# Assignment 2

# 1 FPGA Implementation
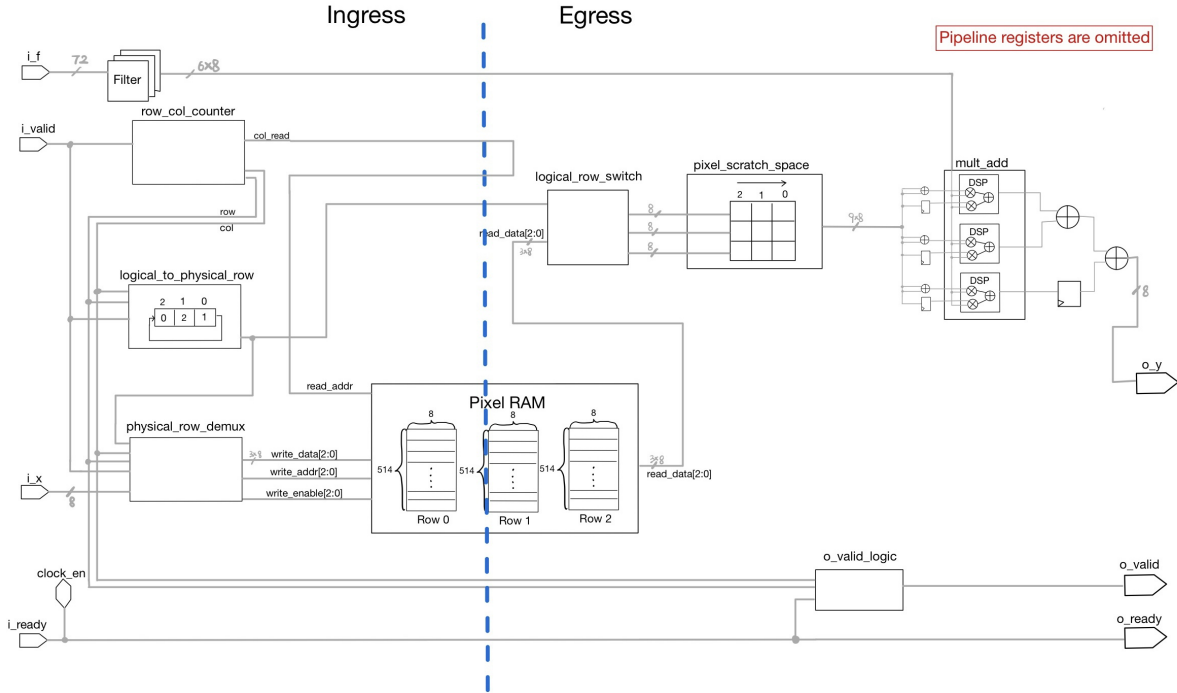
## 1.1 Design Description



Figure 1: A simplified block diagram of the design, which runs at Restricted Fmax 549.45 MHz, Fmax 687.76 MHz, with a throughput of 1 output pixel per cycle

The design can be logically separated into two main stages, ingress and egress, with respect to the RAM that stores the input image pixels (pixel RAM).

Ingress path is responsible for keeping track of the column and row value of the input pixel, which is used as both the write address to the pixel RAM, and for determining the current convolution window (i.e., pointing to the bottom right corner of the convolution window). Another important responsibility for ingress path is to control the row shifting logic for the pixel RAM.

Pixel RAM is a 2D array of pixels, with a dimension of $3 \times (512 + 2) \times 8$. Only three rows are stored because that is the size of the convolution filter. Write is always done at the bottom row, and the row shifts upward when column counter reaches the bound. To save area, this shifting register file is implemented in BRAM, where the shifting is realized by shifting the mapping from logical row index to physical row index circularly, so that after the shifting, new incoming pixels are essentially overwriting the oldest row. Pixel RAM is implemented with three independent BRAM (M20K) blocks, each representing a row, with a shape of $(512 + 2) \times 8$. For ease of testing and maximal compatibility, the BRAM block
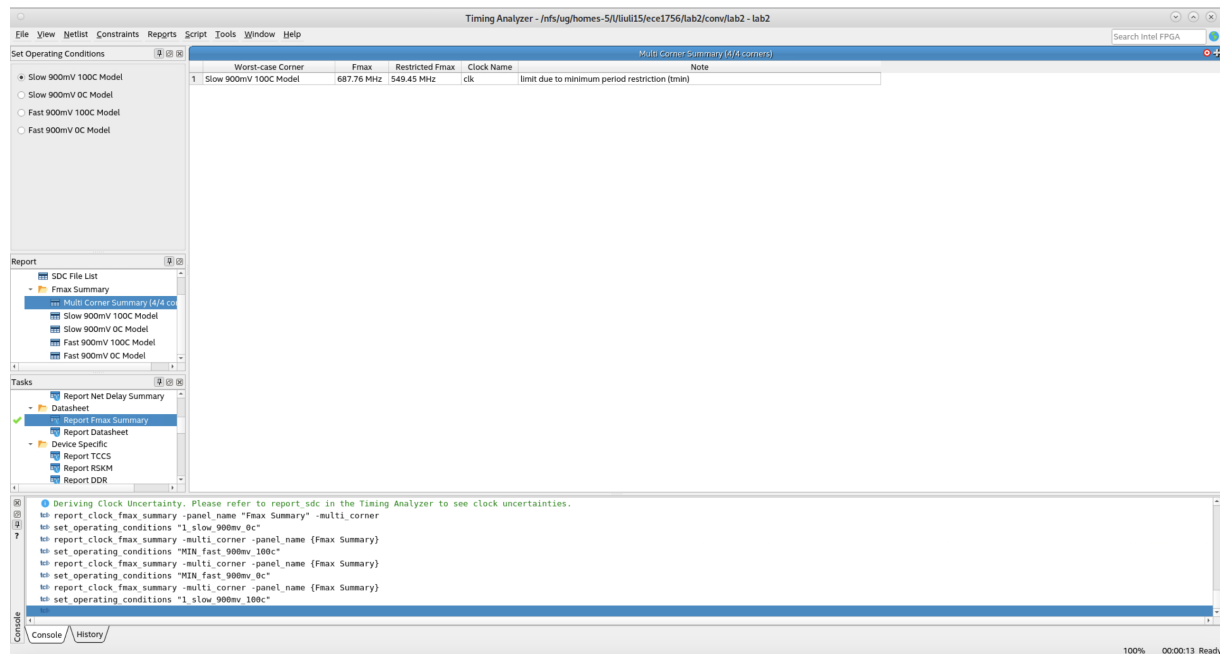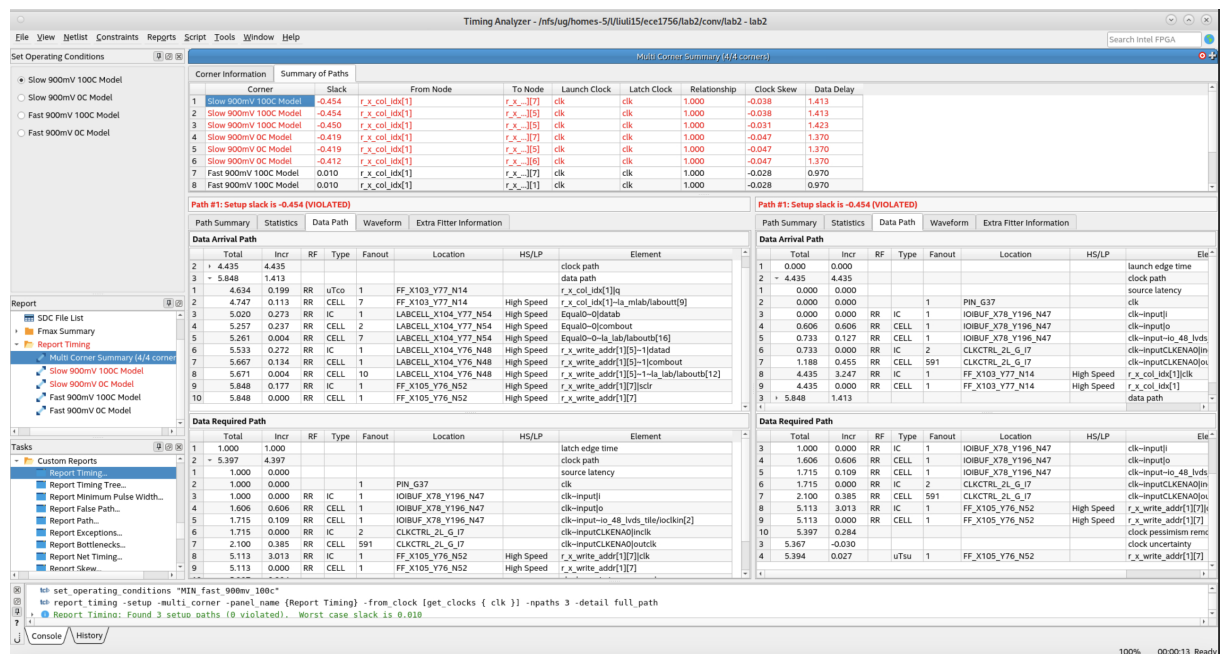
# Assignment 2

Figure 2: Fmax Summary



Figure 3: Critical Paths

Figure 4: Fitter Summary

is not instantiated from IP, but rather using an HDL template that can be inferred into BRAM (M20K) in simple dual-port mode by Quartus Synthesis, with read and write width of 8. With three independent BRAM (M20K) blocks, the pixel RAM is capable of doing a simultaneous read and write on each row in any clock cycle.

Egress path is responsible for performing the actual convolution computation for the convolution window that is pointed by the row and column counter aforementioned. For every window, it reads a single pixel for each row from the pixel RAM, and stores it into a scratch space that is $3 \times 3$ in size. The scratch space would shift horizontally to represent the fact that the convolution windows is also moving horizontally. So at any convolution window, only the right-most column is freshly read from the pixel RAM, while the other columns are cached from previous convolution windows. This effectively reduces the number of simultaneous read demand from 9 in a naive implementation, down to 3. To better exploit the fact that the filter is symmetric in x-direction, pixels at column 0 and column 2 are added before multiplying the filter, which effectively saves 3 multipliers in total. The total of 6 multipliers are divided into 3 sum of 2 $18 \times 18$ multipliers groups, where each group can nicely fit into a DSP block on chip. Again, an HDL template that can be inferred to a fully registered DSP is used rather than instantiating from IP. Afterwards, a 2-level reduction (sum) tree is used to accumulate from the 3 DSP outputs. In total, 8 adders and 6 multipliers are used saving 3 multipliers from a naive implementation of 8 adders and 9 multipliers. A lot of

pipeline stages are inserted to this arithmetic chain to boost Fmax. The entire arithmetic logic runs continuously, but the output is valid only if the convolution window pointer is at valid location and has not been computed before.

## 1.2 FPGA Implementation Results

|  | Result |
|---|---|
| **ALM Utilization** | 166 out of 427200 |
| **DSP Utilization** | 3 out of 1518 |
| **BRAM (M20K) Utilization** | 3 out of 2713 |
| **Maximum Operating Frequency (MHz)** | Restricted Fmax 549.45 MHz, Fmax 687.76 MHz |
| **Cycles for Test 7a (Hinton)** | 264222 |
| **Dynamic Power for one module @ maximum frequency (W)** | 0.048 W |
| **Throughput of one module (GOPS)** | 9.267 GOPS |
| **Throughput of full device (GOPS)** | 4689.198 GOPS |
| **Total Power for full device (W)** | 10.188 W |

**Calculations**

**Dynamic Power for one module @ maximum frequency (W):** $(DSP + M20K + LUT + FF + ClockNetwork) \div 50MHz \times RestrictedFmax = (0.08mW + 0.45mW + 0.20mW + 0.79mW + 2.82mW) \div 50MHz \times 549.45MHz$

**Throughput of one module (GOPS):** $512 \times 512 \times (9 + 8) \div (264222 \div 549.45MHz)$

**Throughput of full device (GOPS):** $506 \times 9.267GOPS$

**Total Power for full device (W):** $(DSP + M20K + LUT + FF) \div 50MHz \times RestrictedFmax \times 506 + ClockNetwork \div 50MHz \times RestrictedFmax + DeviceStatic(1704.56mW) + IO(0.19mW)$

## 1.3 Simulation Waveforms and Testbench Output

Refer to Figure 5 and Figure 6.

# 2 Efficiency Comparison

## 2.1 CPU Implementation

CPU convolution implementations are quite straight forward. Both single-threaded and multithreaded, and both basic and hand-vectorized versions loop for filter_id, output pixel

Figure 5: ModelSim waveform and simulation log for Test7a as expected



Figure 6: ModelSim waveform and simulation log for Test7b as expected

row and column respectively. Then a dot product is computed for the $3 \times 3$ convolution window involved. The basic implementation traverses through all 9 elements and accumulates the products, whereas the hand-vectorized version computes the dot product of the entire row in one shot (load, dot product and store) and then accumulate the the results of three rows together. Multi-threading is done using OpenMP directives, and is applied in the outermost filter_id dimension.

## 2.2   GPU Implementation

GPU convolution kernel executes in two stages. Firstly, each thread loads a few pixels of the current tile into shared memory. Then after all threads within the current block are done, each thread starts to calculate the dot product for the output pixel and filter_id assigned to it. The filters are stored in the GPU constant memory.

## 2.3   Runtime Results of CPU and GPU Implementations

|                                          | Runtime (ms) | | | |
|------------------------------------------|--------|--------|--------|--------|
| **No. of filters**                       | **1**  | **4**  | **16** | **64** |
| **GPU**                                  | 0.0151 | 0.0401 | 0.1374 | 0.5060 |
| **CPU (basic - no opt - 1 thread)**      | 6.1022 | 24.306 | 97.235 | 389.02 |
| **CPU (vectorized - no opt - 1 thread)** | 3.1988 | 13.027 | 51.896 | 207.30 |
| **CPU (basic - O2 - 1 thread)**          | 1.2414 | 4.9722 | 19.807 | 79.261 |
| **CPU (vectorized - O2 - 1 thread)**     | 0.8441 | 3.3751 | 13.493 | 54.013 |
| **CPU (basic - O3 - 1 thread)**          | 0.5334 | 2.1221 | 8.4699 | 33.843 |
| **CPU (vectorized - O3 - 1 thread)**     | 0.8823 | 3.4746 | 13.898 | 55.689 |
| **CPU (basic - O3 - 4 threads)**         | 0.5770 | 0.6108 | 1.2119 | 4.8529 |
| **CPU (vectorized - O3 - 4 threads)**    | 0.9082 | 1.0459 | 1.8682 | 7.7334 |

When there is no compiler optimization, the vectorized version has a 2x performance than the basic version. When the compiler uses O2 optimization, the gap becomes smaller but still in favor of the vectorized version. Both versions have significant improvements comparing to the no opt version, but basic version gains more out of it suggesting that O2 optimizes both the dot product and iteration overheads for the basic version whereas the vectorized version only gains improvement from iteration overheads. This is even more obvious under O3 optimization, where the basic version has much better performance than the vectorized version and the performance of the vectorized version is the same as under O2 optimization.

After turning on multi-threading, comparing to the single-threaded result (with O3 optimization), the relative performance gap between the basic version and vectorized version stays the same. It is worth pointing out the speed up comparing to single-threaded version is dependent on the number of filters, there is no speed up when there is only 1 filter, and

about 3.5x speed up for 4 filters. The speed up caps at around 8x for both 16 and 64 filters primarily due to the fact that the CPU under test has 8 cores.

## 2.4 Comparison of Results Between the 3 Compute Platforms

| | Throughput (GOPS) | Power (W) | Energy Efficiency (GOPS/W) | Area Efficiency (GOPS/mm$^2$) |
|---|---|---|---|---|
| **FPGA (20 nm)** | 4689.2 GOPS | 10.188 W | 460.27 GOPS/W | 11.723 GOPS/mm$^2$ |
| **CPU (14 nm)** | 58.771 GOPS | 65 W | 0.904 GOPS/W | 0.213 GOPS/mm$^2$ |
| **GPU (8 nm)** | 563.63 GOPS | 220 W | 2.562 GOPS/W | 1.434 GOPS/mm$^2$ |
| **FPGA (scaled to 8nm)** | 7502.7 GOPS | 6.875 W | 1091.3 GOPS/W | 75.027 GOPS/mm$^2$ |
| **CPU (scaled to 8nm)** | 73.464 GOPS | 56.875 W | 1.292 GOPS/W | 0.532 GOPS/mm$^2$ |

**Calculations**

**Scaled Throughput (GOPS):** $OriginalThroughput \times ClockSpeedScaling$

**Scaled Power:** $OriginalPower \times ClockSpeedScaling \times Power@sameclockrateScaling$

**Scaled Area:** $OriginalArea \times AreaScaling$

When scaling to the same process technology, throughput wise, FPGA is one order ahead of GPU due primarily to that FPGA is programmed and dedicated to compute the convolution. One advantage with the FPGA is that it outputs a pixel at each clock cycle, where as GPU takes several. The GPU version is also a magnitude ahead of CPU, primarily due to the massive parallelism that convolution operation exhibits and that GPU can better exploit it.

Energy efficiency wise, FPGA is almost 1000x ahead of CPU and 400x ahead of GPU. This time, CPU and GPU are within the same magnitude due to them both being a general purpose computing platform, whereas the programmed FPGA is more application specific thus avoiding a lot of the overhead cost associated with generality. GPU is still ahead of CPU because of its wider SIMD, reducing overhead cost of instruction fetching and decoding.

Area efficiency wise, FPGA is almost 150x ahead of CPU and 50x ahead of GPU. The reason is again related to the overhead cost with generality, and overhead cost of instruction fetching and decoding between CPU and GPU.

## 3   Appendix

### 3.1   HDL

```verilog
1  // This module implements 2D covolution between a 3x3 filter and a
   //   512-pixel-wide image of any height.
2  // It is assumed that the input image is padded with zeros such that the
   //   input and output images have
3  // the same size. The filter coefficients are symmetric in the x-direction
   //   (i.e. f[0][0] = f[0][2],
4  // f[1][0] = f[1][2], f[2][0] = f[2][2] for any filter f) and their values
   //   are limited to integers
5  // (but can still be positive of negative). The input image is grayscale
   //   with 8-bit pixel values ranging
6  // from 0 (black) to 255 (white).
7  module lab2 (
8         input  clk,                      // Operating clock
9         input  reset,                    // Active-high reset signal
   //   (reset when set to 1)
10        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.)
11        input  i_valid,                  // Set to 1 if input pixel
   //   is valid
12        input  i_ready,                  // Set to 1 if consumer
   //   block is ready to receive a new pixel
13        input  [7:0] i_x,                // Input pixel value (8-bit
   //   unsigned value between 0 and 255)
14        output o_valid,                  // Set to 1 if output pixel
   //   is valid
15        output o_ready,                  // Set to 1 if this block is
   //   ready to receive a new pixel
16        output [7:0] o_y                 // Output pixel value (8-bit
   //   unsigned value between 0 and 255)
17 );
18
19 localparam FILTER_SIZE = 3;       // Convolution filter dimension (i.e.
   //   3x3)
20 localparam PIXEL_DATAW = 8;       // Bit width of image pixels and filter
   //   coefficients (i.e. 8 bits)
21
22 integer col, row, i; // variables to use in the for loop
```

```
23
24   // The following code is intended to show you an example of how to use
     ↪   paramaters and
25   // for loops in SytemVerilog. It also arrages the input filter
     ↪   coefficients for you
26   // into a nicely-arranged and easy-to-use 2D array of registers. However,
     ↪   you can ignore
27   // this code and not use it if you wish to.
28
29   logic signed [PIXEL_DATAW-1:0] r_f [FILTER_SIZE-1:0][FILTER_SIZE-1:0]; // 2D
     ↪   array of registers for filter coefficients
30   always_ff @ (posedge clk) begin
31           // If reset signal is high, set all the filter coefficient
     ↪   registers to zeros
32           // We're using a synchronous reset, which is recommended style for
     ↪   recent FPGA architectures
33           if(reset)begin
34                   for(row = 0; row < FILTER_SIZE; row = row + 1) begin
35                           for(col = 0; col < FILTER_SIZE; col = col + 1) begin
36                                   r_f[row][col] <= 0;
37                           end
38                   end
39           // Otherwise, register the input filter coefficients into the 2D
     ↪   array signal
40           end else begin
41                   for(row = 0; row < FILTER_SIZE; row = row + 1) begin
42                           for(col = 0; col < FILTER_SIZE; col = col + 1) begin
43                                   // Rearrange the 72-bit input into a 3x3
     ↪   array of 8-bit filter coefficients.
44                                   // signal[a +: b] is equivalent to
     ↪   signal[a+b-1 : a]. You can try to plug in
45                                   // values for col and row from 0 to 2, to
     ↪   understand how it operates.
46                                   // For example at row=0 and col=0:
     ↪   r_f[0][0] = i_f[0+:8] = i_f[7:0]
47                                   //              at row=0 and col=1:
     ↪   r_f[0][1] = i_f[8+:8] = i_f[15:8]
48                                   r_f[row][col] <= i_f[(row * FILTER_SIZE *
     ↪   PIXEL_DATAW)+(col * PIXEL_DATAW) +: PIXEL_DATAW];
49                           end
50                   end
51           end
```

```systemverilog
52  end
53
54  // Start of your code
55  logic enable;
56  assign enable = i_ready;
57
58  // *********************
59  // INGRESS
60  // *********************
61  // If pipelined, need to pipeline i_valid
62
63  // Logics for buffer of x
64  localparam IMAGE_WIDTH = 512;
65  localparam R_X_ROWS = FILTER_SIZE; // Always store 3 rows of i_x
66  localparam R_X_COL_WIDTH = IMAGE_WIDTH + 2;
67  localparam R_X_COL_ADDR_WIDTH = 10;
68
69
70  // Pixel RAM
71  // RAM input, need to be registered except for r_x_read_addr
72  logic [R_X_COL_ADDR_WIDTH-1:0] r_x_write_addr [R_X_ROWS-1:0]; // 0..511+2
    ↪   (10 bit)
73  logic r_x_write_enable [R_X_ROWS-1:0];
74  logic unsigned [PIXEL_DATAW-1:0] r_x_write_data [R_X_ROWS-1:0];
75  // Registered inside module, read a pixel from 3 rows
76  logic [R_X_COL_ADDR_WIDTH-1:0] r_x_read_addr;
77  // RAM output
78  logic unsigned [PIXEL_DATAW-1:0] r_x_read_data [R_X_ROWS-1:0]; // is
    ↪   registered inside module, 3 rows of 1 pixel
79  pixelram pixel_ram
80  (
81          .clk(clk),
82          .reset(reset),
83          .enable(enable),
84          // RAM input, unregistered inside module
85          .i_write_addr(r_x_write_addr),
86          .i_write_enable(r_x_write_enable),
87          .i_write_data(r_x_write_data),
88          .i_read_addr(r_x_read_addr),
89          // RAM output
90          .o_read_data(r_x_read_data)
91  );
```

```
92
93
94   // Registers
95   // 0: [] [] [] [] [] [] [] []                           512 + 2
96   // 1: [] [] [] [] [] [] [] []                           512 + 2
97   // 2: [] [] [] [] [] [] [] []                           512 + 2
98   logic unsigned [1:0] r_x_row_logical_idx; // Count from 0 to R_X_ROWS - 1
     ↪  (incl), logical order, not necessarily physical
99   logic unsigned [R_X_COL_ADDR_WIDTH-1:0] r_x_col_idx; // Count from 0 to
     ↪  R_X_COL_WIDTH (incl)
100  logic unsigned [R_X_COL_ADDR_WIDTH-1:0] r_x_col_idx_read_addr_adjusted; //
     ↪  Count from 0 to R_X_COL_WIDTH (incl)
101  // Count from 0 to R_X_ROWS - 1 (incl), physical order
102  logic unsigned [R_X_ROWS-1:0][1:0] r_x_row_logical_to_physical_index;
103
104  // INGRESS: Stage -1
105  always_ff @ (posedge clk) begin
106          if(reset) begin
107                  r_x_row_logical_idx <= 0;
108                  r_x_col_idx <= 0;
109                  r_x_col_idx_read_addr_adjusted <= 0;
110
111                  for(i = 0; i < R_X_ROWS; i = i + 1) begin
112                          r_x_row_logical_to_physical_index[i] <= i;
113                  end
114
115                  for(row = 0; row < R_X_ROWS; row = row + 1) begin
116                          r_x_write_addr[row] <= 0;
117                          r_x_write_data[row] <= 0;
118                          r_x_write_enable[row] <= 0;
119                  end
120          end else if (enable) begin
121                  // Do not write by default
122                  for(row = 0; row < R_X_ROWS; row = row + 1) begin
123                          r_x_write_addr[row] <= 0;
124                          r_x_write_data[row] <= 0;
125                          r_x_write_enable[row] <= 0;
126                  end
127                  r_x_col_idx_read_addr_adjusted <= 0;
128                   if(i_valid) begin
129                          if(r_x_col_idx == R_X_COL_WIDTH) begin
```

```
130                              // Load input pixel to a new row at the
     ↪  current logical idx 0 (R_X_COL_WIDTH implies 0),
131                              // which would be discarded, then
     ↪  reused/overwritten as the new logical idx 2 in the next cycle
132                              r_x_write_addr[r_x_row_logical_to_physical_index[0]]
     ↪  <= 0;
133                              r_x_write_data[r_x_row_logical_to_physical_index[0]]
     ↪  <= i_x;
134                              r_x_write_enable[r_x_row_logical_to_physical_index[0]]
     ↪  <= 1;
135
136                              // Do the row shifting logic at the first
     ↪  input of the new row,
137                              // rather than at the last input of the old
     ↪  row (will have conflict
138                              // in writing old row and shifting old row)
139
140                              // Instead of shifting the actual data,
     ↪  shift the mapping from logical index to physical index
141                              // Shift the mapping, upward
     ↪  (idx[0]->idx[2])
142                              r_x_row_logical_to_physical_index <=
     ↪  {r_x_row_logical_to_physical_index[0],
     ↪  r_x_row_logical_to_physical_index[R_X_ROWS-1:1]};
143
144                              // Reset r_x_col_idx if necessary,
     ↪  continuing at idx 1.
145                              // Skipping idx 0 because we are at idx 0
     ↪  currently
146                              r_x_col_idx <= 1;
147                              r_x_col_idx_read_addr_adjusted <= 0;
148
149                              // Increment r_x_row_logical_idx_ipipelined
     ↪  only when r_x_row_logical_idx_ipipelined is 0 or 1,
150                              // so that r_x_row_logical_idx_ipipelined
     ↪  will reach to 2 in steady state
151                              if(r_x_row_logical_idx < R_X_ROWS - 1) begin
152                                      r_x_row_logical_idx <=
     ↪  r_x_row_logical_idx + 1;
153                              end
154                      end else begin
155                              // Load data at logical idx 2
```

```
156                                              r_x_write_addr[r_x_row_logical_to_physical_index[R_X_ROW
    ↪  <= r_x_col_idx;
157                                              r_x_write_data[r_x_row_logical_to_physical_index[R_X_ROW
    ↪  <= i_x;
158                                              r_x_write_enable[r_x_row_logical_to_physical_index[R_X_R
    ↪  <= 1;
159
160                                              // Increment r_x_col_idx
161                                              r_x_col_idx <= r_x_col_idx + 1;
162
163                                              // Increment r_x_col_idx_read_addr_adjusted
164                                              r_x_col_idx_read_addr_adjusted <=
    ↪  r_x_col_idx;
165                                      end
166                              end
167                      end
168      end
169
170      // Pipeline registers for ingress
171      logic unsigned [1:0] r_x_row_logical_idx_ipipelined;
172      logic unsigned [R_X_COL_ADDR_WIDTH-1:0] r_x_col_idx_ipipelined;
173      logic unsigned [R_X_ROWS-1:0][1:0]
    ↪  r_x_row_logical_to_physical_index_ipipelined;
174      // INGRESS: Stage 0, to match registered RAM write
175      always_ff @ (posedge clk) begin
176              if (enable) begin
177                      r_x_row_logical_idx_ipipelined <= r_x_row_logical_idx;
178                      r_x_col_idx_ipipelined <= r_x_col_idx;
179                      r_x_row_logical_to_physical_index_ipipelined <=
    ↪  r_x_row_logical_to_physical_index;
180              end
181      end
182
183      // *********************
184      // EGRESS
185      // *********************
186
187      // Pipeline registers for egress
188      localparam NUM_EGRESS_STAGE = 9;
189      logic unsigned [NUM_EGRESS_STAGE-1:0] [R_X_COL_ADDR_WIDTH-1:0]
    ↪  r_x_col_idx_epipelined;
190      logic unsigned [NUM_EGRESS_STAGE-1:0] [1:0] r_x_row_logical_idx_epipelined;
```

```
191  logic unsigned [1:0] [R_X_ROWS-1:0][1:0]
  ↪  r_x_row_logical_to_physical_index_epipelined; // Not needed for full
  ↪  pipeline stage
192  always_ff @ (posedge clk) begin
193         if(enable) begin
194                r_x_col_idx_epipelined <=
  ↪  {r_x_col_idx_epipelined[NUM_EGRESS_STAGE-2:0], r_x_col_idx_ipipelined};
195                r_x_row_logical_idx_epipelined <=
  ↪  {r_x_row_logical_idx_epipelined[NUM_EGRESS_STAGE-2:0],
  ↪  r_x_row_logical_idx_ipipelined};
196                r_x_row_logical_to_physical_index_epipelined <=
  ↪  {r_x_row_logical_to_physical_index_epipelined[0:0],
  ↪  r_x_row_logical_to_physical_index_ipipelined};
197         end
198  end
199
200  // Logics for convolution core
201  // Computation
202
203  // EGRESS: Stage -1
204  always_comb begin
205         r_x_read_addr = r_x_col_idx_read_addr_adjusted;
206  end
207
208  // EGRESS: Stage 0, 1
209  // Signed x unsigned gets unsigned, which is not what we intend.
210  // So convert unsigned to signed by treating unsigned number as positive
  ↪   (by adding a 0 to msb)
211  logic signed [FILTER_SIZE-1:0] [PIXEL_DATAW:0] r_mult_i_pixel
  ↪   [R_X_ROWS-1:0];
212  logic [PIXEL_DATAW:0] r_x_read_data_reg [R_X_ROWS-1:0];
213  always_ff @ (posedge clk) begin
214         if(enable) begin
215                for(row=0; row<R_X_ROWS; row=row+1) begin
216                       r_x_read_data_reg[row] <= r_x_read_data[row];
217                       r_mult_i_pixel[row] <= {{1'b0,
  ↪  r_x_read_data_reg[r_x_row_logical_to_physical_index_epipelined[1][row]]},
  ↪  r_mult_i_pixel[row][FILTER_SIZE-1:1]};
218                end
219         end
220  end
221
```

```
222  // Multiplication
223  // EGRESS: Stage 2, 3, 4, 5
224  logic signed [FILTER_SIZE-1:0] [2*PIXEL_DATAW-1:0] sums_stage_0;
225  genvar gen_row;
226  generate
227          for(gen_row=0; gen_row<FILTER_SIZE; gen_row=gen_row+1) begin: mult
228                  mult8x8p8x8 m0(
229                          .clk(clk),
230                          .enable(enable),
231                          .i_filtera(r_f[gen_row][0]),
232                          .i_pixela0(r_mult_i_pixel[gen_row][0]),
233                          .i_pixela1(r_mult_i_pixel[gen_row][2]),
234                          .i_filterb(r_f[gen_row][1]),
235                          .i_pixelb(r_mult_i_pixel[gen_row][1]),
236                          .o_res(sums_stage_0[gen_row])
237                  );
238          end
239  endgenerate
240
241  // Reduction tree
242  // EGRESS: Stage 6
243  logic signed [FILTER_SIZE-1:0] [2*PIXEL_DATAW-1:0] sums_stage_0_reg;
244  always_ff @ (posedge clk) begin
245          if(enable) begin
246                  sums_stage_0_reg <= sums_stage_0;
247          end
248  end
249
250  logic signed [2*PIXEL_DATAW-1:0] sums_stage_1 [1:0];
251  always_comb begin
252          sums_stage_1[0] = sums_stage_0_reg[0] + sums_stage_0_reg[1];
253          sums_stage_1[1] = sums_stage_0_reg[2];
254  end
255
256  logic signed [2*PIXEL_DATAW-1:0] sums_stage_2;
257  always_comb begin
258          sums_stage_2 = sums_stage_1[0] + sums_stage_1[1];
259  end
260  // EGRESS: Stage 7
261  logic signed [2*PIXEL_DATAW-1:0] sums_stage_2_reg;
262  always_ff @ (posedge clk) begin
263          if(enable)begin
```

```
264                       sums_stage_2_reg <= sums_stage_2;
265             end
266    end
267
268    logic unsigned [PIXEL_DATAW-1:0] y;
269    always_comb begin
270            if(sums_stage_2_reg>255) begin
271                    y = 255;
272            end else if (sums_stage_2_reg<0) begin
273                    y = 0;
274            end else begin
275                    y = sums_stage_2_reg[PIXEL_DATAW-1:0];
276            end
277    end
278
279    // Output interface logics
280    // EGRESS: Stage 8
281    logic unsigned [PIXEL_DATAW-1:0] r_y;
282    logic r_y_valid;
283    logic unsigned [R_X_COL_ADDR_WIDTH-1:0] r_x_col_idx_prev;
284    always_ff @ (posedge clk) begin
285            if(reset) begin
286                    r_y <= 0;
287                    r_y_valid <= 0;
288                    r_x_col_idx_prev <= 0;
289            end else if(enable) begin
290                    r_x_col_idx_prev <=
    ↪   r_x_col_idx_epipelined[NUM_EGRESS_STAGE-1];
291                    // By the time r_x_col_idx is 3, pixel at idx 2 is already
    ↪   written with i_x
292                    if(r_x_col_idx_prev !=
    ↪   r_x_col_idx_epipelined[NUM_EGRESS_STAGE-1] &&
293                            r_x_col_idx_epipelined[NUM_EGRESS_STAGE-1] >=
    ↪   FILTER_SIZE &&
294                            r_x_row_logical_idx_epipelined[NUM_EGRESS_STAGE-1]
    ↪   == R_X_ROWS - 1) begin
295                            r_y <= y;
296                            r_y_valid <= 1;
297                    end else begin
298                            r_y <= 0;
299                            r_y_valid <= 0;
300                    end
```

```
301            end
302    end
303
304    assign o_y = r_y;
305    // Ready for inputs as long as receiver is ready for outputs
306    assign o_ready = i_ready;
307    assign o_valid = r_y_valid & i_ready;
308
309    // End of your code
310
311    endmodule
312
313    /********************************************************************************
314
315    // Multiplier module for 8x8 multiplications + 8x8 multiplications
316    module mult8x8p8x8 (
317            input clk,
318            input enable,
319            input signed [7:0] i_filtera,
320            input signed [8:0] i_pixela0,
321            input signed [8:0] i_pixela1,
322            input signed [7:0] i_filterb,
323            input signed [8:0] i_pixelb,
324            output logic signed [15:0] o_res
325    );
326
327    // Pipeline 0
328    logic signed [8:0] i_pixela0_reg, i_pixela1_reg, i_pixelb_reg;
329    logic signed [7:0] i_filtera_reg, i_filterb_reg;
330    always_ff @(posedge clk) begin
331            if(enable)begin
332                    i_pixela0_reg <= i_pixela0;
333                    i_pixela1_reg <= i_pixela1;
334                    i_pixelb_reg <= i_pixelb;
335                    i_filtera_reg <= i_filtera;
336                    i_filterb_reg <= i_filterb;
337            end
338    end
339
340    // Pipeline 1
341    logic signed [9:0] i_pixela01_reg_reg;
342    logic signed [8:0] i_pixelb_reg_reg;
```

```systemverilog
343    logic signed [7:0] i_filtera_reg_reg, i_filterb_reg_reg;
344    always_ff @ (posedge clk) begin
345            if(enable) begin
346                    i_pixela01_reg_reg <= i_pixela0_reg + i_pixela1_reg;
347                    i_pixelb_reg_reg <= i_pixelb_reg;
348                    i_filtera_reg_reg <= i_filtera_reg;
349                    i_filterb_reg_reg <= i_filterb_reg;
350            end
351    end
352
353    // Pipeline 2
354    logic signed [9:0] i_pixela01_reg_reg_reg;
355    logic signed [8:0] i_pixelb_reg_reg_reg;
356    logic signed [7:0] i_filtera_reg_reg_reg, i_filterb_reg_reg_reg;
357    always_ff @ (posedge clk) begin
358            if(enable) begin
359                    i_pixela01_reg_reg_reg <= i_pixela01_reg_reg;
360                    i_pixelb_reg_reg_reg <= i_pixelb_reg_reg;
361                    i_filtera_reg_reg_reg <= i_filtera_reg_reg;
362                    i_filterb_reg_reg_reg <= i_filterb_reg_reg;
363            end
364    end
365
366    // Pipeline 3
367    always_ff @ (posedge clk) begin
368            if(enable) begin
369                    o_res <= i_pixela01_reg_reg_reg * i_filtera_reg_reg_reg +
       ↪  i_pixelb_reg_reg_reg * i_filterb_reg_reg_reg;
370            end
371    end
372    endmodule
373
374    /************************************************************************
375
376    module pixelram #
377    (
378            parameter FILTER_SIZE = 3,
379            parameter PIXEL_DATAW = 8,
380            parameter IMAGE_WIDTH = 512,
381            parameter R_X_ROWS = FILTER_SIZE,
382            parameter R_X_COL_ADDR_WIDTH = 10
383    )
```

```verilog
384  (
385          input clk,
386          input reset,
387          input enable,
388          // RAM input, unregistered inside the module except for
     ↪   i_read_addr
389          input [R_X_COL_ADDR_WIDTH-1:0] i_write_addr [R_X_ROWS-1:0],
390          input i_write_enable [R_X_ROWS-1:0],
391          input unsigned [PIXEL_DATAW-1:0] i_write_data [R_X_ROWS-1:0],
392          input [R_X_COL_ADDR_WIDTH-1:0] i_read_addr, // registered inside
     ↪   the module
393          // RAM output
394          output unsigned [PIXEL_DATAW-1:0] o_read_data [R_X_ROWS-1:0] //
     ↪   registered, 3 rows of 1 pixel
395  );
396          // Wrap as RAM
397          // 0: [] [] [] [] [] [] [] []                          512 + 2
398          // 1: [] [] [] [] [] [] [] []                          512 + 2
399          // 2: [] [] [] [] [] [] [] []                          512 + 2
400          genvar gen_row;
401          generate
402                  for(gen_row=0; gen_row<R_X_ROWS; gen_row=gen_row+1) begin:
     ↪   pixel_ram_row
403                          pixelrowram pixel_row_ram
404                          (
405                                  .clk(clk),
406                                  .reset(reset),
407                                  .enable(enable),
408                                  // RAM input, unregistered
409                                  .i_write_addr(i_write_addr[gen_row]),
410                                  .i_write_enable(i_write_enable[gen_row]),
411                                  .i_write_data(i_write_data[gen_row]),
412                                  .i_read_addr(i_read_addr), // read 1 pixel
413                                  // RAM output
414                                  .o_read_data(o_read_data[gen_row]) //
     ↪   registered, 1 pixel
415                          );
416                  end
417          endgenerate
418  endmodule
419
420  /************************************************************************
```

```verilog
421
422  module pixelrowram #
423  (
424          parameter PIXEL_DATAW = 8,
425          parameter IMAGE_WIDTH = 512,
426          parameter R_X_COL_WIDTH = IMAGE_WIDTH + 2,
427          parameter R_X_COL_ADDR_WIDTH = 10
428  )
429  (
430          input clk,
431          input reset,
432          input enable,
433          // RAM input, unregistered inside the module except for
     ↪   i_read_addr
434          input [R_X_COL_ADDR_WIDTH-1:0] i_write_addr,
435          input i_write_enable,
436          input unsigned [PIXEL_DATAW-1:0] i_write_data,
437          input [R_X_COL_ADDR_WIDTH-1:0] i_read_addr, // registered inside
     ↪   the module
438          // RAM output
439          output logic unsigned [PIXEL_DATAW-1:0] o_read_data // registered
440  );
441          // Wrap as RAM
442          // 0: [] [] [] [] [] [] [] []                          512 + 2
443          // 2D array of registers for input pixels, row major
444          // set_global_assignment -name
     ↪   ADD_PASS_THROUGH_LOGIC_TO_INFERRED_RAMS OFF
445          logic unsigned [PIXEL_DATAW-1:0] mem [R_X_COL_WIDTH-1:0];
446
447          logic [R_X_COL_ADDR_WIDTH-1:0] i_read_addr_reg;
448
449          integer i;
450          initial begin
451                  for (i=0; i<R_X_COL_WIDTH; i=i+1) begin
452                          mem[i] = 0;
453                  end
454          end
455
456          always_ff @ (posedge clk) begin
457                  if(enable) begin
458                          i_read_addr_reg <= i_read_addr;
459                  end
```

```
460         end
461
462         always_ff @ (posedge clk) begin
463                 if(enable) begin
464                         if(i_write_enable) begin
465                                 mem[i_write_addr] <= i_write_data;
466                         end
467                         o_read_data <= mem[i_read_addr_reg];
468                 end
469         end
470 endmodule
```