

ELEC522 – Fall 2022

**Project 1- Xilinx Vitis Model Composer
/ Nvidia CUDA Tutorial Assignment**

By

Hsuan-You (Shaun) Lin

Rice University, Houston, TX

Sep. 16, 2022

Supervised by

Dr. Joseph R. Cavallaro

Contents

Contents	2
List of Figures	3
List of Tables.....	4
Xilinx Lab:	5
Lab1: Introduction to Vitis Model Composer HDL Library	5
Lab1_1:	5
Lab1_2:	6
Lab1_3:	7
Lab1_4: (Part 1: Designing with Floating-Point Data Types)	8
Lab1_4: (Part 2: Designing with Fixed-Point Data Types).....	9
Lab2: Importing Code into a Vitis Model Composer HDL Design	11
Lab2_1:	11
Lab2_2:	13
Lab3: Timing and Resource Analysis	14
Lab3_1:	14
Lab3_2:	14
CUDA Lab:	16

List of Figures

Figure 1.....	5
Figure 2.....	5
Figure 3.....	6
Figure 4.....	6
Figure 5.....	7
Figure 6.....	7
Figure 7.....	8
Figure 8.....	8
Figure 9.....	8
Figure 10.....	9
Figure 11.....	9
Figure 12.....	10
Figure 13.....	10
Figure 14.....	12
Figure 15.....	12
Figure 16.....	13
Figure 17.....	13
Figure 18.....	14
Figure 19.....	14
Figure 20.....	15
Figure 21.....	16
Figure 22.....	17

List of Tables

Table 1. arraySize 10 ~ 1250 performance result	16
---	----

Xilinx Lab:

Lab1: Introduction to Vitis Model Composer HDL Library

Step 1: Creating a Design in an FPGA

Lab1_1:

(a) Screen capture of spectrum plots of the initial waveforms

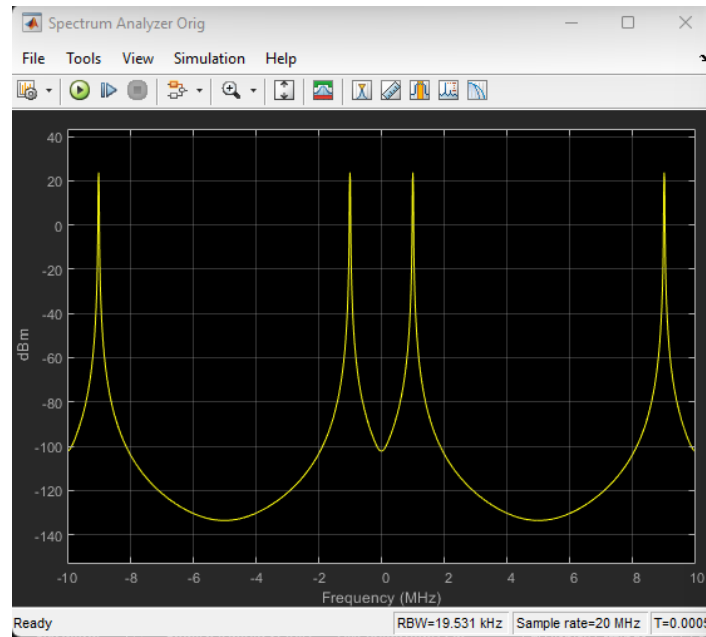


Figure 1.

(b) Screen capture of spectrum plots after adding Digital FIR Filter block

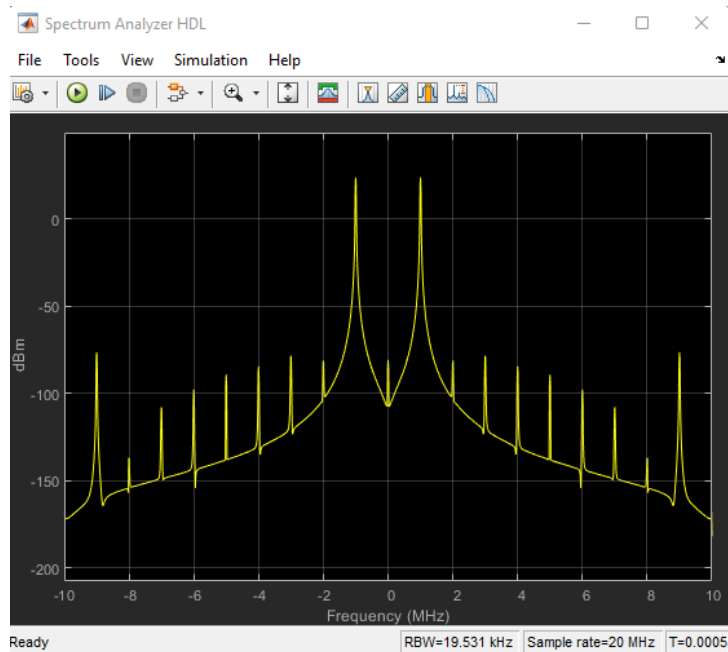


Figure 2.

(c) Screen capture of resource utilization output



Name	BRAMs (140)	DSPs (220)	LUTs (53200)	Registers (106400)
Lab1_1	0	0	5	281
Digital FIR Filter	0	0	5	281

Figure 3.

(d) Screen capture of final block diagram

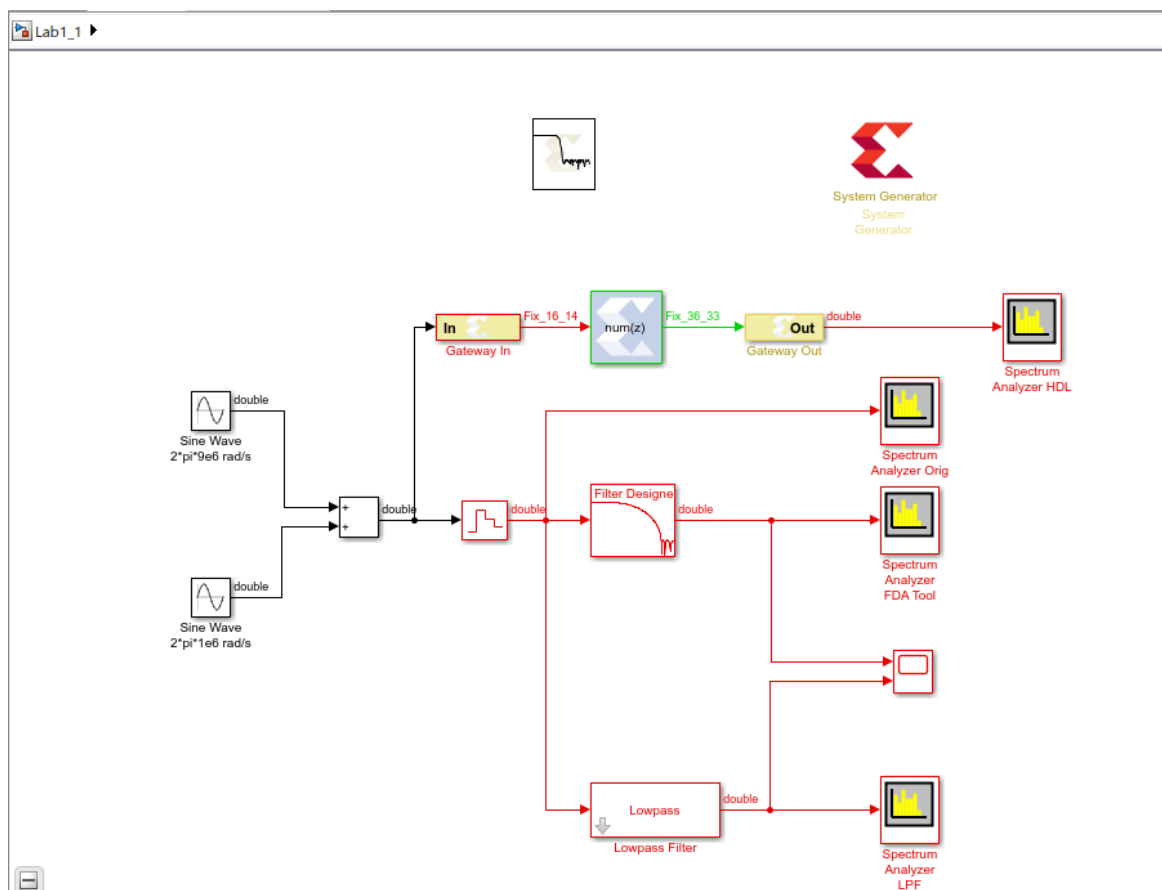


Figure 4.

(e) Compare your results with the tutorial results and note any differences or if the same.

My answer: No differences, I got the same results with the tutorial results.

Step 2: Creating an Optimized Design in an FPGA

Lab1_2:

(a) Screen capture of resource utilization output for higher frequency in Generate step

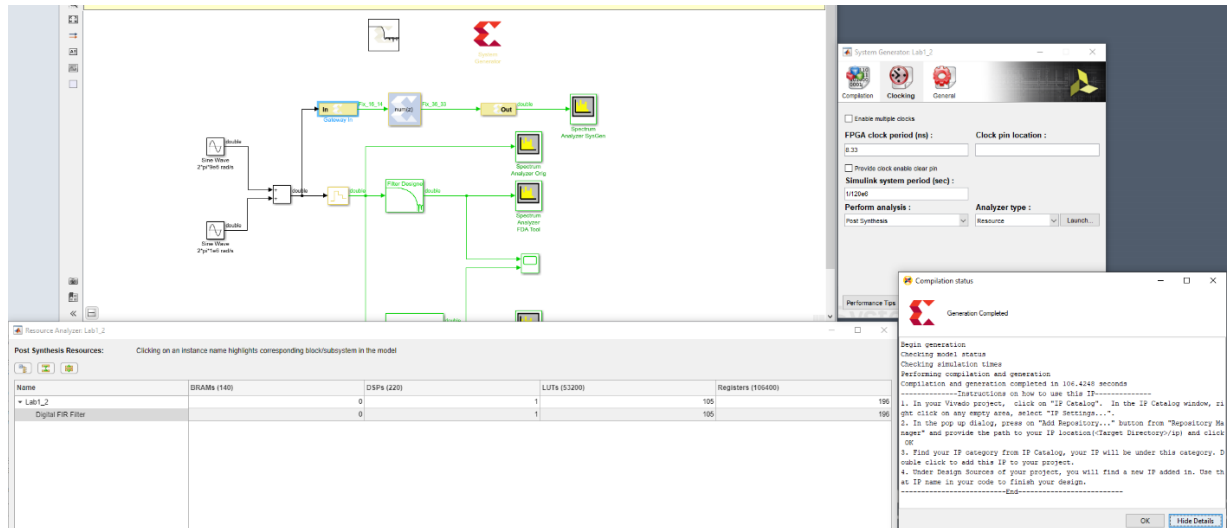


Figure 5.

(b) Compare your results with the tutorial results and note any differences or if the same.

My answer: I got the differences results with the tutorial results, but it seems like the tutorial results is wrong, because the title of the resource utilization output is “Lab1_4_1_sol” with the tutorial results not “Lab1_2”.

Step 3: Creating a Design using Discrete Components

Lab1_3:

(a) Screen capture of final block diagram using discrete components

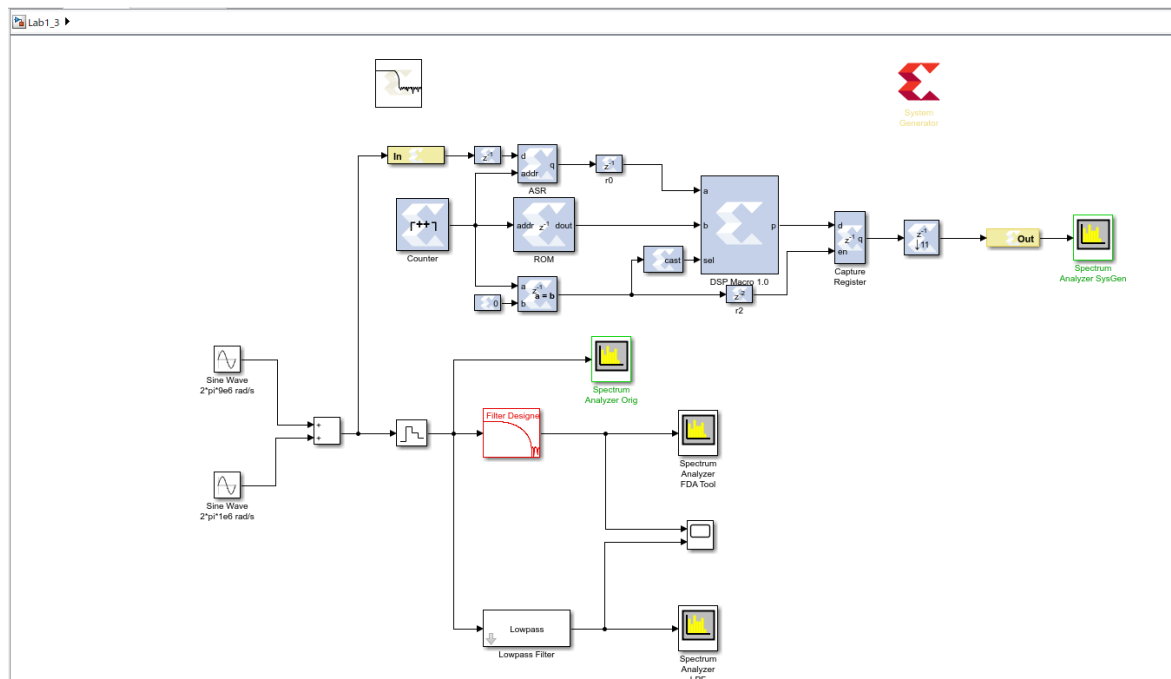


Figure 6.

(b) Screen capture of spectrum plots

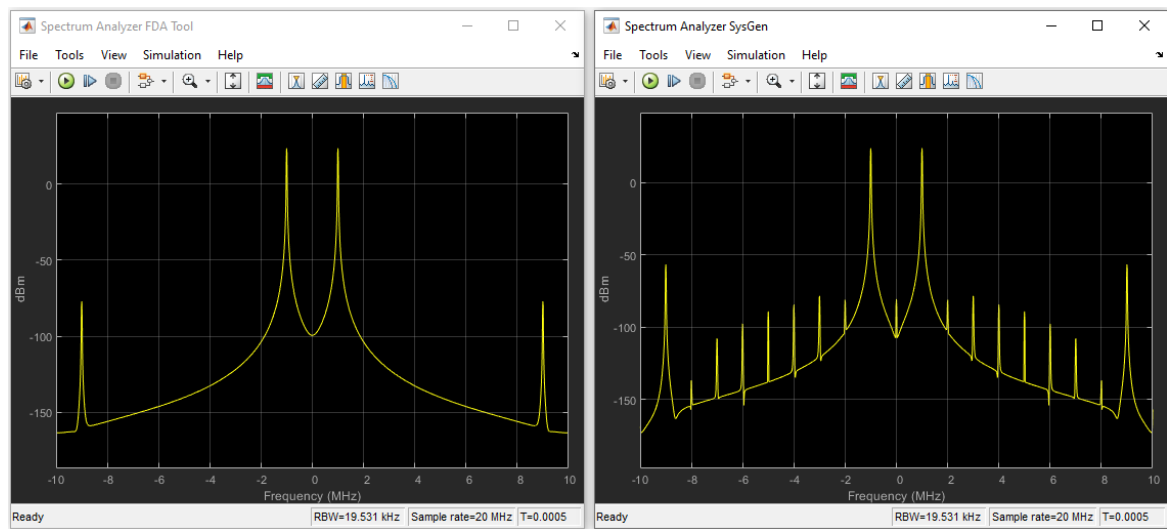


Figure 7.

(c) Screen capture of resource utilization output

Resource Analyzer: Lab1_3

Post Synthesis Resources: Clicking on an instance name highlights corresponding block/system in the model

Name	SRAMs (140)	DSFs (220)	LTUs (51200)	Registers (105400)	
Lab1_3	0.0000	0.0000	1	22	130
RD00	0.0000	0.0000	0	0	0
Reconstruction1	0	0	1	0	1
r3	0	0	0	0	10
r2	0	0	0	1	1
r1	0	0	0	0	10
DSP Macro 1.0	0	0	1	2	20
Down Sampler1	0	0	0	0	40
Counter	0	0	0	2	4
CaptureRegister	0	0	0	0	40
AD00	0	0	0	10	0

Figure 8.

(d) Compare your results with the tutorial results and note any differences or if the same.

My answer: No differences, I got the same results with the tutorial results.

Step 4: Working with Data Types

Lab1_4: (Part 1: Designing with Floating-Point Data Types)

(a) Screen capture of spectrum plots of **floating-point** version

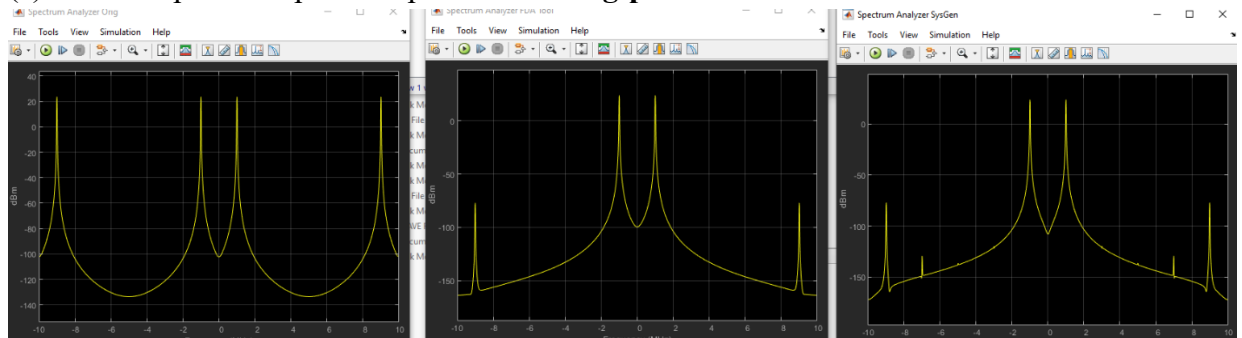
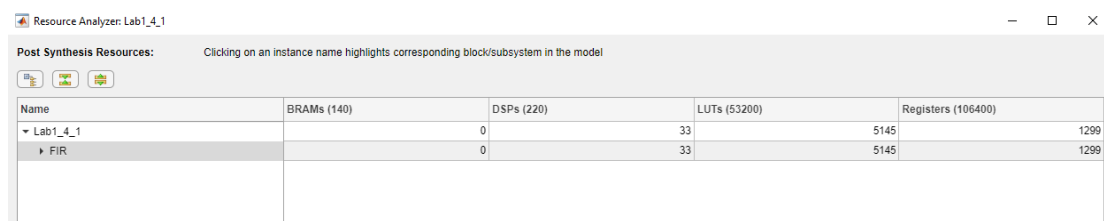


Figure 9.

(b) Screen capture of **floating-point** resource utilization output



Resource Analyzer: Lab1_4_1

Post Synthesis Resources: Clicking on an instance name highlights corresponding block/subsystem in the model

Name	BRAMs (140)	DSPs (220)	LUTs (53200)	Registers (106400)
Lab1_4_1	0	0	33	5145
FIR	0	0	33	5145

Figure 10.

(c) Compare your results with the tutorial results and note any differences or if the same.

My answer: No differences, I got the same results with the tutorial results.

Lab1_4: (Part 2: Designing with Fixed-Point Data Types)

(d) Screen capture of scope results for **fixed-point** version – initial version

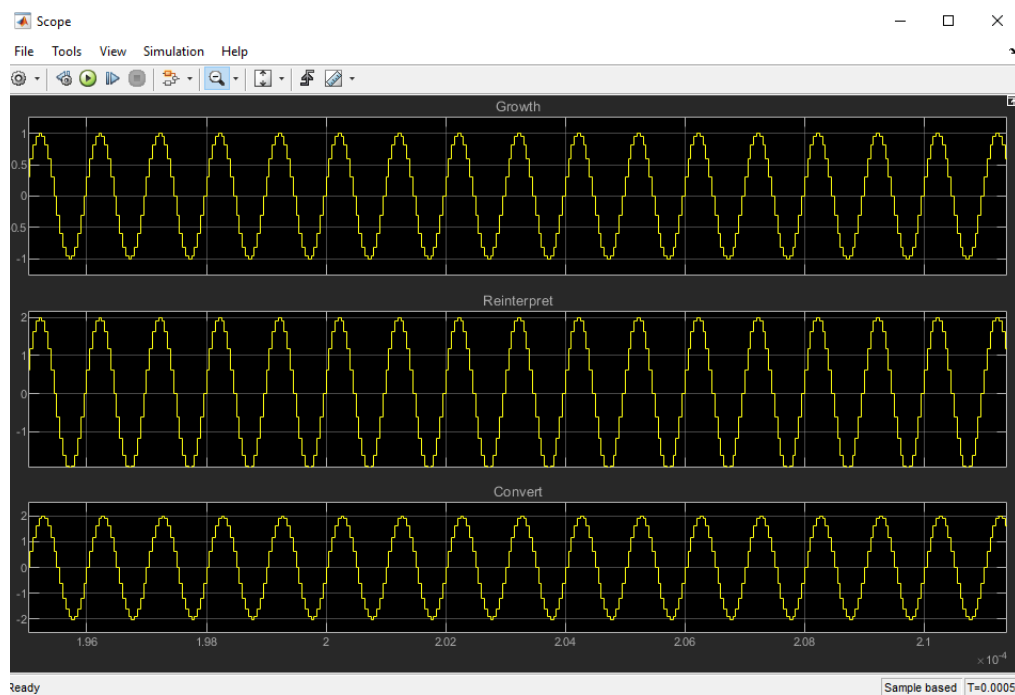


Figure 11.

(e) Screen capture of scope results for **fixed-point** version – after changing reinterpret and convert blocks

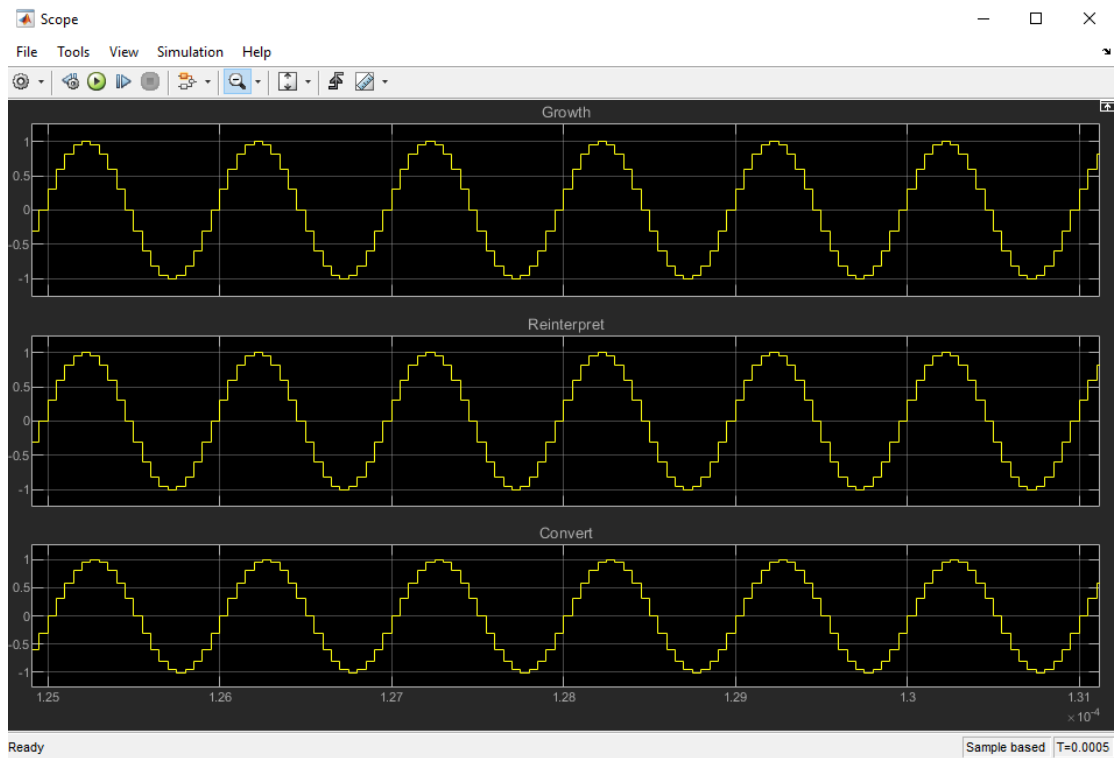


Figure 12.

(f) Screen capture of **fixed-point** resource utilization output

Resource Analyzer: Lab1_4_2

Post Synthesis Resources: Clicking on an instance name highlights corresponding block/subsystem in the model

Name	BRAMs (140)	DSPs (220)	LUTs (53200)	Registers (106400)
Lab1_4_2	0	44	5690	1893
FIR-Fixed-Point	0	11	545	578
FIR	0	33	5145	1299
Convert	0	0	0	16

Figure 13.

(g) Compare your results with the tutorial results and note any differences or if the same.

My answer: No differences, I got the same results with the tutorial results.

Lab2: Importing Code into a Vitis Model Composer HDL Design

Step 1: Modeling Control with M-Code

Lab2_1:

(a) state_machine.m code:

```
function matched = state_machine(din)
persistent state, state = xl_state(0,{xlUnsigned, 3, 0});
switch state
    case 0
        if din == 1
            state = 1;
        else
            state = 0;
        end
        matched = 0;
    case 1
        if din == 0
            state = 2;
        else
            state = 1;
        end
        matched = 0;
    case 2
        if din == 1
            state = 3;
        else
            state = 0;
        end
        matched = 0;
    case 3
        if din == 1
            state = 4;
        else
            state = 2;
        end
        matched = 0;
    otherwise
        if din == 1
            state = 1;
        else
            state = 0;
        end
        matched = 1;
end
```

(b) Screen capture of final block diagram

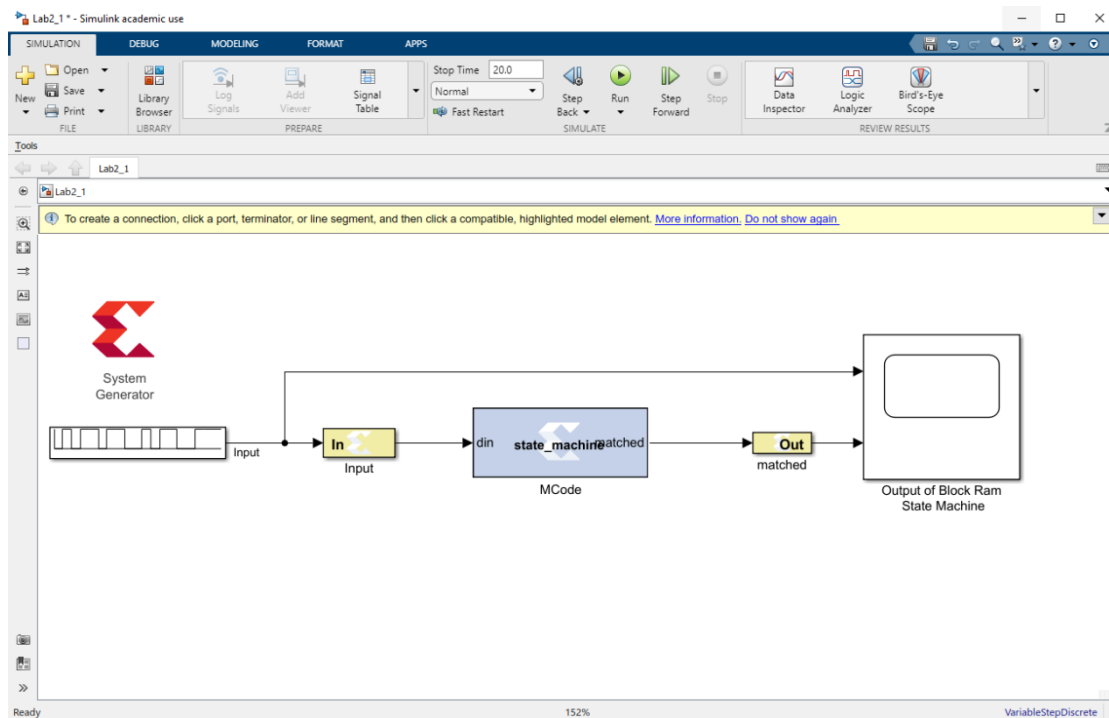


Figure 14.

(c) Screen capture of output waveform

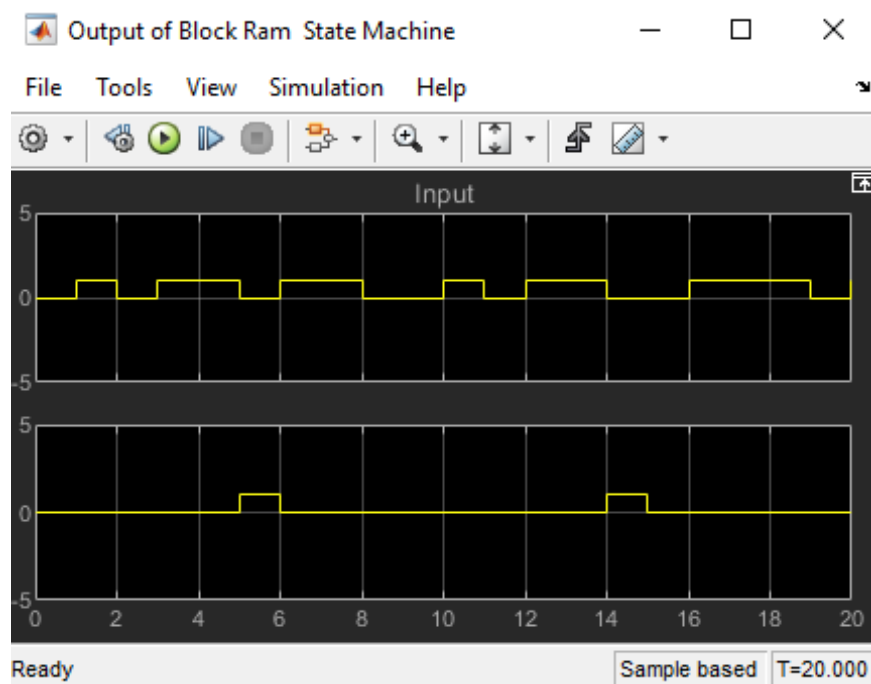


Figure 15.

(d) Compare your results with the tutorial results and note any differences or if the same.

My answer: No differences, I got the same results with the tutorial results.

Lab2_2:

Step 2: Modeling Blocks with HDL

Lab2_2:

(a) Screen capture of final block diagram

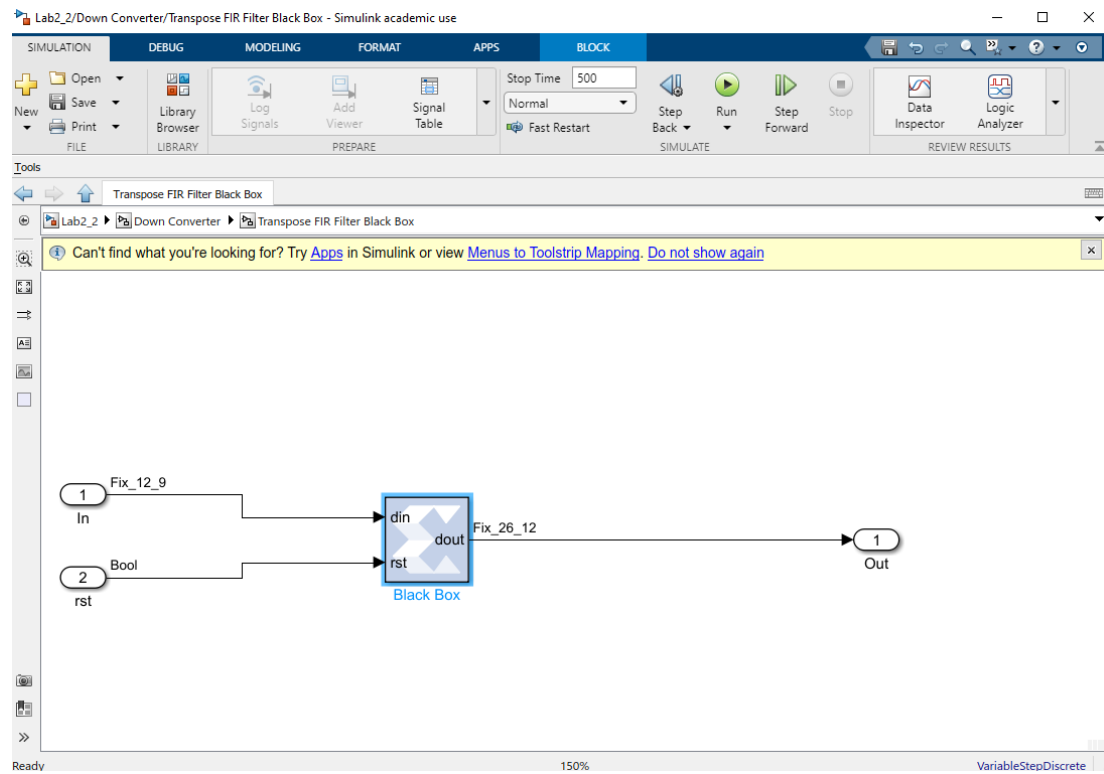


Figure 16.

(b) Screen capture of final output waveform

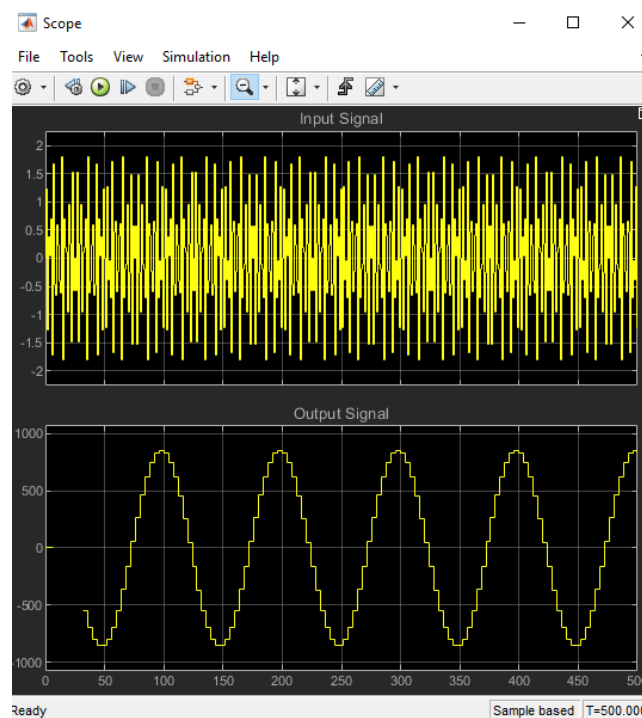


Figure 17.

(c) Compare your results with the tutorial results and note any differences or if the same.

My answer: No differences, I got the same results with the tutorial results.

Lab3: Timing and Resource Analysis

Step 1: Timing Analysis in Vitis Model Composer

Lab3_1:

(a) Screen capture of timing analyzer results – realize that the Kintex -3 part will have less timing error than the Zynq Artix -1 part

Post Synthesis Timing Paths										
Clicking on an instance name highlights corresponding block/subsystem in the model										
Violation type: Setup Hold Y Setup Counts Status: FAILED										
Stack (ns)	Delay (ns)	Logic Delay (ns)	Routing Delay (ns)	Levels of Logic	Source	Destination	Source Clock	Destination Clock	Path Constraints	
1	3.2000	9.2200	4.4810	0.6900	lsc3subsystem1M0A1	lsc3subsystem1M0A2	ck	ck	create_clock name ck period 2 [get_ports clk]	
2	-0.6590	2.0840	2.3550	0.4790	lsc3sub_genRegister4	lsc3sub_genRegister4	ck	ck	create_clock name ck period 2 [get_ports clk]	
3	-0.6590	2.0840	1.9540	0.4910	lsc3sub_genRegister10	lsc3sub_genRegister10	ck	ck	create_clock name ck period 2 [get_ports clk]	
4	-0.7250	2.0010	1.9520	0.5020	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]	
5	-0.7150	1.9010	0.9010	0.0000	lsc3subsystem1M0A1	lsc3subsystem1M0A1	ck	ck	create_clock name ck period 2 [get_ports clk]	
6	-0.5150	1.7150	0.5150	0.0000	lsc3sub_genRegister4	lsc3subsystem1M0A1	ck	ck	create_clock name ck period 2 [get_ports clk]	
7	-0.6000	1.5500	0.5150	1.0010	lsc3sub_genRegister7	lsc3subsystem1M0A1	ck	ck	create_clock name ck period 2 [get_ports clk]	
8	0.4170	1.5170	1.5170	0.0000	lsc3sub_genRegister3	lsc3sub_genRegister3	ck	ck	create_clock name ck period 2 [get_ports clk]	
9	0.4480	1.5770	1.5880	0.0000	lsc3sub_genRegister7	lsc3sub_genRegister7	ck	ck	create_clock name ck period 2 [get_ports clk]	
10	0.7190	1.2330	0.9130	0.4790	lsc3sub_genRegister7	lsc3sub_genRegister7	ck	ck	create_clock name ck period 2 [get_ports clk]	
11	0.9250	0.8940	0.4730	0.4790	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]	
12	0.9490	0.8520	0.5180	0.2340	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]	
13	0.9100	0.8760	0.5180	0.3480	lsc3sub_genRegister10	lsc3sub_genRegister10	ck	ck	create_clock name ck period 2 [get_ports clk]	
14	0.8630	0.8750	0.5180	0.3480	lsc3sub_genRegister3	lsc3sub_genRegister3	ck	ck	create_clock name ck period 2 [get_ports clk]	
15	0.8660	0.8510	0.5180	0.3440	lsc3subsystem1M0A2	lsc3sub_genRegister3	ck	ck	create_clock name ck period 2 [get_ports clk]	
16	0.8770	0.8520	0.5180	0.3440	lsc3sub_genRegister3	lsc3sub_genRegister3	ck	ck	create_clock name ck period 2 [get_ports clk]	
17	0.8170	0.8520	0.5180	0.3440	lsc3sub_genRegister3	lsc3sub_genRegister3	ck	ck	create_clock name ck period 2 [get_ports clk]	

Figure 18.

(b) Modify latency in various blocks in order to pass timing. This may be different than the suggestion in the tutorial. Explain your process.

My answer:

After I changed the latency from 1 to 2, I still got a bunch of errors, but I thought it was because the Zedboard couldn't handle so many calculations, so I changed the board to Zynq UltraScale+ and it passed.

(c) Screen capture of timing analyzer results – passing – after modifications

Post Synthesis Timing Paths

Violation Type

Setup

Hold

Y

Setup Counts

Status: Passing

Clicking on an instance name highlights corresponding block/subsystem in the model

Stack (ns)	Delay (ns)	Logic Delay (ns)	Routing Delay (ns)	Levels of Logic	Source	Destination	Source Clock	Destination Clock	Path Constraints
2	1.0380	0.8840	0.8840	0	lsc3sub_genRegister1M0A1	lsc3subsystem1M0A2	ck	ck	create_clock name ck period 2 [get_ports clk]
3	1.1290	0.9150	0.9150	2	lsc3sub_genRegister1M0A1	lsc3subsystem1M0A1	ck	ck	create_clock name ck period 2 [get_ports clk]
4	1.1590	0.7890	0.4090	2	lsc3sub_genRegister1M0A1	lsc3sub_genRegister1M0A1	ck	ck	create_clock name ck period 2 [get_ports clk]
5	1.2090	0.6890	0.4090	5	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]
6	1.3270	0.6770	0.3070	9	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]
7	1.3750	0.7750	0.3750	9	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]
8	1.4540	0.5770	0.3070	3	lsc3sub_genRegister5	lsc3sub_genRegister5	ck	ck	create_clock name ck period 2 [get_ports clk]
9	1.4890	0.4270	0.1090	1	lsc3sub_genRegister7	lsc3sub_genRegister7	ck	ck	create_clock name ck period 2 [get_ports clk]
10	1.5090	0.3910	0.0790	9	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]
11	1.6130	0.2050	0.0790	9	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]
12	1.7580	0.2760	0.0790	9	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]
13	1.7590	0.2760	0.0790	9	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]
14	1.7590	0.2890	0.1090	1	lsc3sub_genRegister1	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]
15	1.7170	0.2440	0.0790	9	lsc3subsystem1M0A2	lsc3sub_genRegister1	ck	ck	create_clock name ck period 2 [get_ports clk]
16	1.7590	0.2650	0.0790	9	lsc3sub_genRegister3	lsc3sub_genRegister3	ck	ck	create_clock name ck period 2 [get_ports clk]
17	1.7180	0.2650	0.0790	9	lsc3sub_genRegister3	lsc3sub_genRegister3	ck	ck	create_clock name ck period 2 [get_ports clk]
18	1.7180	0.2650	0.0790	9	lsc3sub_genRegister3	lsc3sub_genRegister3	ck	ck	create_clock name ck period 2 [get_ports clk]
19	1.7180	0.2650	0.0790	9	lsc3sub_genRegister3	lsc3sub_genRegister3	ck	ck	create_clock name ck period 2 [get_ports clk]

Figure 19.

Step 2: Resource Analysis in Vitis Model Composer

Lab3_2:

(a) Screen capture of resource analyzer results for the Zynq Artix chip

Resource Analyzer: Lab3

Post Synthesis Resources: Clicking on an instance name highlights corresponding block/subsystem in the model

Name	BRAMs (365)	DSPs (740)	LUTs (134600)	Registers (269200)
Lab3	0.5000	1	153	321
↳ subsystem1	0.5000	1	97	97
Register1	0	0	0	1
Register	0	0	0	48
Delay7	0	0	1	1
Delay4	0	0	0	20
Delay3	0	0	1	1
Delay2	0	0	0	48
↳ addr_gen	0	0	54	105

Figure 20.

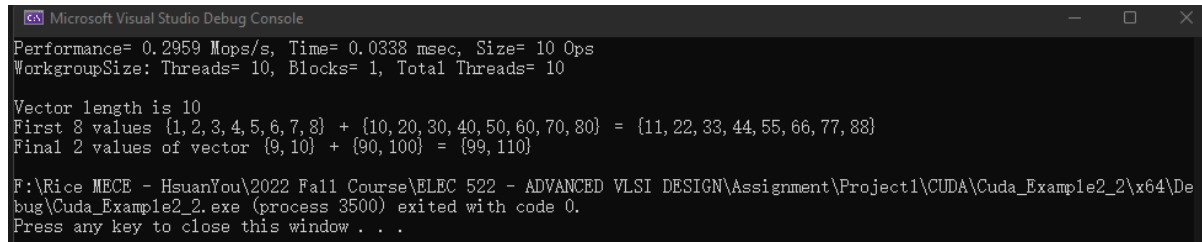
(b) Compare your results with the tutorial results and note any differences or if the same.

My answer: Basically no differences, there are just a few slightly different in the registers, I got almost the same results with the tutorial results.

CUDA Lab:

<Experiments with the Project 1 CUDA program>

(1) Screen capture of debug results window for arraySize = 10, threads = size, blocks = 1 initial parameters.



```
Microsoft Visual Studio Debug Console
Performance= 0.2959 Mops/s, Time= 0.0338 msec, Size= 10 Ops
WorkgroupSize: Threads= 10, Blocks= 1, Total Threads= 10

Vector length is 10
First 8 values {1, 2, 3, 4, 5, 6, 7, 8} + {10, 20, 30, 40, 50, 60, 70, 80} = {11, 22, 33, 44, 55, 66, 77, 88}
Final 2 values of vector {9, 10} + {90, 100} = {99, 110}

F:\Rice MECE - HsuanYou\2022 Fall Course\ELEC 522 - ADVANCED VLSI DESIGN\Assignment\Project1\CUDA\Cuda_Example2_2\x64\De
bug\Cuda_Example2_2.exe (process 3500) exited with code 0.
Press any key to close this window . . .
```

Figure 21.

(2) Include table of results for all of the arraySize values listed in the problem: 100, 200, 500, 750, 1000, 1250, 1500, 2000, 5000, 10000, 25000, and 64000. Table should record: Performance in Mops/s, Time, Size, Threads, Blocks, and Total Threads.

Table 1. arraySize 10 ~ 1250 performance result

arraySize	Performance (Mops/s)	Time (msec)	Size (Ops)	Threads (Set)	Blocks (Set)	Threads (Total)
10	0.2959	0.0338	10	10	1	10
100	3.6295	0.0276	100	100	1	100
200	6.7495	0.0296	200	200	1	200
500	20.4248	0.0245	500	500	1	500
750	26.1579	0.0287	750	750	1	750
1000	37.5601	0.0266	1000	1000	1	1000
1250	50.7965	0.0246	1250	1024	2	2048
1500	54.2535	0.0276	1500	1024	2	2048
2000	69.7545	0.0287	2000	1024	2	2048
5000	195.5570	0.0256	5000	1024	5	5120
10000	338.9371	0.0295	10000	1024	10	10240
25000	915.8852	0.0273	25000	1024	25	25600
64000	2040.8164	0.0314	64000	1024	63	64512

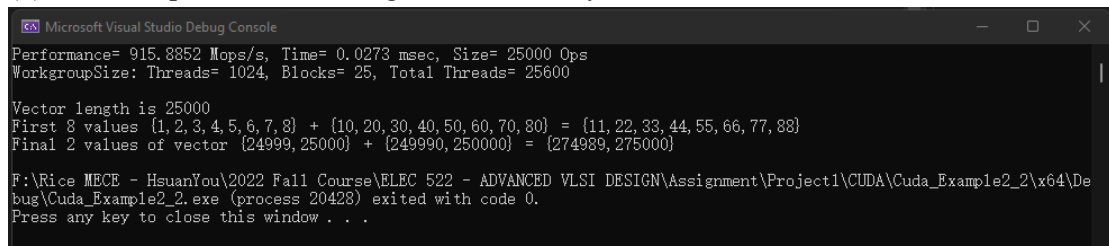
(3) What happens at arraySize = 1250 with the original threads and blocks value? Explain.

My answer:

It will show “Maximum threads/block is 1024. Increase block size and reduce threads/block”, as the result in *Figure 1*. This is because the maximum number of threads per block is 1024 in the CUDA specifications, if we set arraySize = 1250, it

will exceed that value, after which there is no register space to save thread state and to limit memory and synchronization contention.

(4) Screen capture of the Debug results for arraySize = 25000



```
Microsoft Visual Studio Debug Console
Performance= 915.8852 Mops/s, Time= 0.0273 msec, Size= 25000 Ops
WorkgroupSize: Threads= 1024, Blocks= 25, Total Threads= 25600

Vector length is 25000
First 8 values {1, 2, 3, 4, 5, 6, 7, 8} + {10, 20, 30, 40, 50, 60, 70, 80} = {11, 22, 33, 44, 55, 66, 77, 88}
Final 2 values of vector {24999, 25000} + {249990, 250000} = {274989, 275000}

F:\Rice MECE - HsuanYou\2022 Fall Course\ELEC 522 - ADVANCED VLSI DESIGN\Assignment\Project1\CUDA\Cuda_Example2_2\x64\De
bug\Cuda_Example2_2.exe (process 20428) exited with code 0.
Press any key to close this window . . .
```

Figure 22.

(5) Final two values of variable c at 64000 with minimum threads and blocks set for 25000. How did you fix these values.

My answer:

If we used the minimum threads and blocks size for the arraySize = 25000, which means threads = 1024 and blocks = 25, then increase the arraySize = 64000 will shows that the final 2 values of vector {63999, 64000} + {639990, 640000} = {0, 0}, this means you start blocks with thread count that is not divisible by the warp size, the hardware will simply execute the last warp with some of the threads "masked out" (i.e. they do have to execute, but without any effect on the state of the GPU/memory).

The easiest solution is increasing the number of blocks, let threads*blocks >= arraySize, in this case we set blocks = 63, which means we'll have total 64512 threads as shown in *Table 1*.

(6) Analyze your table for the percentage increase in performance, time, and size.

My answer:

From *Table 1*, we can see that as the arraySize increases, performance also increases, which means performance and size has the proportional relationship. In contrast, time is almost unchanged, because parallel threads to execute on the GPU, so time and size are irrelevant.