

Juan Arturo Garza

September 29th 2022

Elec 522

High Level Design

The purpose of this project is to multiply two matrices systolically; however instead of using simulink and laying out the hardware blocks, it is done in HLS in C++. The challenges I faced were trade offs in design and in the timing that came up.

Base layout (pre optimization):

Modules & Loops	Issue Type	Violation Type	Distance	Slack	Latency(cycles)	Latency(ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	DSR	FF	LUT	URAM
matrixmul	H Violation		-	-	24	240.000	-	25	-	no	0	2	66	365	0
Row_Col	H Violation	Resource Limitation	-	-	22	220.000	7	2	9	yes	-	-	-	-	-

This code is the baseline which I will compare my solution to.

Final Design Choices

I plan to use the following

- HLS Parameters: HLS Array_Reshape, with dim=n
 - This will allow for faster read times of the arrays
- HLS Interface mode= ap_fifo
 - This will allow inputs to be fed systolic way when it is in simulink
- HLS Pipeline
 - This will reduce the clock cycles however can cause timing violations when reading and writing

Design 1:

- HLS Parameters: HLS Array_Reshape, with dim=2 for both arrays
- HLS Interface mode= ap_fifo

Design 2

- HLS Parameters: HLS Array_Reshape, with dim=3 for both arrays
- HLS Interface mode= ap_fifo
- Pipeline II=1
- Matrix output was an Array of arrays which was outputted as one

Design 3 ** FINAL DESIGN

- HLS Parameters: HLS Array_Reshape, with dim=2 for both arrays
- HLS Interface mode= ap_fifo
- Pipeline II=4
 - No Violations

Target	Estimated	Uncertainty	
10.00 ns	5.627 ns	2.70 ns	

Performance & Resource Estimates															
Modules & Loops	Issue Type	Violation Type	Distance	Slack	Latency(cycles)	Latency(ns)	Iteration Latency	Interval	Trip Count	Pipelined	BRAM	DSP	FF	LUT	URAM
matrixmul				-	11	110.000	-	4	-	yes	0	64	1223	598	0

HW Interfaces	
AP_FIFO	
Interface	Data Width
a	64
b	64
res	64
res1	16
res2	16
res3	16
res4	16

TOP LEVEL CONTROL		
Interface	Type	Ports
ap_clk	clock	ap_clk
ap_rst	reset	ap_rst
ap_ctrl	ap_ctrl_hs	ap_done ap_idle ap_ready ap_start

This design has matrix A and Matrix B fed into as one big array and then the output fed out as several rows of the resultant matrix. My design uses Pipelining and a fifo structure. It is able to get results after 14 cycles.

I will not that changing the number type in HLS for matrix A and matrix B would decrease DSP and latency however it would not work in model composer and was thus changed to be shorts. This is a drawback to my design and needs further change.

This design uses more DSP but works in Model composer and passes all test cases

Co-sim Passing:

Cosimulation Report for 'matrixmul'

General Information

Date: Tue Oct 25 08:14:34 CDT 2022

Version: 2022.1 (Build 3526262 on Mon Apr 18 15:48:16 MDT 2022)

Project: solution1

Status: Pass

Solution: solution1 (Vivado IP Flow Target)

Product family: zynq

Target device: xc7z020-clg484-1

Cosim Options

Tool: Vivado XSIM

Dump Trace: all

RTL: Verilog

Performance Estimates

Modules & Loops	Avg II	Max II	Min II	Avg Latency	Max Latency	Min Latency
matrixmul				11	11	11

Console

Errors

Warnings

Guidance

Properties

Man Pages

Git Repositories

Modules/Loops

Vitis HLS Console

add_wave /apath_matrixmul_top/resi_din -into \$tb_return_group -radix hex

add_wave /apath_matrixmul_top/resi_write -into \$tb_return_group -color #ffff00 -radix hex

add_wave /apath_matrixmul_top/resi_full_n -into \$tb_return_group -color #ffff00 -radix hex

add_wave /apath_matrixmul_top/resi_din -into \$tb_return_group -radix hex

set theInputgroup [add_wave_group "C Input" -into \$testbenchgroup]

set tb_return_group [add_wave_group return(fifo) -into \$tbcinoutgroup]

add_wave /apath_matrixmul_top/a_read -into \$tb_return_group -color #ffff00 -radix hex

add_wave /apath_matrixmul_top/a_empty_n -into \$tb_return_group -color #ffff00 -radix hex

add_wave /apath_matrixmul_top/a_dout -into \$tb_return_group -radix hex

add_wave /apath_matrixmul_top/b_read -into \$tb_return_group -color #ffff00 -radix hex

add_wave /apath_matrixmul_top/b_empty_n -into \$tb_return_group -color #ffff00 -radix hex

add_wave /apath_matrixmul_top/b_dout -into \$tb_return_group -radix hex

save_wave_config matrixmul.wcfg

run all

// Inter-Transaction Progress: Completed Transaction / Total Transaction

// Intra-Transaction Progress: Measured Latency / Latency Estimation * 100%

//

// RTL Simulation : "Inter-Transaction Progress" ["Intra-Transaction Progress"] @ "Simulation Time"

//

// RTL Simulation : 0 / 1 [0.00%] @ "125000"

// RTL Simulation : 1 / 1 [100.00%] @ "255000"

//

finish called at time : 315 ns : File "C:/EL6322/jag31/lab1/solution1/solution1/sim/verilog/matrixmul.autotb.v" Line 540

quit

INFO: [Common 17-206] Exiting xsim at Tue Oct 25 08:14:33 2022...

INFO: [COSIM 212-316] Starting C post checking ...

{

{10,10,10,10}

{20,20,20,20}

{30,30,30,30}

{40,40,40,40}

}

Test passed.

INFO: [COSIM 212-1000] *** C/RTL co-simulation finished: PASS ***

INFO: [COSIM 212-211] It is measurable only when transaction number is greater than 1 in RTL simulation. Otherwise, they will be marked as all NA. If user wants to calculate them, please make sure there are at least 2 transactions in RTL simulation.

INFO: [HLS 200-111] Finished Command cosim_design CPU user time: 1 seconds. CPU system time: 0 seconds. Elapsed time: 15.711 seconds; current allocated memory: 2.598 MB.

INFO: [HLS 200-112] Total CPU user time: 2 seconds. Total CPU system time: 1 seconds. Total elapsed time: 17.013 seconds; peak allocated memory: 1.110 GB.

Finished C/RTL cosimulation.

88°F Mostly clear

8:14 AM 10/25/2022

Test Bench

Passed with no errors from hardware and software. Both Aspects agreed on output

```

#include <iostream>
#include "matrixmul.h"

using namespace std;

int main(int argc, char **argv)
{
    mat_a_t in_mat_a[4][4] = {
        {1, 1, 1, 1},
        {2, 2, 2, 2},
        {3, 3, 3, 3},
        {4, 4, 4, 4}
    };
    mat_b_t in_mat_b[4][4] = {
        {1, 1, 1, 1},
        {2, 2, 2, 2},
        {3, 3, 3, 3},
        {4, 4, 4, 4}
    };
    result_t hw_result[4][4], sw_result[4][4];
    //short sw_result[4][4];
    result_t_r1 hw_result_1[4];
    result_t_r2 hw_result_2[4];
    result_t_r3 hw_result_3[4];
    result_t_r4 hw_result_4[4];
    int err_cnt = 0;

    for(int i = 0; i < MAT_A_ROWS; i++) {
        for(int j = 0; j < MAT_B_COLS; j++) {
            sw_result[i][j] = 0;

            for(int k = 0; k < MAT_B_ROWS; k++) {
                sw_result[i][j] += in_mat_a[i][k] * in_mat_b[k][j];
            }
        }
    }
}

```

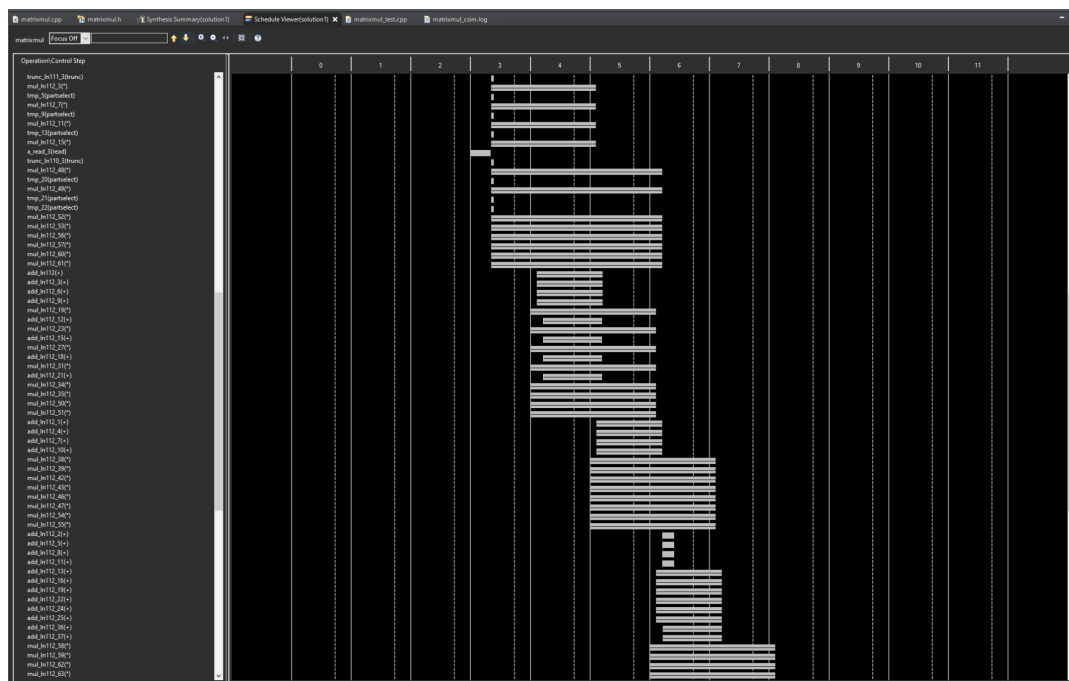
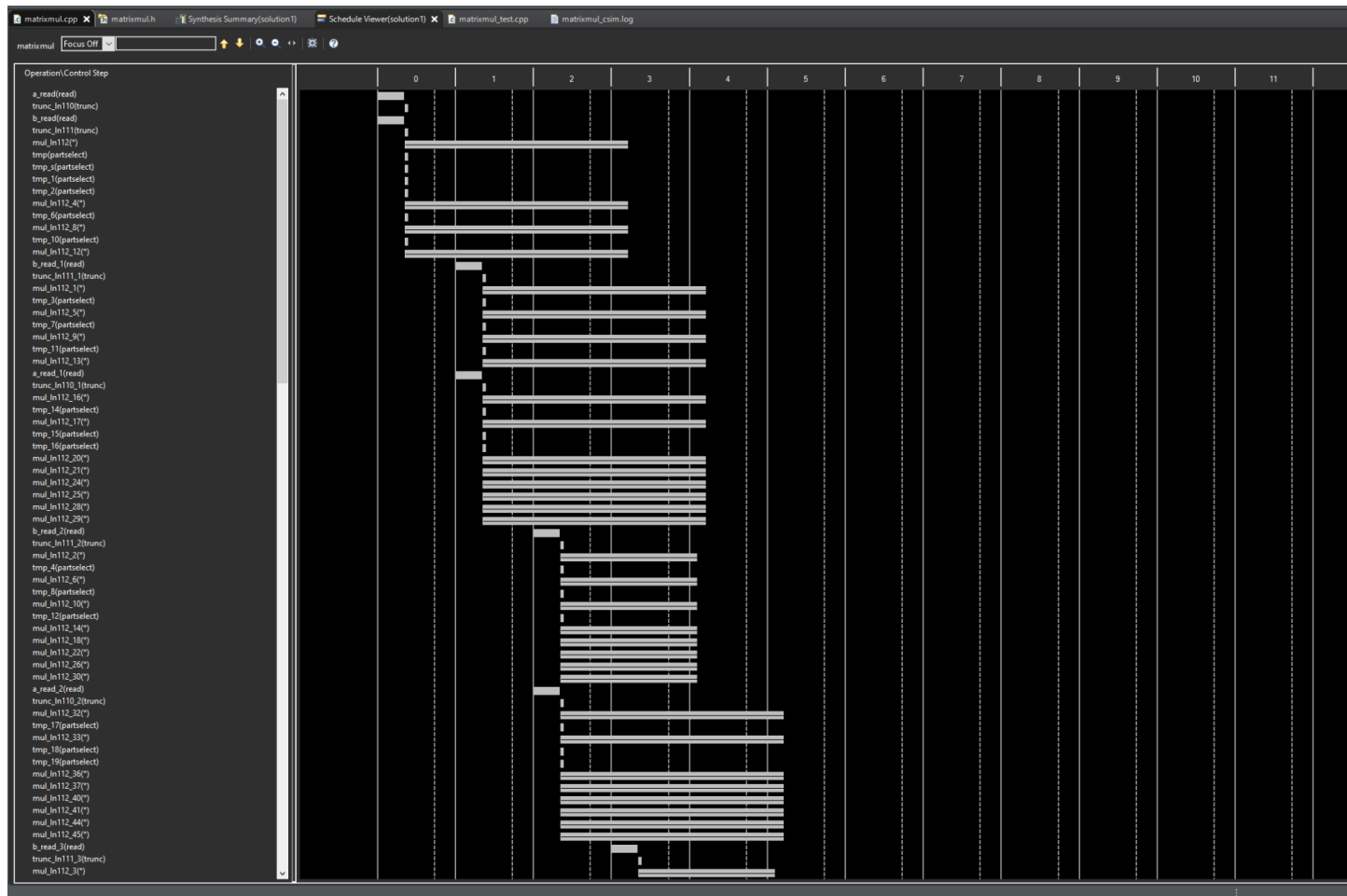
```

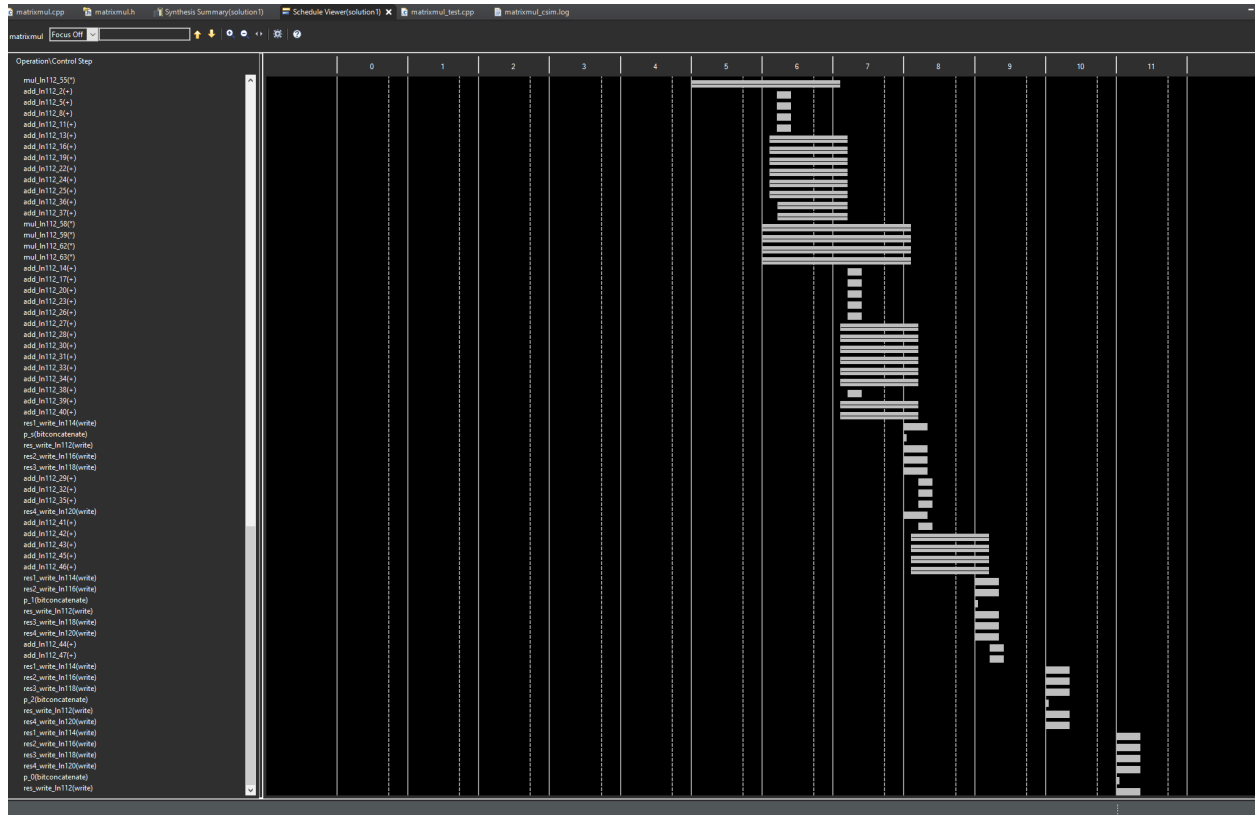
1 INFO: [SIM 2] ***** CSIM start *****
2 INFO: [SIM 4] CSIM will launch GCC as the compiler.
3   Compiling ../../../../matrixmul_test.cpp in debug mode
4   Generating csim.exe
5 {
6 {10,10,10,10}
7 {20,20,20,20}
8 {30,30,30,30}
9 {40,40,40,40}
10 }
11 Test passed.
12 INFO: [SIM 1] CSim done with 0 errors.
13 INFO: [SIM 3] ***** CSIM finish *****
14

```

Note that my testbench compares the row row output which are four arrays to the software result which is one array.

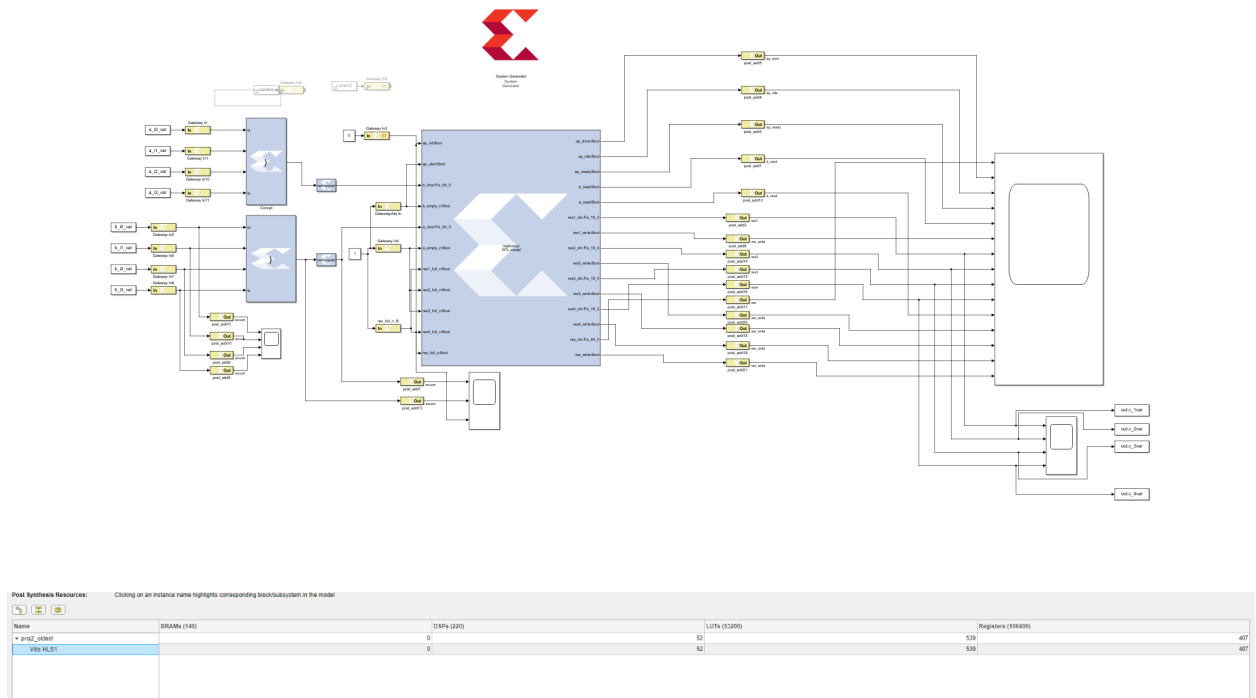
Timing:





The timing sheet shows that there are multiple adds and multiples which most cycle times and overlap, meaning that it is done in parallel. The drawback of this is that it uses more resources to finish faster as more is done in parallel.

Model Composer



Inputs:

```

Editor - C:\ELEC522\jag33\lab1\example_var_data.m
out_disp.m example_var_data.m
1  a_i0_var.time = [0:20];
2  a_i0_var.signals.dimensions = 1;
3  a_i1_var.time = [0:20];
4  a_i1_var.signals.dimensions = 1;
5  a_i2_var.time = [0:20];
6  a_i2_var.signals.dimensions = 1;
7  a_i3_var.time = [0:20];
8  a_i0_var.signals.values = [1 2 3 4 0 0 0 0 5 6 7 8 0 0 0 0 0 0 0 0 0];
9  a_i1_var.signals.values = [1 2 3 4 2 3 0 7 5 6 7 8 0 0 0 0 0 0 0 0 0];
10 a_i2_var.signals.values = [1 2 3 4 6 5 9 1 5 6 7 8 0 0 0 0 0 0 0 0 0];
11 a_i3_var.signals.values = [1 2 3 4 3 5 5 4 5 6 7 8 0 0 0 0 0 0 0 0 0];
12 a_i3_var.signals.dimensions = 1;
13
14 b_i0_var.time = [0:20];
15 b_i0_var.signals.dimensions = 1;
16 b_i1_var.time = [0:20];
17 b_i1_var.signals.dimensions = 1;
18 b_i2_var.time = [0:20];
19 b_i2_var.signals.dimensions = 1;
20 b_i3_var.time = [0:20];
21
22 b_i0_var.signals.values = [1 2 3 4 0 0 0 1 1 2 3 4 0 0 0 0 0 0 0 0 0];
23 b_i1_var.signals.values = [1 2 3 4 0 0 0 1 0 1 2 3 4 0 0 0 0 0 0 0 0];
24 b_i2_var.signals.values = [1 2 3 4 0 1 0 0 1 2 3 4 0 0 0 0 0 0 0 0 0];
25 b_i3_var.signals.values = [1 2 3 4 1 0 0 0 1 2 3 4 0 0 0 0 0 0 0 0 0];
26 b_i3_var.signals.dimensions = 1;
27

```

Note that there is no padding for the inputs of the matrix, every 4 clock cycles there is a new matrix, The matrix is also fed by column, not by row, thus an example: the bottom row of b_i3 is the leftmost column of the matrix, however the output is fed by row instead of column. For Model composer there are three different matrices used and tested one after the other

Output:

```

>> out_disp
NaN NaN NaN NaN NaN NaN NaN 10 10 10 10 3 6 2 0 26 26 26 26 0 0 0 0 0 0 0 0
NaN NaN NaN NaN NaN NaN NaN NaN 20 20 20 20 5 5 3 0 52 52 52 52 0 0 0 0 0 0 0 0
NaN NaN NaN NaN NaN NaN NaN NaN 30 30 30 30 5 9 0 0 78 78 78 78 0 0 0 0 0 0 0 0
NaN NaN NaN NaN NaN NaN NaN NaN 40 40 40 40 4 1 7 0 104 104 104 104 0 0 0 0 0 0 0 0
>> |

```

It takes a total of 14 cycles to calculate the initial result however each new matrix come directly after the other meaning that the max output will be 1 maxtrix result every 4 cycles, which is 40 ns, and means 250 million results per second assuming no stalling and no delay

Comparison to Project 2

Project 2 took only 11 cycles however there were zeros padding between each set of new coefficients and results were staggered compared to Proj3 implementation which gets all results at once. This means Project 2 has only 200 Million results per second as it takes 5 cycles between every result updating. The Drawback of my design is the resource usages.