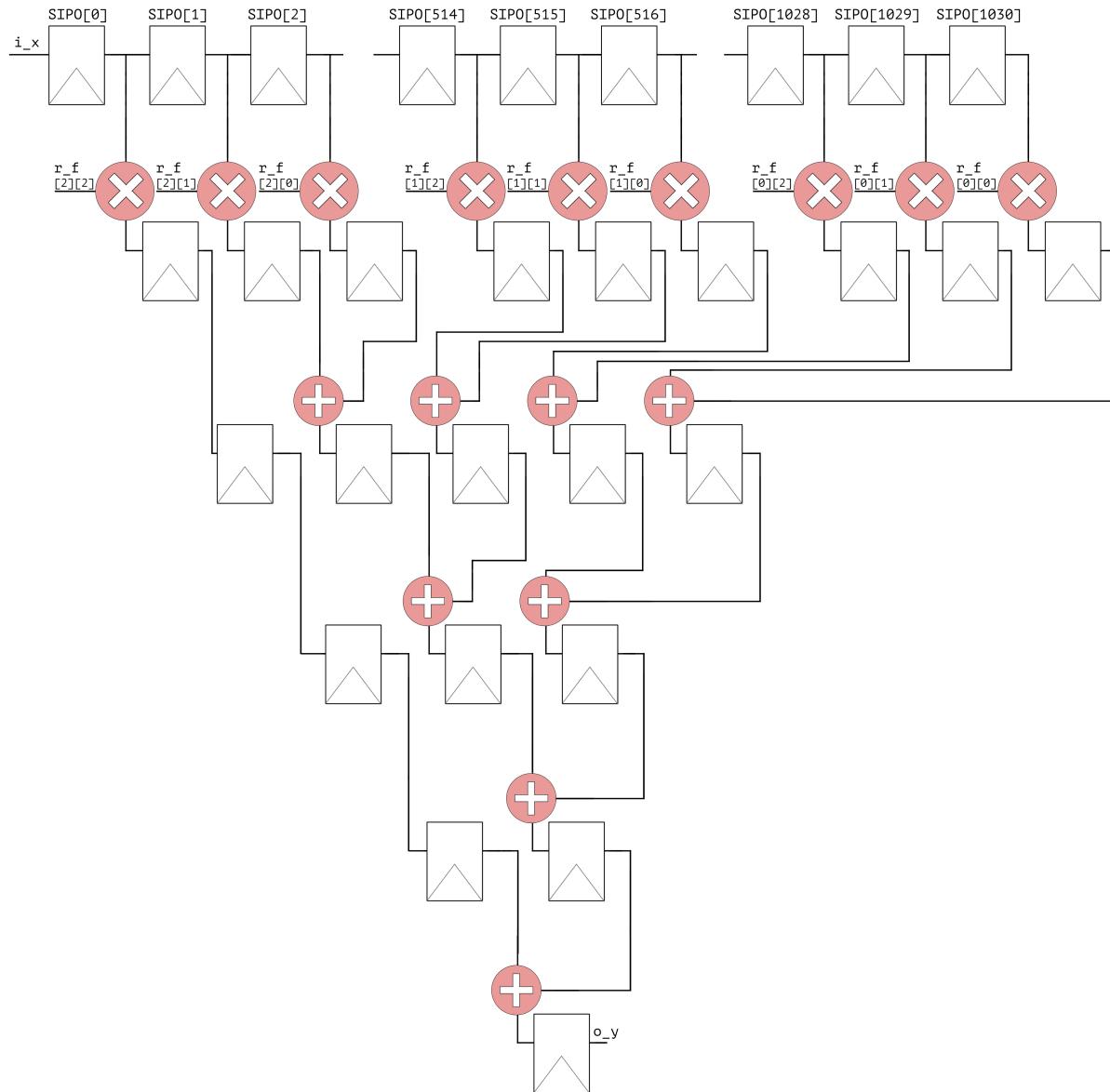


Figure 3: Data datapath

In a similar fashion, Fig. 2 shows a SIPO shift register for the data valid bit. From this, 9 bits corresponding to the 9 pixels involved in the convolution are connected to an AND gate along with a counter that keeps track of the current input's column position. While not shown explicitly, the `data_valid_mult` output signal from the AND gate is propagated along with the data from multiplier and adder stages in sync.

## 2.3 Operation

During operation, input data and its valid bit is fed into the SIPO shift register, beginning with the 1st position (SIPO[0], data\_valid[0]), as long as it is valid and the downstream circuit is ready (i.e. i\_valid and i\_ready both high). With each clock cycle, data in the shift register propagates down and the same happens for the valid bits that correspond to each pixel (i.e. SIPO[i] = SIPO[i-1], data\_valid[i] = data\_valid[i-1]).

	SIPO[1030] r_f[0][0] mult_out[0]	SIPO[1029] r_f[0][1] mult_out[1]	SIPO[1028] r_f[0][2] mult_out[2]
	SIPO[516] r_f[1][0] mult_out[3]	SIPO[515] r_f[1][1] mult_out[4]	SIPO[514] r_f[1][2] mult_out[5]
	SIPO[2] r_f[2][0] mult_out[6]	SIPO[1] r_f[2][1] mult_out[7]	SIPO[0] r_f[2][2] mult_out[8]

Figure 4: Convolution filter and pixel mapping

The multiply and accumulate portion of the circuit samples the 9 required pixels and the respective data valid bits from the SIPO shift registers according to Fig. 4 and performs the computations continuously regardless of the validity of the pixels sampled. Instead, the valid bits are checked separately to ensure that all 9 pixels are valid. In addition, an additional check is performed on a counter, one that keeps track of the column position of the current input. This is to ensure that convolution is not being done when the first two pixels of each row is inserted, since that would be invalid (i.e. instead of performing a 3x3 convolution, a 3x2 and 3x1 or 3x1 and 3x2 set of pixels spanning 4 rows will be used instead). Once the AND gate from Fig. 2 returns true, the output of the multipliers are now valid and this knowledge/status is passed along with the computed result via a separate set of valid bit registers that accompanies the multiplier and adder stages in sync with the computed values. Due to the fact that filter values are signed and pixel values are not. Pixel values are padded

```
// Pipeline reg for adder stage 4 out, final result, clipping performed here
if (addr_stage_4_out > 32'd255) begin
    addr_stage_4_out_pipe <= 8'd255;
end else if (addr_stage_4_out < 32'd0) begin
    addr_stage_4_out_pipe <= 8'd0;
end else begin
    addr_stage_4_out_pipe <= addr_stage_4_out[7:0];
end
end
```

Figure 5: Systemverilog clipping

with a 0 in the MSB position to convert it to signed and ensure that the multiplication performed is signed as well. Results are also stored in a sufficiently wide register to ensure no overflow occurs throughout the multiplier and adder stages. Once the final addition is performed, the values are clipped to be between 0 and 255 as specified by the assignment (Fig. 5).

## 2.4 Correctness

The FPGA implementation successfully completes all test cases (Appendix A).

Test	Completed?
1	Yes
2	Yes
3a	Yes
3b	Yes
4a	Yes
4b	Yes
5a	Yes
5b	Yes
6a	Yes
6b	Yes
7a	Yes
7b	Yes

Table 1: FPGA implementation correctness

## 2.5 Design Space Exploration

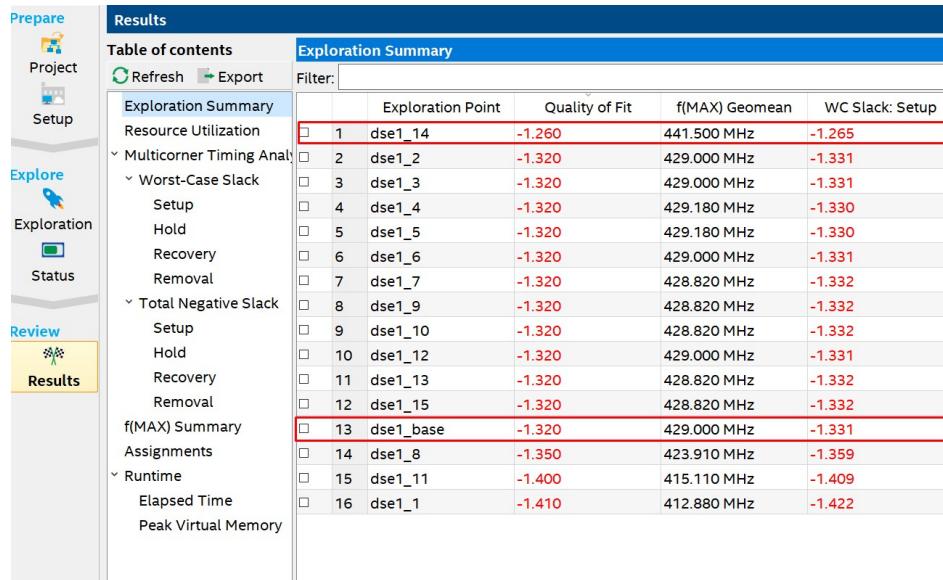


Figure 6: Design space exploration optimization

The design is further optimized using Intel's Design Space Explorer. A total of 15 different

seeds were explored in search of a better performance. The final design chosen is des1\_14 as shown in Fig. , which has a faster maximum operating frequency than the base compile.

## 2.6 Performance Metrics

	Result (Appendix B)
ALM Utilization	2497
DSP Utilization	9
BRAM (M20K) Utilization	0
Maximum Operating Frequency (MHz)	441.5
Cycles for Test 7a (Hinton)	264216
Dynamic Power for one module @ maximum frequency (W)	$111.7878 \times 10^{-3}$
Throughput of one module (GOPS)	7.4466
Throughput of a full device (GOPS)	$1.2510 \times 10^3$
Total Power for a full device (W)	11.98

Table 2: FPGA implementation results

Details regarding resource utilization and performance for the FPGA implementation is highlighted in Table 2, calculations are included in Appendix B.

## 3.0 CPU and GPU Implementations

### 3.1 Operations

The CPU and GPU implementations are similar in function in that both computes the convolution of the image and filter given. There are minute differences in the way that the computation is implemented due to the different hardware that the code will be executed on. On the CPU, since there are no special hardware available to perform the multiply and accumulate operations, for loops are used extensively to loop through the number of filters, rows and columns of the image and subsequently, depending on whether if the operation is vectorized, the filter as well. On the other hand, the implementation on GPU involves more optimization in terms of memory management as well as utilizing the CUDA cores efficiently, leading to less nested loops.

### 3.2 Performance Metrics

No. of filters	Runtime (ms) (Appendix C)			
	1	4	16	64
GPU	0.0297056	0.0994058	0.372583	1.47965
CPU (basic - no opt - 1 thread)	8.82178	35.0252	139.794	558.436
CPU (vectorized - no opt - 1 thread)	5.34393	21.2444	87.1801	340.192
CPU (basic - O2 - 1 thread)	1.59544	6.37035	25.4466	101.752
CPU (vectorized - O2 - 1 thread)	1.26072	5.05124	23.0627	83.6116
CPU (basic - O3 - 1 thread)	0.679928	2.72494	10.8303	43.27
CPU (vectorized - O3 - 1 thread)	1.25945	5.03768	23.019	83.5155
CPU (basic - O3 - 4 threads)	0.728284	0.823869	2.93434	11.8149
CPU (vectorized - O3 - 4 threads)	1.19086	1.71489	7.82391	22.2321

Table 3: Runtimes of CPU and GPU convolution implementations

From Table 3, with respect to CPU runtimes, we can see that for any given test case, runtimes increases with the number of filters roughly linearly. For any given optimization level, the vectorized version of convolution computation took less time to complete compared to the basic version unless the -O3 flag is enabled. This is due to more efficient calculations being done as shown by the fewer number of nested loops in Fig. 7. With respect to compilation optimizations (i.e. -O2 and -O3 flags), a higher level of optimization leads more aggressive optimizations being done by the compiler. These include unrolling more loops, and performing more aggressive instruction optimizations late in the compiler pipeline, all leading to lower runtimes. The basic version had more loops and likely benefited from the -O3 optimization more than the vectorized version, thus leading to the lower runtimes. Finally, when multithreading is enabled, significant performance gain is seen in runs with higher

**Assignment 2**

October 30, 2021

```

// CALCULATE THE Padded image size
int padded_image_width = image_width + (2*FILTER_RADIUS);

for(int filter_id = 0; filter_id < num_filters; filter_id++){ // For each filter
    for(int image_y = 0; image_y < image_height; image_y++){ // For each y location in the image
        for(int image_x = 0; image_x < image_width; image_x++){ // For each x location in the image
            // Calculate the dot-product of the 3 rows of the filter using the hand-vectorized dot3 function
            r1_dp = dot3(&filter[(filter_id * FILTER_SIZE * FILTER_SIZE)], &input_image[((image_y * padded_image_width) + image_x)];
            r2_dp = dot3(&filter[(filter_id * FILTER_SIZE * FILTER_SIZE) + FILTER_SIZE], &input_image[((image_y+1) * padded_image_width) + image_x]);
            r3_dp = dot3(&filter[(filter_id * FILTER_SIZE * FILTER_SIZE) + 2 * FILTER_SIZE], &input_image[((image_y+2) * padded_image_width) + image_x]);

            // Store the summation of the results of the 3 filter rows in the corresponding pixel location
            output_image[(filter_id * image_width * image_height) + (image_y * image_width) + image_x] = r1_dp + r2_dp + r3_dp;
        }
    }
}

// Calculate the padded image size
int padded_image_width = image_width + (2*FILTER_RADIUS);

for(int filter_id = 0; filter_id < num_filters; filter_id++){ // For each filter
    for(int image_y = 0; image_y < image_height; image_y++){ // For each y location of the output image
        for(int image_x = 0; image_x < image_width; image_x++){ // For each x location of the output image
            // Start a new accumulation result
            accum = 0;

            // Compute the dot-product of the filter and the input image
            for(int filter_y = 0; filter_y < FILTER_SIZE; filter_y++){
                for(int filter_x = 0; filter_x < FILTER_SIZE; filter_x++){
                    xx = image_x + filter_x;
                    yy = image_y + filter_y;
                    accum += input_image[(yy * padded_image_width) + xx] * filter[(filter_id * FILTER_SIZE * FILTER_SIZE) + (filter_y * FILTER_SIZE) + filter_x];
                }
            }

            // Store the convolution result in the corresponding pixel location
            output_image[(filter_id * image_width * image_height) + (image_y * image_width) + image_x] = accum;
        }
    }
}

```

Figure 7: CPU basic (right) vs vectorized (left) implementations

number of filters, reducing the runtime by up to 4 times. With more CPU resources being used, the amount of work done by each thread is reduced, thus leading to faster completion. Comparing between GPU and CPU runtimes, we can see that GPU outperformed CPU by a significant amount due to the GPU being much more optimized for the calculations needed in a convolution calculation. There are also many more processing units upon which work can be divided compared to the CPU (i.e. this is equivalent to running on CPU with many many more threads), thus the lower runtime.

## 4.0 Evaluations

Device	Model	Process Tech.	Die Size (mm <sup>2</sup> )	TDP (W)
CPU	Intel core i7-4790 (4 cores)	22nm	177	84
GPU	Nvidia GeForce GTX 980	28nm	396	165
FPGA	Arria 10 GX 1150	20nm	~350	70

Table 4: Specification of the three compute platforms

	Throughput (GOPS)	Power (W)	Energy Efficiency (GOPS/W)	Area Efficiency (GOPS/mm <sup>2</sup> )
FPGA (20nm)	$1.2510 \times 10^3$	11.98	104.42	3.5743
GPU (28nm)	15.786	165	0.095673	0.39864
GPU (scaled 20nm)	22.1	247.5	0.089293	0.055808
CPU (22nm)	163.72	84	1.9490	0.92497

Table 5: Device evaluations

Based on the results in Table 5, FPGA outperforms CPU and GPU by a wide margin. This is mostly due to the fact that the calculations for throughput is computed with only the operations that directly contribute to the convolutions. As a result, the FPGA which only performs computation directly related to the convolution has a significant advantage (due to the streaming dataflow model it uses) over both the CPU and GPU which involves lots of other operations that handles the control flow (loops, ifs) and other setup required for the convolution computations to be carried out.

In addition, power data used for the calculation in the CPU and GPU cases are the theoretical upper limit of those devices and it is very unlikely that these devices actually consumed this much power for the operation. The power data used for FPGA however, is much more accurate as it comes from a power analysis with realistic toggle rates based off of an ModelSim simulation.

Furthermore, the area numbers also contributed to the significant advantage shown by the FPGA. Relative to the CPU and GPU, much more of the area used in the calculation for FPGA is actually involved in the convolution. This is due to the lack of overhead hardware required by the CPU and GPU to stay more generic/general purpose.

Finally, the throughput and power numbers used for FPGA may also be more optimistic than reality since they are estimates from a single module. Performance will likely not scale linearly, especially as the FPGA approaches full utilization. More power might be consumed and clock frequency will decrease due to non optimal routing/layout.

With regards to estimation for GPU at a 20nm process, it is first worth while to note that while a smaller process means more transistors can be fit within the same area, companies

rarely keep die sizes the same. Die sizes often change going from process to process as well as within each process as the technology matures and architecture design demand changes. However, assuming that we keep die sizes the same, and that the decrease in process does not lead to any additional complications with power (such as significantly increased current leakage), we can then expect the number of transistors to increase approximately linearly. Assuming proper utilization, we can expect this increase to be passed on to the throughput as well. Power, on the other hand, might scale less than ideally than linearly due to Dennard scaling breaking down beginning in 2006. With an arbitrary factor of 50% more power than linear scaling, estimates in Table 5 are calculated in Appendix D.

## 5.0 Appendix

### Appendix A: FPGA Correctness

The screenshot shows the ModelSim simulation environment. The top half displays the test bench code in a text editor, and the bottom half shows the simulation transcript.

**Test Bench Code (lab2\_tb.v)**

```

// Define some filters for testing
logic [7:0] identity_filter; // Identity filter (output image is the same as input image)
initial identity_filter = {8'd0, 8'd0, 8'd0, 8'd0, 8'd0, 8'd0, 8'd0, 8'd0};

logic [7:0] edge_filter;
initial edge_filter = {-8'd1, -8'd1, -8'd1, -8'd1, -8'd1, -8'd1, -8'd1, -8'd1};

logic [7:0] sharpen_filter;
initial sharpen_filter = {8'd0, -8'd1, 8'd0, -8'd1, 8'd0, -8'd1, 8'd0, -8'd1};

// Define the test image and the golden results you want to use for testing
// PS: These tests are ordered in increasing difficulty and size. Following the same order when you
// test your design will make it easier to debug. You can uncomment the parameters of the test you
// want to use.

// Test 1: Delta Dot (3 x 512) & Identity filter
// This test uses a 3 x 512 black image with one white pixel in the middle
localparam TEST_IMAGE = {PROJECT_DIR, "/tests/01_delta_point.pgm"};
localparam GOLDEN_RES = {PROJECT_DIR, "/tests/01_delta_point_identity_golden.pgm"};
initial i_f = identity_filter;

// Test 2: Delta Line (3 x 512) & Identity filter
// This test uses a 3 x 512 black image with a white vertical line of pixels in the middle
localparam TEST_IMAGE = {PROJECT_DIR, "/tests/02_delta_line.pgm"};
localparam GOLDEN_RES = {PROJECT_DIR, "/tests/02_delta_line_identity_golden.pgm"};
initial i_f = identity_filter;

// Test 3a: Step (3 x 512) & Identity filter
// This test uses a 3 x 512 half black and half white image
localparam TEST_IMAGE = {PROJECT_DIR, "/tests/03_step.pgm"};
localparam GOLDEN_RES = {PROJECT_DIR, "/tests/03_step_identity_golden.pgm"};
initial i_f = identity_filter;

// Test 3b: Step (3 x 512) & Edge filter
// This test uses a 3 x 512 half black and half white image
localparam TEST_IMAGE = {PROJECT_DIR, "/tests/03_step.pgm"};
localparam GOLDEN_RES = {PROJECT_DIR, "/tests/03_step_edge_golden.pgm"};
initial i_f = edge_filter;

```

**Transcript**

```

# Pixel 1527 matching! Expected: 0 Got: 0
# Pixel 1528 matching! Expected: 0 Got: 0
# Pixel 1529 matching! Expected: 0 Got: 0
# Pixel 1530 matching! Expected: 0 Got: 0
# Pixel 1531 matching! Expected: 0 Got: 0
# Pixel 1532 matching! Expected: 0 Got: 0
# Pixel 1533 matching! Expected: 0 Got: 0
# Pixel 1534 matching! Expected: 0 Got: 0
# Pixel 1535 matching! Expected: 0 Got: 0
# TEST PASSED! ALL PIXELS ARE MATCHING GOLDEN RESULT!
# Break in Module lab2_tb at lab2_tb.v line 371
VSM 4> run -over

```

Figure 8: FPGA test 1 test bench

# Assignment 2

ECE 1756  
October 30, 2021

```

28 logic [71:0] edge_filter; // Edge detection filter
29 initial edge_filter = {-8'd1, -8'd1, -8'd1, -8'd1, 8'd8, -8'd1, -8'd1, -8'd1};
30 logic [71:0] sharpen_filter; // Image sharpening filter
31 initial sharpen_filter = {8'd0, -8'd1, 8'd0, -8'd1, 8'd5, -8'd1, 8'd0, -8'd1, 8'd0};
32
33 // Define the test image and the golden results you want to use for testing
34 // PS: These tests are ordered in increasing difficulty and size. Following the same order when you
35 // test your design will make it easier to debug. You can uncomment the parameters of the test you
36 // want to use.
37
38 // Test 1: Delta Dot (3 x 512) * Identity filter
39 // This test uses a 3 x 512 black image with one white pixel in the middle
40 localparam TEST_IMAGE = {PROJECT_DIR, "/tests/01_delta_point.pgm"};
41 localparam GOLDEN_RES = {PROJECT_DIR, "/tests/01_delta_point_identity_golden.pgm"};
42 initial i_f = identity_filter;
43
44 // Test 2: Delta Line (3 x 512) * Identity filter
45 // This test uses a 3 x 512 black image with a white vertical line of pixels in the middle
46 localparam TEST_IMAGE = {PROJECT_DIR, "/tests/02_delta_line.pgm"};
47 localparam GOLDEN_RES = {PROJECT_DIR, "/tests/02_delta_line_identity_golden.pgm"};
48 initial i_f = identity_filter;
49
50 // Test 3a: Step (3 x 512) * Identity filter
51 // This test uses a 3 x 512 half black and half white image
52 localparam TEST_IMAGE = {PROJECT_DIR, "/tests/03_step.pgm"};
53 localparam GOLDEN_RES = {PROJECT_DIR, "/tests/03_step_identity_golden.pgm"};
54 initial i_f = identity_filter;
55
56 // Test 3b: Step (3 x 512) * Edge filter
57 // This test uses a 3 x 512 half black and half white image
58 localparam TEST_IMAGE = {PROJECT_DIR, "/tests/03_step.pgm"};
59 localparam GOLDEN_RES = {PROJECT_DIR, "/tests/03_step_edge_golden.pgm"};
60 initial i_f = edge_filter;
61
62 // Test 4a: Pixel Checker Board (3 x 512) * Identity filter
63 // This test uses a 3 x 512 image with alternating black and white pixels in a checker board pattern

```

Transcript:

```

# Pixel 1525 matching! Expected: 0 Got: 0
# Pixel 1526 matching! Expected: 0 Got: 0
# Pixel 1527 matching! Expected: 0 Got: 0
# Pixel 1528 matching! Expected: 0 Got: 0
# Pixel 1529 matching! Expected: 0 Got: 0
# Pixel 1530 matching! Expected: 0 Got: 0
# Pixel 1531 matching! Expected: 0 Got: 0
# Pixel 1532 matching! Expected: 0 Got: 0
# Pixel 1533 matching! Expected: 0 Got: 0
# Pixel 1534 matching! Expected: 0 Got: 0
# Pixel 1535 matching! Expected: 0 Got: 0
# TEST PASSED! ALL PIXELS ARE MATCHING GOLDEN RESULT!
# Break in Module lab2_tb at lab2_tb.sv line 371
VSIM 10>
Now: $1.798 ns Delta: 0 sim:lab2_tb#INITIAL#267

```

Figure 9: FPGA test 2 test bench

```

43 // Test 2: Delta Line (3 x 512) * Identity filter
44 // This test uses a 3 x 512 black image with a white vertical line of pixels in the middle
45 localparam TEST_IMAGE = {PROJECT_DIR, "/tests/02_delta_line.pgm"};
46 localparam GOLDEN_RES = {PROJECT_DIR, "/tests/02_delta_line_identity_golden.pgm"};
47 initial i_f = identity_filter;
48
49 // Test 3a: Step (3 x 512) * Identity filter
50 // This test uses a 3 x 512 half black and half white image
51 localparam TEST_IMAGE = {PROJECT_DIR, "/tests/03_step.pgm"};
52 localparam GOLDEN_RES = {PROJECT_DIR, "/tests/03_step_identity_golden.pgm"};
53 initial i_f = identity_filter;
54
55 // Test 3b: Step (3 x 512) * Edge filter
56 // This test uses a 3 x 512 half black and half white image
57 localparam TEST_IMAGE = {PROJECT_DIR, "/tests/03_step.pgm"};
58 localparam GOLDEN_RES = {PROJECT_DIR, "/tests/03_step_edge_golden.pgm"};
59 initial i_f = edge_filter;
60
61 // Test 4a: Pixel Checker Board (3 x 512) * Identity filter
62 // This test uses a 3 x 512 image with alternating black and white pixels in a checker board pattern

```

Transcript:

```

# Pixel 1525 matching! Expected: 255 Got: 255
# Pixel 1526 matching! Expected: 255 Got: 255
# Pixel 1527 matching! Expected: 255 Got: 255
# Pixel 1528 matching! Expected: 255 Got: 255
# Pixel 1529 matching! Expected: 255 Got: 255
# Pixel 1530 matching! Expected: 255 Got: 255
# Pixel 1531 matching! Expected: 255 Got: 255
# Pixel 1532 matching! Expected: 255 Got: 255
# Pixel 1533 matching! Expected: 255 Got: 255
# TEST PASSED! ALL PIXELS ARE MATCHING GOLDEN RESULT!
# Break in Module lab2_tb at lab2_tb.sv line 371
VSIM 13>
Now: $1.798 ns Delta: 0 sim:lab2_tb#INITIAL#267

```

Figure 10: FPGA test 3a test bench

# Assignment 2

ECE 1756  
October 30, 2021

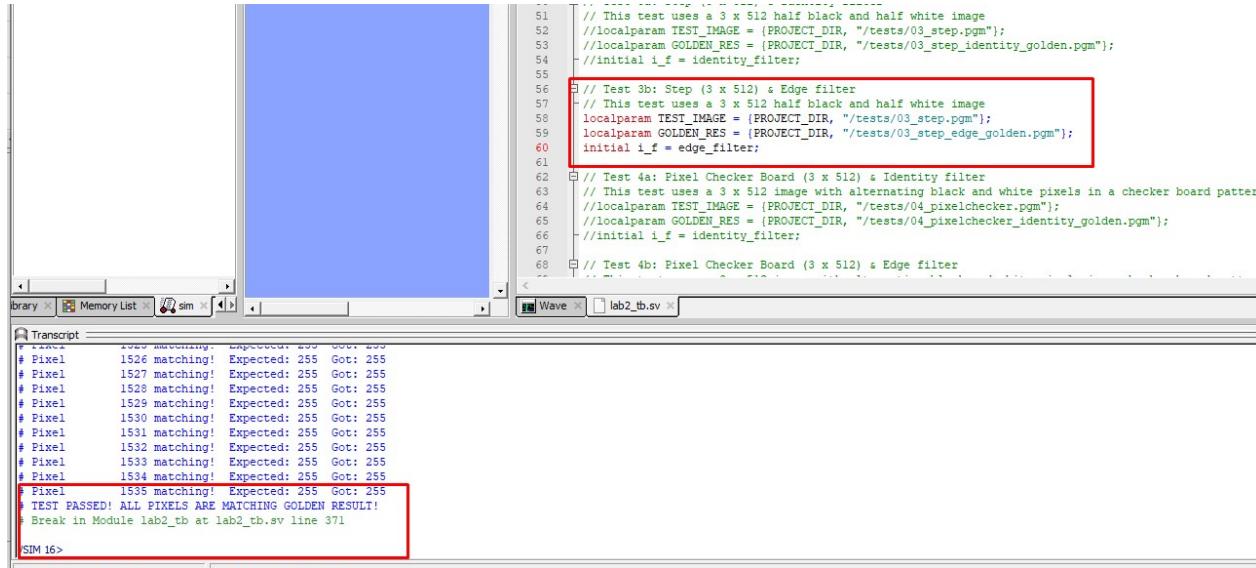


Figure 11: FPGA test 3b test bench

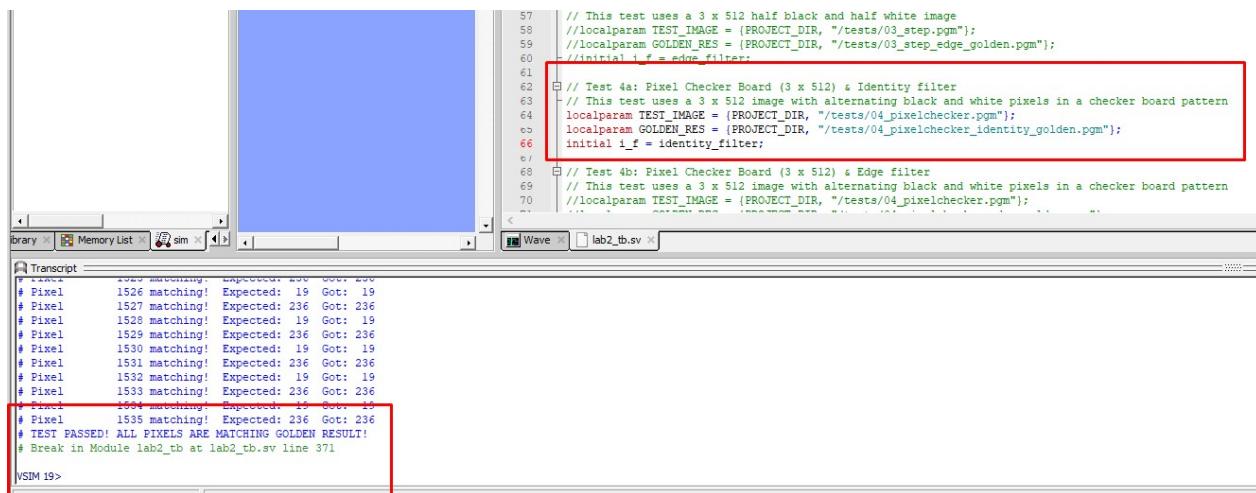
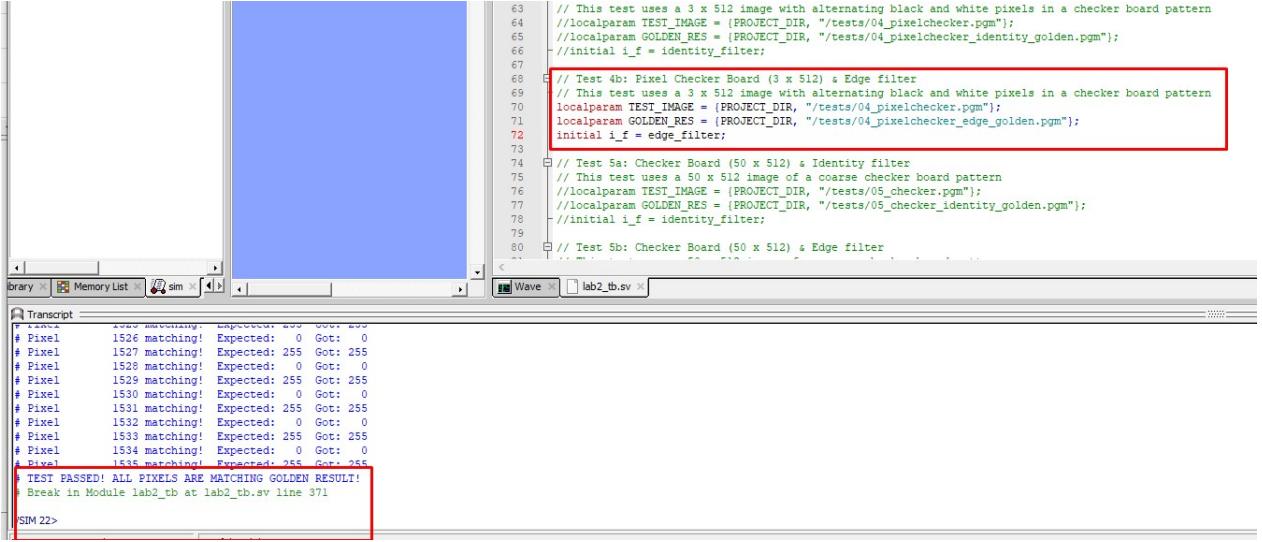


Figure 12: FPGA test 4a test bench

**Assignment 2**

October 30, 2021



The screenshot shows a simulation interface with a waveform viewer and a transcript window. The transcript window displays the results of a pixel comparison between the test image and the golden result. The results show that all pixels matched correctly.

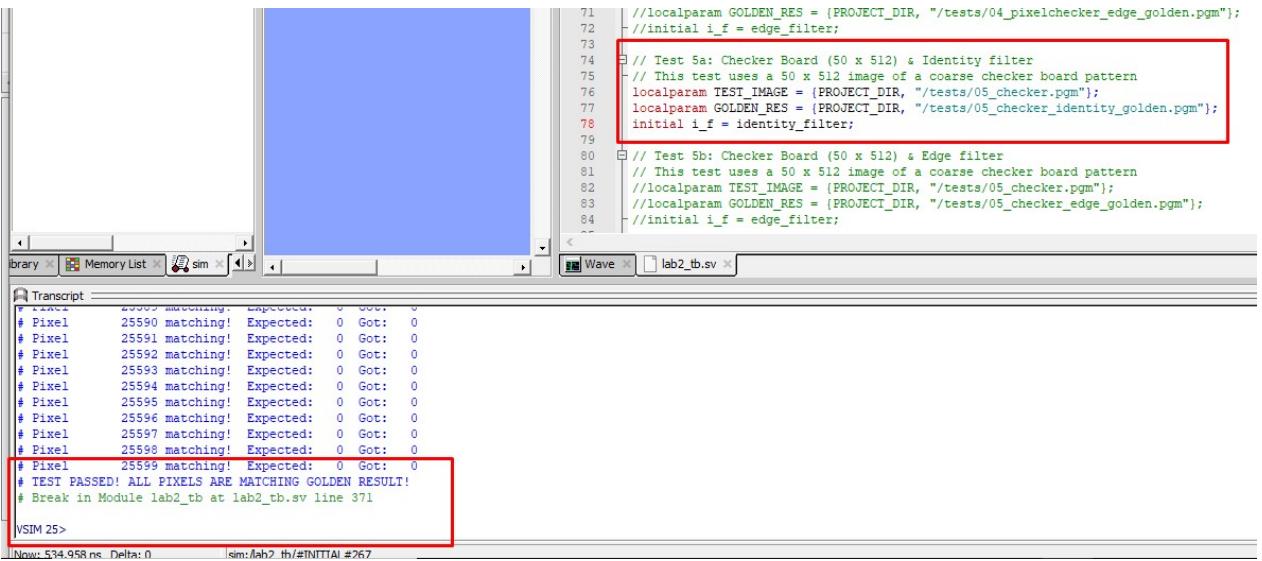
```

63 // This test uses a 3 x 512 image with alternating black and white pixels in a checker board pattern
64 //localparam TEST_IMAGE = [PROJECT_DIR, "/tests/04_pixelchecker.pgm"];
65 //localparam GOLDEN_RES = [PROJECT_DIR, "/tests/04_pixelchecker_identity_golden.pgm"];
66 //initial i_f = identity_filter;
67
68 // Test 4b: Pixel Checker Board (3 x 512) & Edge filter
69 // This test uses a 3 x 512 image with alternating black and white pixels in a checker board pattern
70 localparam TEST_IMAGE = [PROJECT_DIR, "/tests/04_pixelchecker.pgm"];
71 localparam GOLDEN_RES = [PROJECT_DIR, "/tests/04_pixelchecker_edge_golden.pgm"];
72 initial i_f = edge_filter;
73
74 // Test 5a: Checker Board (50 x 512) & Identity filter
75 // This test uses a 50 x 512 image of a coarse checker board pattern
76 //localparam TEST_IMAGE = [PROJECT_DIR, "/tests/05_checker.pgm"];
77 //localparam GOLDEN_RES = [PROJECT_DIR, "/tests/05_checker_identity_golden.pgm"];
78 //initial i_f = identity_filter;
79
80 // Test 5b: Checker Board (50 x 512) & Edge filter
81 // This test uses a 50 x 512 image of a coarse checker board pattern
82 //localparam TEST_IMAGE = [PROJECT_DIR, "/tests/05_checker.pgm"];
83 //localparam GOLDEN_RES = [PROJECT_DIR, "/tests/05_checker_edge_golden.pgm"];
84 //initial i_f = edge_filter;

# Pixel 1526 matching! Expected: 0 Got: 0
# Pixel 1527 matching! Expected: 255 Got: 255
# Pixel 1528 matching! Expected: 0 Got: 0
# Pixel 1529 matching! Expected: 255 Got: 255
# Pixel 1530 matching! Expected: 0 Got: 0
# Pixel 1531 matching! Expected: 255 Got: 255
# Pixel 1532 matching! Expected: 0 Got: 0
# Pixel 1533 matching! Expected: 255 Got: 255
# Pixel 1534 matching! Expected: 0 Got: 0
# Pixel 1535 matching! Expected: 255 Got: 255
TEST PASSED! ALL PIXELS ARE MATCHING GOLDEN RESULT!
Break in Module lab2_tb at lab2_tb.sv line 371
SIM 22>

```

Figure 13: FPGA test 4b test bench



The screenshot shows a simulation interface with a waveform viewer and a transcript window. The transcript window displays the results of a pixel comparison between the test image and the golden result. The results show that all pixels matched correctly.

```

71 //localparam GOLDEN_RES = [PROJECT_DIR, "/tests/04_pixelchecker_edge_golden.pgm"];
72 //initial i_f = edge_filter;
73
74 // Test 5a: Checker Board (50 x 512) & Identity filter
75 // This test uses a 50 x 512 image of a coarse checker board pattern
76 localparam TEST_IMAGE = [PROJECT_DIR, "/tests/05_checker.pgm"];
77 localparam GOLDEN_RES = [PROJECT_DIR, "/tests/05_checker_identity_golden.pgm"];
78 initial i_f = identity_filter;
79
80 // Test 5b: Checker Board (50 x 512) & Edge filter
81 // This test uses a 50 x 512 image of a coarse checker board pattern
82 //localparam TEST_IMAGE = [PROJECT_DIR, "/tests/05_checker.pgm"];
83 //localparam GOLDEN_RES = [PROJECT_DIR, "/tests/05_checker_edge_golden.pgm"];
84 //initial i_f = edge_filter;

# Pixel 25500 matching! Expected: 0 Got: 0
# Pixel 25501 matching! Expected: 0 Got: 0
# Pixel 25502 matching! Expected: 0 Got: 0
# Pixel 25503 matching! Expected: 0 Got: 0
# Pixel 25504 matching! Expected: 0 Got: 0
# Pixel 25505 matching! Expected: 0 Got: 0
# Pixel 25506 matching! Expected: 0 Got: 0
# Pixel 25507 matching! Expected: 0 Got: 0
# Pixel 25508 matching! Expected: 0 Got: 0
# Pixel 25509 matching! Expected: 0 Got: 0
TEST PASSED! ALL PIXELS ARE MATCHING GOLDEN RESULT!
Break in Module lab2_tb at lab2_tb.sv line 371
VSIM 25>
Now: 534.958 ns Delta: 0 sim:lab2_tb#INITIAL #267

```

Figure 14: FPGA test 5a test bench

# Assignment 2

ECE 1756  
October 30, 2021

```

76 // localparam TEST_IMAGE = [PROJECT_DIR, "/tests/05_checker.pgm"];
77 // localparam GOLDEN_RES = [PROJECT_DIR, "/tests/05_checker_identity_golden.pgm"];
78 //initial i_f = identity_filter;
79
80 // // Test 5b: Checker Board (50 x 512) & Edge filter
81 // This test uses a 50 x 512 image of a coarse checker board pattern
82 localparam TEST_IMAGE = [PROJECT_DIR, "/tests/05_checker.pgm"];
83 localparam GOLDEN_RES = [PROJECT_DIR, "/tests/05_checker_identity_golden.pgm"];
84 initial i_f = edge_filter;
85
86 // // Test 6a: Small Hinton (50 x 512) & Edge filter
87 // This test uses a 50 x 512 snippet of the Geoffrey Hinton's image
88 //localparam TEST_IMAGE = [PROJECT_DIR, "/tests/06_hinton_small.pgm"];
89 //localparam GOLDEN_RES = [PROJECT_DIR, "/tests/06_hinton_small_edge_golden.pgm"];
90 //initial i_f = edge_filter;
91
92 // // Test 6b: Small Hinton (50 x 512) & Sharpen filter
93 // This test uses a 50 x 512 snippet of the Geoffrey Hinton's image
94 //localparam TEST_IMAGE = [PROJECT_DIR, "/tests/06_hinton_small.pgm"];
95 //localparam GOLDEN_RES = [PROJECT_DIR, "/tests/06_hinton_small_sharpen_golden.pgm"];
96 //initial i_f = sharpen_filter;

```

**Transcript**

```

# ----- 25590 matching! Expected: 0 Got: 0
# Pixel 25591 matching! Expected: 0 Got: 0
# Pixel 25592 matching! Expected: 0 Got: 0
# Pixel 25593 matching! Expected: 0 Got: 0
# Pixel 25594 matching! Expected: 0 Got: 0
# Pixel 25595 matching! Expected: 0 Got: 0
# Pixel 25596 matching! Expected: 0 Got: 0
# Pixel 25597 matching! Expected: 0 Got: 0
# Pixel 25598 matching! Expected: 0 Got: 0
# Pixel 25599 matching! Expected: 0 Got: 0
# TEST PASSED! ALL PIXELS ARE MATCHING GOLDEN RESULT!
# Break in Module lab2_tb at lab2_tb.sv line 371
VSIM 28>

```

Figure 15: FPGA test 5b test bench

```

83 // localparam TEST_IMAGE = [PROJECT_DIR, "/tests/05_checker.pgm"];
84 //initial i_f = edge_filter;
85
86 // // Test 6a: Small Hinton (50 x 512) & Edge filter
87 // This test uses a 50 x 512 snippet of the Geoffrey Hinton's image
88 localparam TEST_IMAGE = [PROJECT_DIR, "/tests/06_hinton_small.pgm"];
89 localparam GOLDEN_RES = [PROJECT_DIR, "/tests/06_hinton_small_edge_golden.pgm"];
90 initial i_f = edge_filter;
91
92 // // Test 6b: Small Hinton (50 x 512) & Sharpen filter
93 // This test uses a 50 x 512 snippet of the Geoffrey Hinton's image
94 //localparam TEST_IMAGE = [PROJECT_DIR, "/tests/06_hinton_small.pgm"];
95 //localparam GOLDEN_RES = [PROJECT_DIR, "/tests/06_hinton_small_sharpen_golden.pgm"];
96 //initial i_f = sharpen_filter;

```

**Transcript**

```

# ----- 25592 matching! Expected: 255 Got: 255
# Pixel 25593 matching! Expected: 255 Got: 255
# Pixel 25594 matching! Expected: 255 Got: 255
# Pixel 25595 matching! Expected: 255 Got: 255
# Pixel 25596 matching! Expected: 255 Got: 255
# Pixel 25597 matching! Expected: 255 Got: 255
# Pixel 25598 matching! Expected: 255 Got: 255
# Pixel 25599 matching! Expected: 255 Got: 255
# TEST PASSED! ALL PIXELS ARE MATCHING GOLDEN RESULT!
# Break in Module lab2_tb at lab2_tb.sv line 371
VSIM 31> run -over
# Next activity is in 2 ns.
VSIM 32>

```

Figure 16: FPGA test 6a test bench

**Assignment 2**

ECE 1756  
October 30, 2021

```

87 // This test uses a 50 x 512 snippet of the Geoffrey Hinton's image
88 //localparam TEST_IMAGE = (PROJECT_DIR, "/tests/06_hinton_small.pgm");
89 //localparam GOLDEN_RES = (PROJECT_DIR, "/tests/06_hinton_small_edge_golden.pgm");
90 //initial i_f = edge_filter;
91
92 // Test 6b: Small Hinton (50 x 512) & Sharpen filter
93 // This test uses a 50 x 512 snippet of the Geoffrey Hinton's image
94 localparam TEST_IMAGE = (PROJECT_DIR, "/tests/06_hinton_small.pgm");
95 localparam GOLDEN_RES = (PROJECT_DIR, "/tests/06_hinton_small_sharpen_golden.pgm");
96 initial i_f = sharpen_filter;
97
98 // Test 7a: Hinton (512 x 512) & Edge filter
99 // This test uses a full 512 x 512 image of Geoffrey Hinton
100 localparam TEST_IMAGE = (PROJECT_DIR, "/tests/07_hinton.pgm");
101 localparam GOLDEN_RES = (PROJECT_DIR, "/tests/07_hinton_edge_golden.pgm");
102 //initial i_f = edge_filter;
103
104 // Test 7b: Hinton (512 x 512) & Sharpen filter
105 // This test uses a full 512 x 512 image of Geoffrey Hinton
106 localparam TEST_IMAGE = (PROJECT_DIR, "/tests/07_hinton.pgm");
107 localparam GOLDEN_RES = (PROJECT_DIR, "/tests/07_hinton_sharpen_golden.pgm");
108 initial i_f = sharpen_filter;
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
589>

```

Figure 17: FPGA test 6b test bench

```

93 //localparam GOLDEN_RES = (PROJECT_DIR, "/tests/06_hinton_small_sharpen_golden.pgm");
94 initial i_f = sharpen_filter;
95
96 // Test 7a: Hinton (512 x 512) & Edge filter
97 // This test uses a full 512 x 512 image of Geoffrey Hinton
98 localparam TEST_IMAGE = (PROJECT_DIR, "/tests/07_hinton.pgm");
99 localparam GOLDEN_RES = (PROJECT_DIR, "/tests/07_hinton_edge_golden.pgm");
100 initial i_f = edge_filter;
101
102 // Test 7b: Hinton (512 x 512) & Sharpen filter
103 // This test uses a full 512 x 512 image of Geoffrey Hinton
104 localparam TEST_IMAGE = (PROJECT_DIR, "/tests/07_hinton.pgm");
105 localparam GOLDEN_RES = (PROJECT_DIR, "/tests/07_hinton_sharpen_golden.pgm");
106 initial i_f = sharpen_filter;
107
108
109
110
111
112
113
114
115
116
117
118
119
119>

```

Figure 18: FPGA test 7a test bench

The screenshot shows a ModelSim simulation interface. The top window displays a portion of a Verilog testbench script. Lines 101 through 115 are highlighted with a red box. The code includes parameters for golden results and test images, and logic for generating a 50MHz clock.

```
101 //localparam GOLDEN_RES = {PROJECT_DIR, "/tests/07_hinton_edge_golden.pgm"};
102 //initial i_f = edge_filter;
103
104 // Test 7b: Hinton (512 x 512) & Sharpen filter
105 // This test uses a full 512 x 512 image of Geoffrey Hinton
106 localparam TEST_IMAGE = {PROJECT_DIR, "/tests/07_hinton.pgm"};
107 localparam GOLDEN_RES = {PROJECT_DIR, "/tests/07_hinton_sharpen_golden.pgm"};
108 initial i_f = sharpen_filter;
109
110
111 // Generate a 50MHz clock
112 initial clk = 1'b1;
113 always #(CLK_PERIOD/2) clk = ~clk;
114
115 // This is the function to calculate the error between the produced output and gold
```

The bottom window is a transcript window showing the results of the pixel-by-pixel comparison. A red box highlights the final test result and the overall test status.

```
# Pixel 262135 matching! Expected: 46 Got: 46
# Pixel 262134 matching! Expected: 55 Got: 55
# Pixel 262135 matching! Expected: 46 Got: 46
# Pixel 262136 matching! Expected: 51 Got: 51
# Pixel 262137 matching! Expected: 63 Got: 63
# Pixel 262138 matching! Expected: 38 Got: 38
# Pixel 262139 matching! Expected: 44 Got: 44
# Pixel 262140 matching! Expected: 46 Got: 46
# Pixel 262141 matching! Expected: 42 Got: 42
# Pixel 262142 matching! Expected: 46 Got: 46
# Pixel 262143 matching! Expected: 50 Got: 50
# TEST PASSED! ALL PIXELS ARE MATCHING GOLDEN RESULT!
# Break in Module lab2_tb at lab2_tb.sv line 371
```

VSIM 41>

Figure 19: FPGA test 7b test bench

## Appendix B: FPGA Implementation Results Calculations

### ALM, DSP, BRAM Utilization

Table of Contents		Fitter Summary	
> Multiplexer Statistics		<<Filter>>	
> Source Assignments		Fitter Status	Successful - Thu Oct 28 19:19:47 2021
Post-Synthesis Netlist Statistics for Top Entity		Quartus Prime Version	20.1.1 Build 720 11/11/2020 SJ Standard Edition
Elapsed Time Per Partition		Revision Name	lab2
Messages		Top-level Entity Name	lab2
Suppressed Messages		Family	Arria 10
Fitter		Device	10AX115N1F45I1SG
Summary		Timing Models	Final
Settings		Logic utilization (in ALMs)	2,497 / 427,200 (< 1 %)
Parallel Compilation		Total registers	9572
Netlist Optimizations		Total pins	1 / 992 (< 1 %)
Estimated Delay Added for Hold Timing		Total virtual pins	93
Summary		Total block memory bits	0 / 55,562,240 (0 %)
Details		Total RAM Blocks	0 / 2,713 (0 %)
> Incremental Compilation Section		Total DSP Blocks	9 / 1,518 (< 1 %)
Pin-Out File		Total HSSI RX channels	0 / 48 (0 %)
> Resource Section		Total HSSI TX channels	0 / 48 (0 %)
Device Options		Total PLLs	0 / 112 (0 %)
Operating Settings and Conditions			
Messages			
Suppressed Messages			
> Assembler			

Figure 20: FPGA ALM, DSP and BRAM utilization

**Maximum Operating Frequency**

Type ID Message

- ① Analyzing Slow 900mV 100C Model
- > ▲ 332148 Timing requirements not met
- > ① 332146 Worst-case setup slack is -1.211
- > ① 332146 Worst-case hold slack is 0.046
- ① 332140 No Recovery paths to report
- ① 332140 No Removal paths to report
- > ① 332146 Worst-case minimum pulse width slack is -1.150
- ① Analyzing Slow 900mV 0C Model
- > ① 332146 Worst-case setup slack is -1.265
- > ① 332146 Worst-case hold slack is 0.045
- ① 332140 No Recovery paths to report

Messages System Processing (120)

Figure 21: FPGA setup and hold times

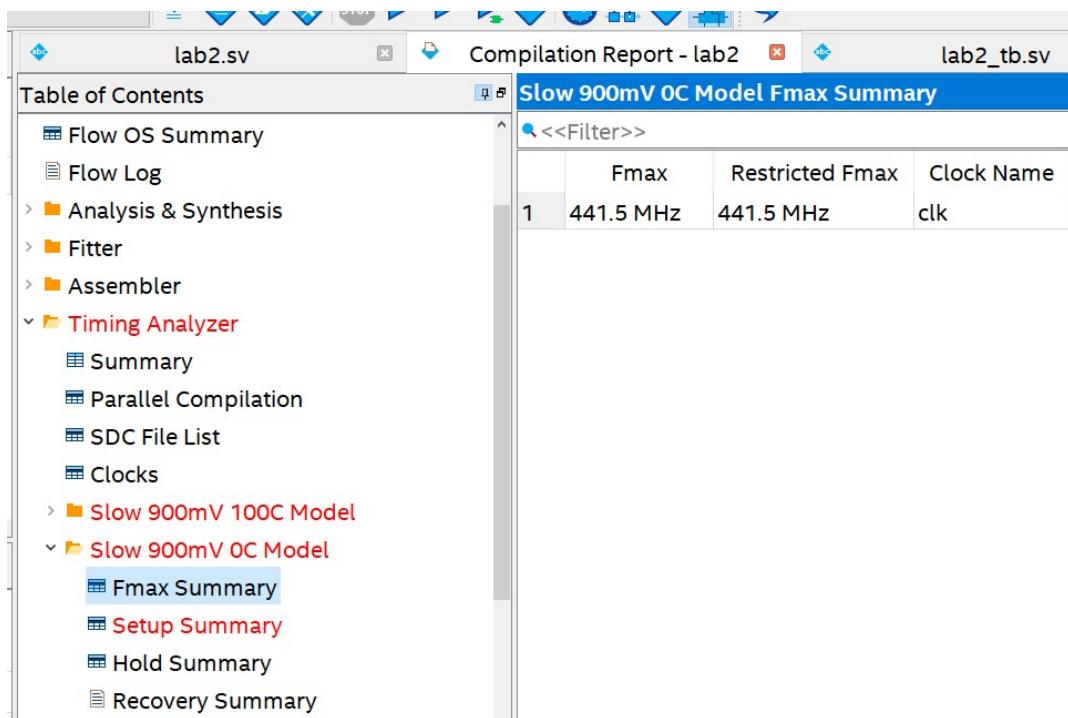


Figure 22: FPGA maximum operating frequency

Cycles for Test 7a (Hinton)

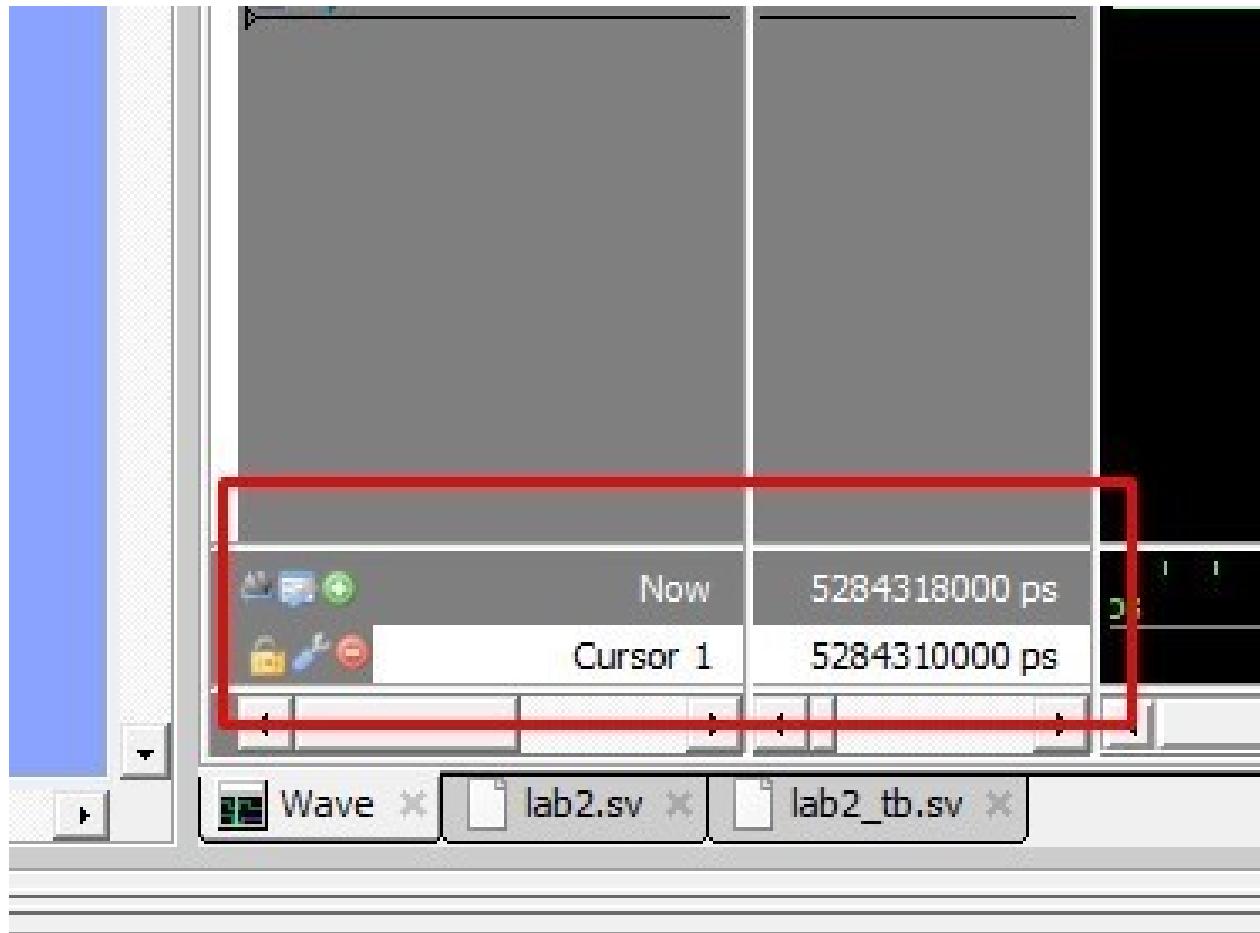


Figure 23: Test 7a total elapsed time

$$\text{No. of cycles} = 5284318000 / 20000 = 264216 \text{ cycles}$$

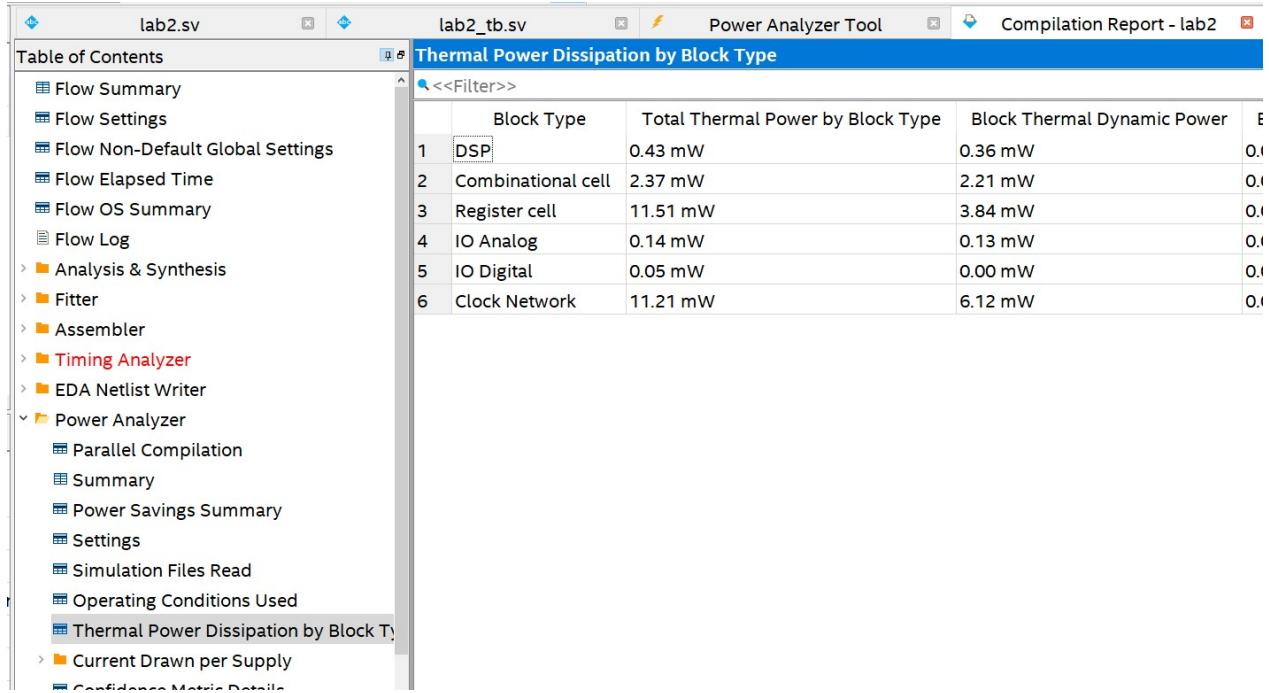
**Dynamic Power for one module @ 441.5MHz**

Figure 24: Power dissipation by block type

$$\text{Total Dynamic Power @}50\text{MHz} = 0.36 + 2.21 + 3.84 + 0.13 + 6.12 = 12.66 \text{ mW}$$

$$\text{Total Dynamic Power @}441.5\text{MHz} = 12.66 * (441.5/50) = 111.7878 \text{ mW} = 111.7878 \times 10^{-3} \text{ W}$$

**Throughput of one module (GOPS)**

To compute the convolution of one pixel, 9 multiplications and a total of 8 additions needs to be performed over a period of 5 cycles (output ready by the 6th cycle). For a 512x512 picture as is given in test 7a, that would take a total of:

$$(512 * 512) * (9 + 8) = 4456448 \text{ operations}$$

With a frequency of 441.5 MHz, 264216 cycles is equal to:

$$(1/441.5 \times 10^6) * 264216 \approx 5.9845 \times 10^{-4} \text{ seconds}$$

Which give us a total of:

$$4456448/5.9845 \times 10^{-4} \approx 7.4466 \times 10^9 = 7.4466 \text{ GOPs}$$

### Throughput of a full device (GOPS)

$$\text{Max. copies (based on ALM count)} = 427200/2497 \approx 171$$

$$\text{Max. copies (based on DSP count)} = 1518/9 \approx 168$$

Since DSP count is the limiting factor, approximate maximum number of copies per device is 168. As a result, the throughput of a full device:

$$7.4466 \times 10^9 * 168 \approx 1.2510 \times 10^{12} = 1.2510 \times 10^3 \text{ GOPs}$$

### Total Power for a full device (W)

Assuming a fully packed device with 168 modules as calculated before (note that IO and clock dynamic power are not multiplied):

$$\text{Total Dynamic Power @50MHz} = (0.36 + 2.21 + 3.84) * 168 + 0.13 + 6.12 = 1083.13 \text{ mW}$$

$$\text{Total Dynamic Power @441.5MHz} = 1083.13 * (441.5/50) = 9564.0379 \text{ mW}$$

$$\begin{aligned} \text{Total Power @441.5MHz} &= 9564.0379 + (0.43 + 2.37 + 11.51) * 168 + 0.14 + 0.05 + 11.21 \\ &= 11979.5179 \text{ mW} \approx 11.980 \text{ W} \end{aligned}$$

## Appendix C: Runtimes for CPU and GPU implementations

```
ug60:~/ece1756/lab2_cpu_gpu% rm conv_cpu
/bin/rm: remove regular file 'conv_cpu'? y
ug60:~/ece1756/lab2_cpu_gpu% g++ -mavx conv_cpu.cpp utils.cpp -o conv_cpu
ug60:~/ece1756/lab2_cpu_gpu% ./run_cpu.sh 1 4 16 64
Runtime of simple convolution for 1 filter(s) averaged accross 20 runs = 8.82178 ms
Runtime of hand vectorized convolution for 1 filter(s) averaged accross 20 runs = 5.34393 ms
-----
Runtime of simple convolution for 4 filter(s) averaged accross 20 runs = 35.0252 ms
Runtime of hand vectorized convolution for 4 filter(s) averaged accross 20 runs = 21.2444 ms
-----
Runtime of simple convolution for 16 filter(s) averaged accross 20 runs = 139.794 ms
Runtime of hand vectorized convolution for 16 filter(s) averaged accross 20 runs = 87.1801 ms
-----
Runtime of simple convolution for 64 filter(s) averaged accross 20 runs = 558.436 ms
Runtime of hand vectorized convolution for 64 filter(s) averaged accross 20 runs = 340.192 ms
-----
ug60:~/ece1756/lab2_cpu_gpu% █
```

Figure 25: CPU no optimization runtimes

```
ug60:~/ece1756/lab2_cpu_gpu% rm conv_cpu
/bin/rm: remove regular file 'conv_cpu'? y
ug60:~/ece1756/lab2_cpu_gpu% g++ -O2 conv_cpu.cpp utils.cpp -o conv_cpu
ug60:~/ece1756/lab2_cpu_gpu% ./run_cpu.sh 1 4 16 64
Runtime of simple convolution for 1 filter(s) averaged accross 20 runs = 1.59554 ms
Runtime of hand vectorized convolution for 1 filter(s) averaged accross 20 runs = 1.26072 ms
-----
Runtime of simple convolution for 4 filter(s) averaged accross 20 runs = 6.37035 ms
Runtime of hand vectorized convolution for 4 filter(s) averaged accross 20 runs = 5.05124 ms
-----
Runtime of simple convolution for 16 filter(s) averaged accross 20 runs = 25.4466 ms
Runtime of hand vectorized convolution for 16 filter(s) averaged accross 20 runs = 23.0627 ms
-----
Runtime of simple convolution for 64 filter(s) averaged accross 20 runs = 101.752 ms
Runtime of hand vectorized convolution for 64 filter(s) averaged accross 20 runs = 83.6116 ms
-----
ug60:~/ece1756/lab2_cpu_gpu% █
```

Figure 26: CPU runtimes with O2 optimization

```

ug60:~/ece1756/lab2_cpu_gpu% rm conv_cpu
/bin/rm: remove regular file 'conv_cpu'? y
ug60:~/ece1756/lab2_cpu_gpu% g++ -fopenmp -O3 conv_cpu.cpp utils.cpp -o conv_cpu
ug60:~/ece1756/lab2_cpu_gpu% ./run_cpu.sh 1 4 16 64
Runtime of simple convolution for 1 filter(s) averaged accross 20 runs = 0.679928 ms
Runtime of hand vectorized convolution for 1 filter(s) averaged accross 20 runs = 1.25945 ms
-----
Runtime of simple convolution for 4 filter(s) averaged accross 20 runs = 2.72494 ms
Runtime of hand vectorized convolution for 4 filter(s) averaged accross 20 runs = 5.03768 ms
-----
Runtime of simple convolution for 16 filter(s) averaged accross 20 runs = 10.8303 ms
Runtime of hand vectorized convolution for 16 filter(s) averaged accross 20 runs = 23.019 ms
-----
Runtime of simple convolution for 64 filter(s) averaged accross 20 runs = 43.27 ms
Runtime of hand vectorized convolution for 64 filter(s) averaged accross 20 runs = 83.5155 ms
-----
ug60:~/ece1756/lab2_cpu_gpu%

```

Figure 27: CPU runtimes with O3 optimization

```

ug50:~/ece1756/lab2_cpu_gpu% g++ -fopenmp -O3 conv_cpu.cpp utils.cpp -o conv_cpu
ug60:~/ece1756/lab2_cpu_gpu% ./run_cpu.sh 1 4 16 64
Runtime of simple convolution for 1 filter(s) averaged accross 20 runs = 0.728284 ms
Runtime of hand vectorized convolution for 1 filter(s) averaged accross 20 runs = 1.19086 ms
-----
Runtime of simple convolution for 4 filter(s) averaged accross 20 runs = 0.823869 ms
Runtime of hand vectorized convolution for 4 filter(s) averaged accross 20 runs = 1.71489 ms
-----
Runtime of simple convolution for 16 filter(s) averaged accross 20 runs = 2.93434 ms
Runtime of hand vectorized convolution for 16 filter(s) averaged accross 20 runs = 7.82391 ms
-----
Runtime of simple convolution for 64 filter(s) averaged accross 20 runs = 11.8149 ms
Runtime of hand vectorized convolution for 64 filter(s) averaged accross 20 runs = 22.2321 ms
-----
ug60:~/ece1756/lab2_cpu_gpu% █

```

Figure 28: CPU runtimes with O3 optimization, multithreaded

```
ug60:~/ece1756/lab2_cpu_gpu% ./run_gpu.sh 1 4 16 64
Runtime for 1 filter(s) averaged accross 20 runs = 0.0297056 ms
-----
Runtime for 4 filter(s) averaged accross 20 runs = 0.0994058 ms
-----
Runtime for 16 filter(s) averaged accross 20 runs = 0.372853 ms
-----
Runtime for 64 filter(s) averaged accross 20 runs = 1.47965 ms
-----
ug60:~/ece1756/lab2_cpu_gpu%
```

Figure 29: GPU runtimes

## Appendix D: Evaluation Calculations

### FPGA related

$$Energy\ Efficiency = 1.2510 \times 10^3 / 11.98 \approx 104.42\ GOPS/W$$

$$Area\ Efficiency = 1.2510 \times 10^3 / 350 \approx 3.5743\ GOPS/mm^2$$

### CPU related

```

for(int filter_id = 0; filter_id < num_filters; filter_id++){ // For each filter
    for(int image_y = 0; image_y < image_height; image_y++){ // For each y location of the output image
        for(int image_x = 0; image_x < image_width; image_x++){ // For each x location of the output image
            // Start a new accumulation result
            accum = 0;

            // Compute the dot-product of the filter and the input image
            for(int filter_y = 0; filter_y < FILTER_SIZE; filter_y++){
                for(int filter_x = 0; filter_x < FILTER_SIZE; filter_x++){
                    xx = image_x + filter_x;
                    yy = image_y + filter_y;
                    accum += input_image[(yy * padded_image_width) + xx] * filter[(filter_id * FILTER_SIZE * FILTER_SIZE) +
                        (filter_y * FILTER_SIZE) + filter_x];
                }
            }

            // Store the convolution result in the corresponding pixel location
            output_image[(filter_id * image_width * image_height) + (image_y * image_width) + image_x] = accum;
        }
    }
}

```

Figure 30: CPU basic code

Based on Fig. 30 above, we can estimate the number of operations that are performed for the convolution in the basic implementation on the CPU. Number of operations within the most inner loop = 11. FILTER\_SIZE = 3 image\_height = image\_width = 512 + 2 num\_filters = 64 Margin of error = 10%

$$\begin{aligned} \text{Total Operations} &= (((11 * 3 * 3) + 5) * 514 * 514 * 64) * 1.10) / 11.8149 \times 10^{-3} \\ &\approx 163.72 \text{ GOPS} \end{aligned}$$

*Energy Efficiency* = 163.72/84 ≈ 1.9490 GOPS/W

$$Area\ Efficiency = 163.72/177 \approx 0.92497\ GOPS/mm^2$$

## GPU related

Based on Fig. 31 above, we can estimate the number of operations that are performed for the convolution in the implementation on the GPU. Number of operations within loop = 15  
 $\text{TILE\_SIZE} = 32$   $\text{TILE\_LOAD\_SIZE} = 32 + 3 - 1 = 34$

```

// Load tiles from Global memory to Shared memory for faster access
for (int itr = 0; itr <= (TILE_LOAD_SIZE * TILE_LOAD_SIZE) / (TILE_SIZE * TILE_SIZE); itr++) {
    // Calculate destination x and y indecies
    thread_id = (threadIdx.y * TILE_SIZE) + threadIdx.x + (itr * TILE_SIZE * TILE_SIZE);
    tile_location_y = thread_id / TILE_LOAD_SIZE;
    tile_location_x = thread_id % TILE_LOAD_SIZE;

    // Calculate source pixel index
    image_location_y = blockIdx.y * TILE_SIZE + tile_location_y - FILTER_RADIUS;
    image_location_x = blockIdx.x * TILE_SIZE + tile_location_x - FILTER_RADIUS;
    pixel_id = (image_location_y * image_width) + image_location_x;

    // Load pixels
    if (tile_location_y < TILE_LOAD_SIZE && tile_location_x < TILE_LOAD_SIZE) {
        if (image_location_y >= 0 && image_location_y < image_height
            && image_location_x >= 0 && image_location_x < image_width) {
            image_tile[tile_location_y][tile_location_x] = input_image[pixel_id];
        } else {
            image_tile[tile_location_y][tile_location_x] = 0;
        }
    }
}
__syncthreads();

```

Figure 31: GPU code

```

// Perform the 2D convolution
int accum = 0;
int y, x, z;
z = blockIdx.z;
for (y = 0; y < FILTER_SIZE; y++) {
    for (x = 0; x < FILTER_SIZE; x++) {
        accum += image_tile[threadIdx.y + y][threadIdx.x + x] * device_filter[(z * FILTER_SIZE * FILTER_SIZE) + (y * FILTER_SIZE) + x];
    }
}

// Write the output
y = blockIdx.y * TILE_SIZE + threadIdx.y;
x = blockIdx.x * TILE_SIZE + threadIdx.x;
if (y < image_height && x < image_width) {
    output_image[(z * image_width * image_height) + (y * image_width) + x] = accum;
}
__syncthreads();

```

Figure 32: GPU code

Number of operations within loop = 9 FILTER\_SIZE = 3 Above repeated for each pixel.  
 Number of operations outside loop = 9 Margin of error = 10%

$$\begin{aligned}
 \text{Total Operations} &= (15 * ((34 * 34) / (32 * 32)) + (9 * 3 * 3) * 512 * 512 + 9) * 1.10 / 1.47965 \times 10^{-3} \\
 &\approx 15.786 \text{ GOPS}
 \end{aligned}$$

$$\text{Energy Efficiency} = 15.786 / 165 \approx 0.095673 \text{ GOPS/W}$$

$$\text{Area Efficiency} = 15.786/396 \approx 0.39864 \text{ GOPS/mm}^2$$

**GPU Estimation**

$$\text{Throughput} = 15.786 * (28/20) \approx 22.1 \text{ GOPS}$$

$$\text{Power} = 165 * 1.5 \approx 247.5 \text{ W}$$

$$\text{Energy Efficiency} = 22.1/247.5 \approx 0.089293 \text{ GOPS/W}$$

$$\text{Area Efficiency} = 22.1/396 \approx 0.055808 \text{ GOPS/mm}^2$$

**Appendix E: .sv file for FPGA Implementation**

```

// This module implements 2D convolution between a 3x3 filter and a
// → 512-pixel-wide image of any height.
// It is assumed that the input image is padded with zeros such that the
// → input and output images have
// the same size. The filter coefficients are symmetric in the x-direction
// → (i.e. f[0][0] = f[0][2],
// f[1][0] = f[1][2], f[2][0] = f[2][2] for any filter f) and their values
// → are limited to integers
// (but can still be positive or negative). The input image is grayscale
// → with 8-bit pixel values ranging
// from 0 (black) to 255 (white).
module lab2 (
    input  clk,                      // Operating clock
    input  reset,                     // Active-high reset signal
    // → (reset when set to 1)
    input  [71:0] i_f,                // Nine 8-bit signed convolution
    // → filter coefficients in row-major format (i.e. i_f[7:0] is
    // → f[0][0], i_f[15:8] is f[0][1], etc.)
    input  i_valid,                  // Set to 1 if input pixel
    // → is valid
    input  i_ready,                  // Set to 1 if consumer
    // → block is ready to receive a new pixel
    input  [7:0] i_x,                // Input pixel value (8-bit
    // → unsigned value between 0 and 255)
    output o_valid,                 // Set to 1 if output pixel
    // → is valid
    output o_ready,                 // Set to 1 if this block is
    // → ready to receive a new pixel
    output [7:0] o_y,                // Output pixel value (8-bit
    // → unsigned value between 0 and 255)
);

localparam FILTER_SIZE = 3;          // Convolution filter dimension (i.e.
// → 3x3)
localparam PIXEL_DATAW = 8;         // Bit width of image pixels and filter
// → coefficients (i.e. 8 bits)

// The following code is intended to show you an example of how to use
// → parameters and

```

```

// for loops in SystemVerilog. It also arranges the input filter
// coefficients for you
// into a nicely-arranged and easy-to-use 2D array of registers. However,
// you can ignore
// this code and not use it if you wish to.

logic signed [PIXEL_DATAW-1:0] r_f [FILTER_SIZE-1:0][FILTER_SIZE-1:0]; // 2D
// array of registers for filter coefficients
integer col, row; // variables to use in the for loop
always_ff @ (posedge clk) begin
    // If reset signal is high, set all the filter coefficient
    // registers to zeros
    // We're using a synchronous reset, which is recommended style for
    // recent FPGA architectures
    if(reset)begin
        for(row = 0; row < FILTER_SIZE; row = row + 1) begin
            for(col = 0; col < FILTER_SIZE; col = col + 1) begin
                r_f[row][col] <= 0;
            end
        end
    // Otherwise, register the input filter coefficients into the 2D
    // array signal
    end else begin
        for(row = 0; row < FILTER_SIZE; row = row + 1) begin
            for(col = 0; col < FILTER_SIZE; col = col + 1) begin
                // Rearrange the 72-bit input into a 3x3
                // array of 8-bit filter coefficients.
                // signal[a +: b] is equivalent to
                // signal[a+b-1 : a]. You can try to plug
                // in
                // values for col and row from 0 to 2, to
                // understand how it operates.
                // For example at row=0 and col=0:
                //      r_f[0][0] = i_f[0+:8] = i_f[7:0]
                //      at row=0 and col=1:
                //      r_f[0][1] = i_f[8+:8] = i_f[15:8]
                r_f[row][col] <= i_f[(row * FILTER_SIZE *
                // PIXEL_DATAW)+(col * PIXEL_DATAW) +:
                // PIXEL_DATAW];
            end
        end
    end
end

```

```
end
```

```
////////// Start of your code
↪ //////////
```

```
localparam PIC_WIDTH = 512+2; // How wide is our picture
```

```
logic [PIXEL_DATAW-1:0] SIPO
↪ [PIC_WIDTH*2+FILTER_SIZE]; // SIPO, shift reg
↪ for data input
```

```

logic data_valid
    ↳ [PIC_WIDTH*2+FILTER_SIZE];
    ↳ Data valid propagate
    //////////////////////////////////////////////////////////////////

// Intermediate outputs
logic signed [31:0] mult_out [9];                                // Output
    ↳ of multiplier
logic
    ↳ data_valid_mult;
    ↳ Valid bit for if multiplier output is valid
    //////////////////////////////////////////////////////////////////

logic signed [31:0] addr_stage_1_out [4];                         // Output of adders
    ↳ in stage 1
logic signed [31:0] addr_stage_2_out [2];                         // Output of adders
    ↳ in stage 2
logic signed [31:0] addr_stage_3_out;                                // Output of
    ↳ adder in stage 3
logic signed [31:0] addr_stage_4_out;                                // Output of
    ↳ adder in stage 4
    //////////////////////////////////////////////////////////////////

// Pipeline registers
logic signed [31:0] mult_out_pipe [9];
logic data_valid_mult_pipe;

logic signed [31:0] addr_stage_1_out_pipe [5];
logic data_valid_addr_stage_1_pipe;
logic signed [31:0] addr_stage_2_out_pipe [3];
logic data_valid_addr_stage_2_pipe;
logic signed [31:0] addr_stage_3_out_pipe [2];
logic data_valid_addr_stage_3_pipe;
logic unsigned [PIXEL_DATAW-1:0] addr_stage_4_out_pipe;
    
```

```

logic data_valid_addr_stage_4_pipe; // Final, same level (in terms of
→ pipelining) as o_valid

// Counter to ensure correctness, no computation is to happen for the
→ first 2 valid inputs
// for each new row of pixels
logic [15:0] counter;

// Instantiating multipliers. All products are computed in one and the
→ same clock cycle
mult8x8 convo_mult_0 (.i_dataaa(SIPO[2+PIC_WIDTH*2]),
→ .i_datab(r_f[0][0]), .o_res(mult_out[0])); // Row 0, Col
→ 0
mult8x8 convo_mult_1 (.i_dataaa(SIPO[1+PIC_WIDTH*2]),
→ .i_datab(r_f[0][1]), .o_res(mult_out[1])); // Row 0, Col
→ 1
mult8x8 convo_mult_2 (.i_dataaa(SIPO[0+PIC_WIDTH*2]),
→ .i_datab(r_f[0][2]), .o_res(mult_out[2])); // Row 0, Col
→ 2
mult8x8 convo_mult_3 (.i_dataaa(SIPO[2+PIC_WIDTH*1]),
→ .i_datab(r_f[1][0]), .o_res(mult_out[3])); // Row 1, Col
→ 0
mult8x8 convo_mult_4 (.i_dataaa(SIPO[1+PIC_WIDTH*1]),
→ .i_datab(r_f[1][1]), .o_res(mult_out[4])); // Row 1, Col
→ 1
mult8x8 convo_mult_5 (.i_dataaa(SIPO[0+PIC_WIDTH*1]),
→ .i_datab(r_f[1][2]), .o_res(mult_out[5])); // Row 1, Col
→ 2
mult8x8 convo_mult_6 (.i_dataaa(SIPO[2]),
→ .i_datab(r_f[2][0]),
→ .o_res(mult_out[6])); // Row 2, Col 0
mult8x8 convo_mult_7 (.i_dataaa(SIPO[1]),
→ .i_datab(r_f[2][1]),
→ .o_res(mult_out[7])); // Row 2, Col 1
mult8x8 convo_mult_8 (.i_dataaa(SIPO[0]),
→ .i_datab(r_f[2][2]),
→ .o_res(mult_out[8])); // Row 2, Col 2

// Instantiating adders. Each grouping of adders in one clock cycle each
addr32p32 convo_addr_stage_1_0 (.i_dataaa(mult_out_pipe[0]),
→ .i_datab(mult_out_pipe[1]), .o_res(addr_stage_1_out[0]));

```

```

addr32p32 convo_addr_stage_1_1 (.i_dataaa(mult_out_pipe[2]),
→ .i_datab(mult_out_pipe[3]), .o_res(addr_stage_1_out[1]));
addr32p32 convo_addr_stage_1_2 (.i_dataaa(mult_out_pipe[4]),
→ .i_datab(mult_out_pipe[5]), .o_res(addr_stage_1_out[2]));
addr32p32 convo_addr_stage_1_3 (.i_dataaa(mult_out_pipe[6]),
→ .i_datab(mult_out_pipe[7]), .o_res(addr_stage_1_out[3]));

addr32p32 convo_addr_stage_2_0 (.i_dataaa(addr_stage_1_out_pipe[0]),
→ .i_datab(addr_stage_1_out_pipe[1]), .o_res(addr_stage_2_out[0]));
addr32p32 convo_addr_stage_2_1 (.i_dataaa(addr_stage_1_out_pipe[2]),
→ .i_datab(addr_stage_1_out_pipe[3]), .o_res(addr_stage_2_out[1]));

addr32p32 convo_addr_stage_3_0 (.i_dataaa(addr_stage_2_out_pipe[0]),
→ .i_datab(addr_stage_2_out_pipe[1]), .o_res(addr_stage_3_out));

addr32p32 convo_addr_stage_4_0 (.i_dataaa(addr_stage_3_out_pipe[0]),
→ .i_datab(addr_stage_3_out_pipe[1]), .o_res(addr_stage_4_out));

// Valid for computed values is only true if all inputs are valid
// otherwise, similar to lab 1, these are just propagated along with
→ computed values
always_comb begin
    // Multiplier output only valid if all 9 data bytes are valid
    data_valid_mult = data_valid[2+PIC_WIDTH*2] &
        → data_valid[1+PIC_WIDTH*2] & data_valid[0+PIC_WIDTH*2] &
            data_valid[2+PIC_WIDTH*1]
            → &
            → data_valid[1+PIC_WIDTH*1]
            → &
            → data_valid[0+PIC_WIDTH*1]
            → &
            data_valid[2]
            → & data_valid[1]
            → & data_valid[0]
            → &
            (counter >=
                → FILTER_SIZE); //
                → Need at least 3
                → element in the
                → row before
                → computing
end

```

```

logic enable;
always_comb begin
    enable = i_ready & i_valid;
end

int i;
always_ff @ (posedge clk) begin

    if(reset)begin
        // Reset counter
        counter <= 16'b0;

        // Reset SIPO and valid propagation
        for(i = 0; i < $size(SIPO); i = i + 1) begin
            SIPO[i] <= 8'b0;
            data_valid[i] <= 1'b0;
        end

    end else begin
        if(enable)begin
            // Counter
            // If counter at 514 means current row complete,
            // → reset to 1 since it's the first entry of the
            // → next row
            if (counter == PIC_WIDTH) begin
                counter <= 16'b1;
            end else begin
                counter <= counter + 1'b1;
            end

            SIPO[0] <= i_x;                                // If
            // → input is valid, load into SIPO[0]
            data_valid[0] <= i_valid;                      //
            // → Similarly with valids
            for(i = 1; i < $size(SIPO); i = i + 1) begin
                SIPO[i] <=
                // → SIPO[i-1];
                // → Propagate signal in shift reg

```

```

        data_valid[i] <= data_valid[i-1];           //  

        → Propagate valid signal  

    end

    // Pipeline reg for multiplier out, adder stage 1
    → in
    mult_out_pipe <= mult_out;
    data_valid_mult_pipe <= data_valid_mult;

    data_valid_addr_stage_1_pipe <=
        → data_valid_mult_pipe;
    data_valid_addr_stage_2_pipe <=
        → data_valid_addr_stage_1_pipe;
    data_valid_addr_stage_3_pipe <=
        → data_valid_addr_stage_2_pipe;
    data_valid_addr_stage_4_pipe <=
        → data_valid_addr_stage_3_pipe;

    // Pipeline reg for adder stage 1 out, adder stage
    → 2 in
    addr_stage_1_out_pipe[0] <= addr_stage_1_out[0];
    addr_stage_1_out_pipe[1] <= addr_stage_1_out[1];
    addr_stage_1_out_pipe[2] <= addr_stage_1_out[2];
    addr_stage_1_out_pipe[3] <= addr_stage_1_out[3];
    addr_stage_1_out_pipe[4] <= mult_out_pipe[8];

    // Pipeline reg for adder stage 2 out, adder stage
    → 3 in
    addr_stage_2_out_pipe[0] <= addr_stage_2_out[0];
    addr_stage_2_out_pipe[1] <= addr_stage_2_out[1];
    addr_stage_2_out_pipe[2] <=
        → addr_stage_1_out_pipe[4];

    // Pipeline reg for adder stage 3 out, adder stage
    → 4 in
    addr_stage_3_out_pipe[0] <= addr_stage_3_out;
    addr_stage_3_out_pipe[1] <=
        → addr_stage_2_out_pipe[2];

    // Pipeline reg for adder stage 4 out, final
    → result, clipping performed here

```

```

        if (addr_stage_4_out > 32'sd255) begin
            addr_stage_4_out_pipe <= 8'd255;
        end else if (addr_stage_4_out < 32'sd0) begin
            addr_stage_4_out_pipe <= 8'd0;
        end else begin
            addr_stage_4_out_pipe <=
                → addr_stage_4_out[7:0];
        end
    end
end

assign o_y = addr_stage_4_out_pipe;
assign o_ready = i_ready;
assign o_valid = data_valid_addr_stage_4_pipe & i_ready;

//////////////////////////////////////////////////////////////// End of your code
→ //////////////////////////////////////////////////////////////////

endmodule

//****************************************************************************

// Multiplier module for 8x8 multiplication
// Taken from lab 1 with modifications for signed multiplication
module mult8x8 ( /* synthesis multstyle = "dsp" */
    input unsigned [7:0] i_dataaa,
    input signed [7:0] i_datab,
    output signed [31:0] o_res
);

logic signed [31:0] result;

always_comb begin
    result = i_datab * $signed({1'b0,i_dataaa});
end

assign o_res = result;

endmodule

```

```
*****  
*****  
// Adder module for 32b+32b addition operations  
// Taken from lab 1 with modifications for signed addition  
module addr32p32 (  
    input signed [31:0] i_dataaa,  
    input signed [31:0] i_datab,  
    output signed [31:0] o_res  
);  
  
assign o_res = i_dataaa + i_datab;  
  
endmodule  
*****
```