

### 香港中文大學 The Chinese University of Hong Kong

CENG3430 Rapid Prototyping of Digital Systems

# Lab08: High Level Synthesis Exercise



### **Task**



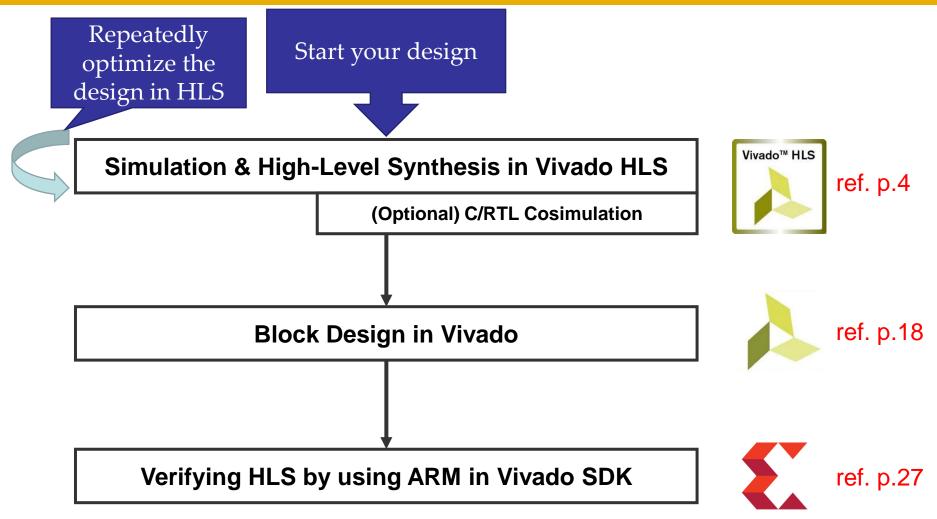
# Accelerating Floating Point Matrix Multiplication with High-Level Synthesis

#### mmult.h

$$\begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \\ a_7 & a_8 & a_9 \end{bmatrix} \begin{bmatrix} b_1 & b_2 & b_3 \\ b_4 & b_5 & b_6 \\ b_7 & b_8 & b_9 \end{bmatrix} = \begin{bmatrix} c_1 & c_2 & c_3 \\ c_4 & c_5 & c_6 \\ c_7 & c_8 & c_9 \end{bmatrix}$$

### Workflow







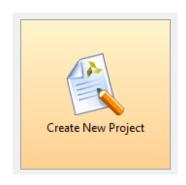
# **High-Level Synthesis**

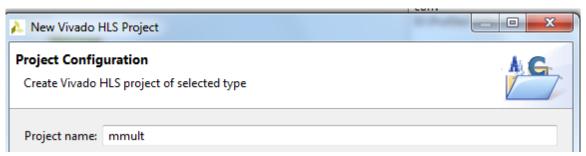
### Step 1 - High-Level Synthesis (HLS)



- Download the lab08 material from Blackboard
- Open Vivado HLS 2016.3 and Create a new project called "mmult"





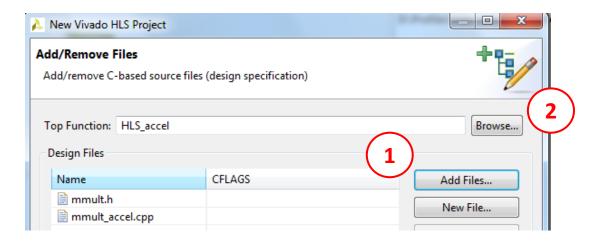


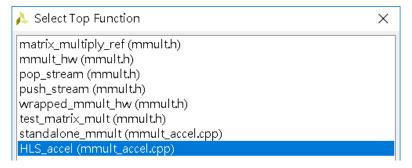
Remarks: Avoid naming your folder with space ""
e.g. "C:\Users\Lab 8" may cause some trouble
Change to Lab 8 or Lab08 to avoid errors from Vivado

### Step 2 - High-Level Synthesis (HLS)



 Add the two C-based source files (mmult.h and mmult\_accel.cpp) into the project and specify the "HLS\_accel" as Top Function by click Browse:

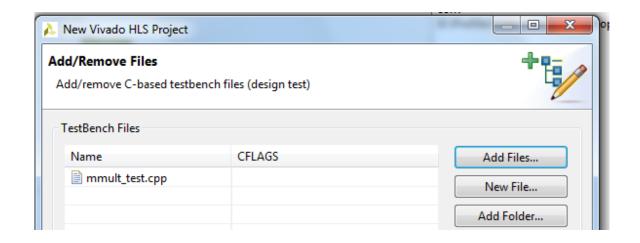




### Step 3 - High-Level Synthesis (HLS)



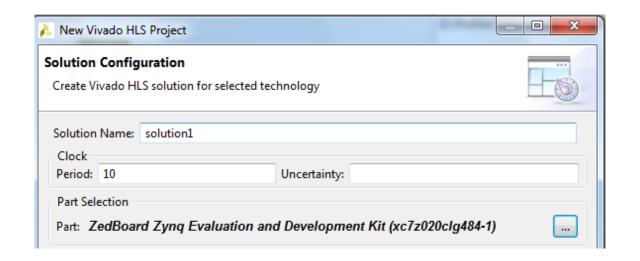
 Add the C-based testbench file (mmult\_test.cpp) into the project:



### Step 4 - High-Level Synthesis (HLS)



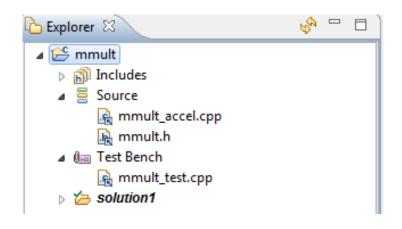
 Specify Zedboard as the "Part", and make sure that the clock period (i.e., the timing constraint) is specified as "10" (ns) since the on-board clock frequency is 100 MHz and the clock uncertainty is left empty.



### Step 5 - High-Level Synthesis (HLS)



Finish it and check the file structure:

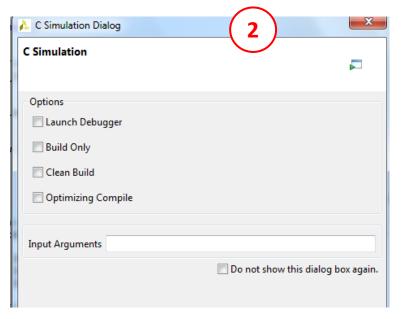


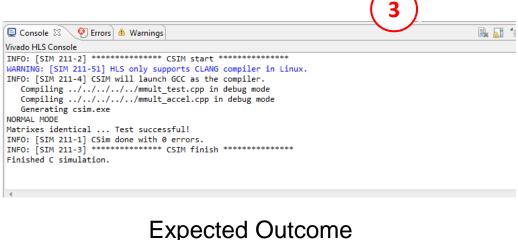
### Step 6 - High-Level Synthesis (HLS)



- Run C Simulation
  - Simply click OK in C Simulation Dialog







### Step 7 - High-Level Synthesis (HLS)



### Run C Synthesis



```
Vivado HLS Console

INFO: [RTGEN 206-100] Generating core module 'HLS_accel_fadd_32bkb': 1 instance(s).

INFO: [RTGEN 206-100] Generating core module 'HLS_accel_fmul_32cud': 1 instance(s).

INFO: [RTGEN 206-100] Finished creating RTL model for 'HLS_accel'.

INFO: [RTGEN 206-100] Finished creating RTL model for 'HLS_accel'.

INFO: [HLS 200-111] Elapsed time: 0.14 seconds; current allocated memory: 94.264 MB.

INFO: [RTMG 210-278] Implementing memory 'HLS_accel_a_ram' using block RAMs.

INFO: [HLS 200-111] Finished generating all RTL models Time (s): cpu = 00:00:04; elapsed = 00:00:23. Memory (MB): peak = 136.422;

INFO: [SYSC 207-301] Generating SystemC RTL for HLS_accel.

INFO: [VHDL 208-304] Generating VHDL RTL for HLS_accel.

INFO: [VLOG 209-307] Generating Verilog RTL for HLS_accel.

INFO: [HLS 200-112] Total elapsed time: 22.766 seconds; peak allocated memory: 94.264 MB.

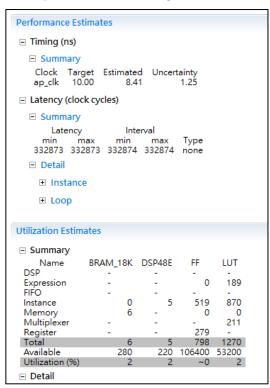
Finished C synthesis.
```

### **Expected Outcome**

### Step 8 - High-Level Synthesis (HLS)



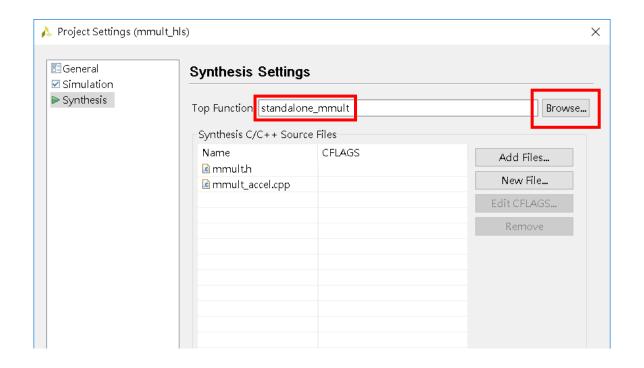
- Check the "Performance Estimates" and "Utilization Estimates" in the "Synthesis Report for HLS\_accel"
  - Solution -> Open Report -> Synthesis



By default, this window will pop out

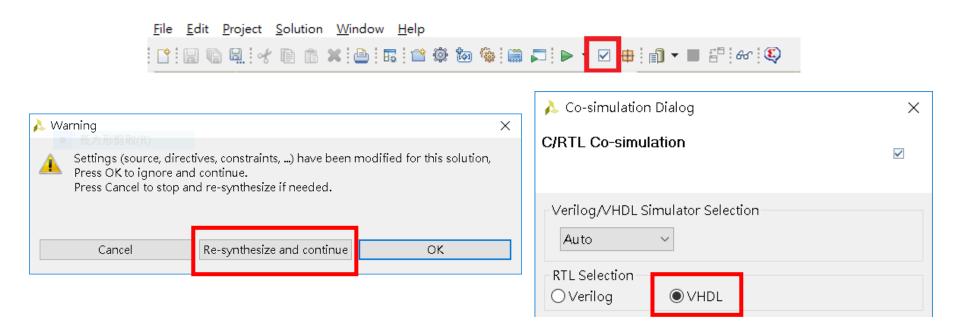
### (Optional) Step 1- C/RTL Cosimulation

- Click Project -> Project Setting -> Synthesis -> Top Function -> Browse
  - Change Top Function to "standalone\_mmult"



## (Optional) Step 2- C/RTL Cosimulation

- Run C/RTL Cosimulation
  - Choose "Re-synthesize and Continue" in the "Warning" pop-up window
  - Choose "VHDL" in the RTL Selection of the "Co-simulation Dialog" as follows)



## (Optional) Step 3- C/RTL Cosimulation

Check the console and wait until C/RTL cosimulation is finished

```
□ Console ⋈
            Frrors & Warnings
Vivado HLS Console
// 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%] @ "3625755000"
$finish called at time : 3625795 ns : File "D:/matmul/newnew/mmult/solution1/sim/verilog/standalone mmult.autotb.v" Line 319
run: Time (s): cpu = 00:00:00; elapsed = 00:00:26 . Memory (MB): peak = 224.191; gain = 0.000
## quit
INFO: [Common 17-206] Exiting xsim at Sat Apr 4 00:05:23 2020...
INFO: [COSIM 212-316] Starting C post checking ...
NORMAL MODE
Matrixes identical ... Test successful!
INFO: [COSIM 212-1000] *** C/RTL co-simulation finished: PASS ***
INFO: [COSIM 212-211] II 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 mal
Finished C/RTL cosimulation.
```

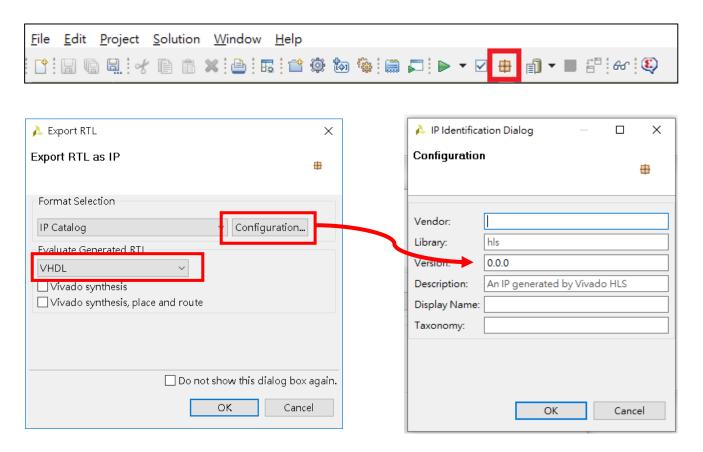
- Click Project -> Project Setting -> Synthesis -> Top Function -> Browse
  - Change Top Function back to "HLS\_accel" and Run C Synthesis again



### Step 9 - High-Level Synthesis (HLS)



- Export RTL
  - Choose "VHDL" in "Evaluate Generated RTL"
  - Click "Configuration", set version to "0.0.0"



### Step 9 - High-Level Synthesis (HLS)



- Export RTL
  - Wait until the export RTL is finished
  - Close Vivado HLS

```
Console S Perrors Warnings

Vivado HLS Console

source C:/Users/mcyang/Documents/mmult_hls/solution1/impl/ip/tmp.hw/webtalk/labtool_webtalk.tcl -notrace INFO: [Common 17-206] Exiting Webtalk at Sat Apr 04 11:15:53 2020...

INFO: [IP_Flow 19-234] Refreshing IP repositories

INFO: [IP_Flow 19-1704] No user IP repositories specified

INFO: [IP_Flow 19-2313] Loaded Vivado IP repository 'C:/Xilinx/Vivado/2016.3/data/ip'.

INFO: [Common 17-206] Exiting Vivado at Sat Apr 04 11:15:54 2020...

Finished export RTL.
```

#### **Expected Outcome**



# **Block Design in Vivado**

### Step 1 - Block Design

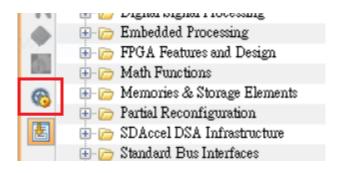


- Open Vivado 2016.3 and Create a new project
  - Don't need to create any files
- Click "IP Catalog" and then click "IP Settings"





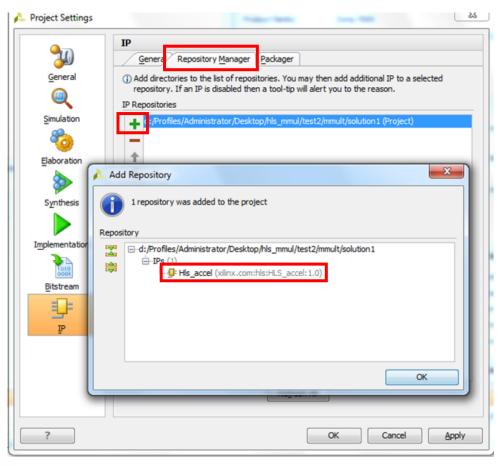




### Step 2 - Block Design



- At the "Repository Manager", click "+"
  - Add the repository of the HLS project
  - Then you can see the IP name "HLS\_accel"



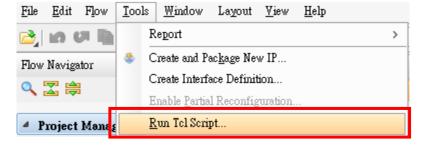
### Step 3 - Block Design

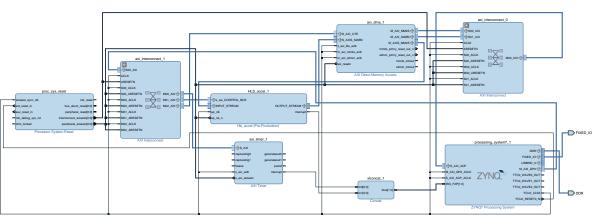


- Click Tools -> Run Tcl Script
  - Select system.tcl from the lab material
  - The block design is generated as follow

Click the Validation Design to ensure the block design

success



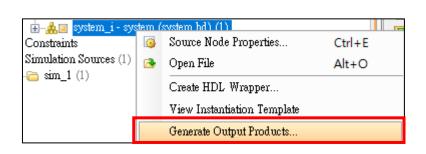


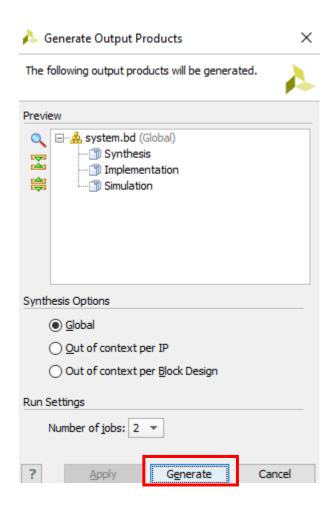


### Step 4 - Block Design



- Right click system.bd, click Generate Output Products
  - Click Generate



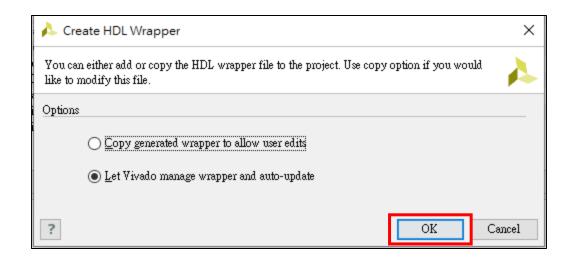


### Step 5 - Block Design



- Right click system.bd
  - Click Create HDL Wrapper and click OK

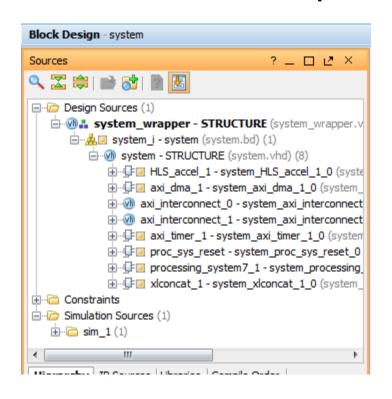


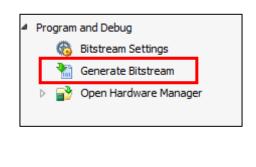


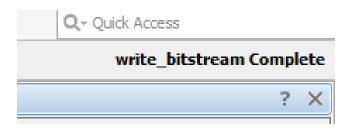
### Step 6 - Block Design



- Then you can see the Design Sources like this: system\_wrapper.vhd is generated
- Click Generate Bitstream and wait until the message "write\_bitstream complete" is shown



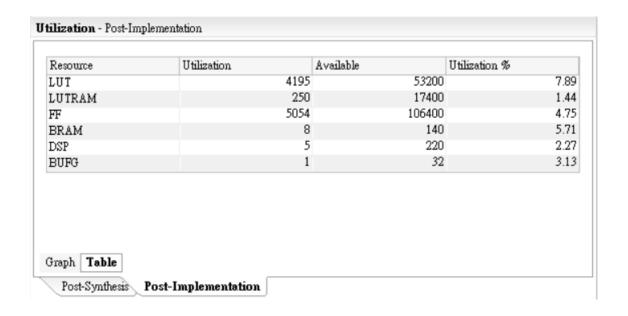




### Step 7 - Block Design



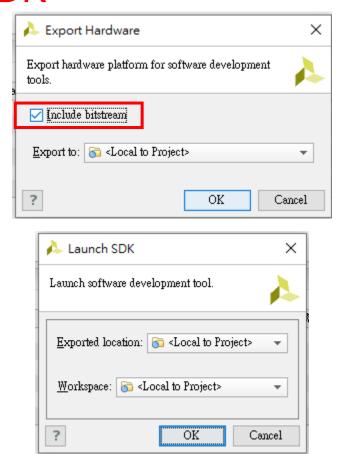
 Check Utilization – Post-Implementation of your design by clicking Window -> Project Summary



### Step 8 - Block Design



- Click File->Export Hardware (Click "Include bitstream")->OK
- Click Launch SDK





# Verifying HLS by using ARM in Vivado SDK

### **Step 1 - Verifying HLS by using ARM**



- After launch Vivado SDK, you have to check the existence of
  - HLS\_accel\_1
  - axi\_dma\_1
  - axi\_timer\_1

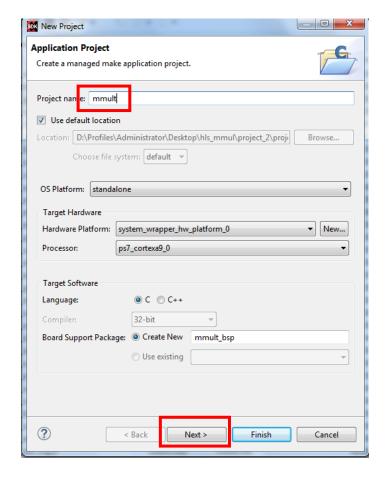
### Click "Cell" to sort the table by cell name

address Map for processor ps7_cortexa9_[0-1]				
Cell	Base Addr	High Addr	Slave I/f	Mem/Reg
HLS_accel_1	0x43c00000	0x43c0ffff	s_axi_CONTR	REGISTER
axi_dma_1	0x40400000	0x4040ffff	S_AXI_LITE	REGISTER
axi_timer_1	0x42800000	0x4280ffff	S_AXI	REGISTER
ps7_afi_0	0xf8008000	0xf8008fff		REGISTER
ps7_afi_1	0xf8009000	0xf8009fff		REGISTER
ps7_afi_2	0xf800a000	0xf800afff		REGISTER
ps7_afi_3	0xf800b000	0xf800bfff		REGISTER
ps7_coresight_comp_0	0xf8800000	0xf88fffff		REGISTER
ps7_ddr_0	0x00100000	0x1fffffff		MEMORY
ps7_ddrc_0	0xf8006000	0xf8006fff		REGISTER
ps7 dev cfa 0	0xf8007000	0xf80070ff		REGISTER

### Step 2 - Verifying HLS by using ARM



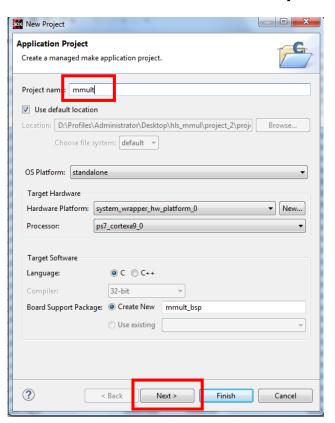
 Create a new application project by File->New->Application Project and input project name and click next:

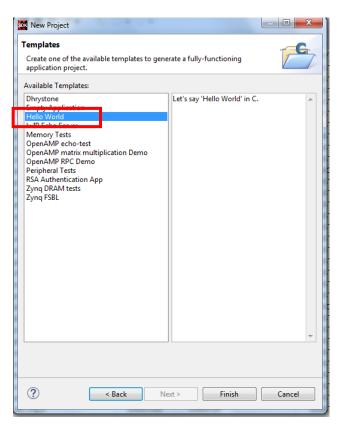


### Step 3 - Verifying HLS by using ARM



- Create a new application project by File->New->Application Project and input project name and click Next
  - Select Hello World Templates

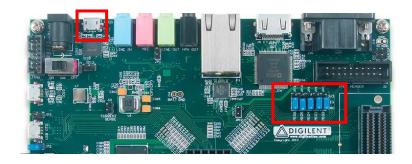




### Step 4 - Verifying HLS by using ARM



- Verify the boot mode are set as follows
- Connect the PROG port of Zed Board (J17) to a PC using the MicroUSB cable and Click Program FPGA

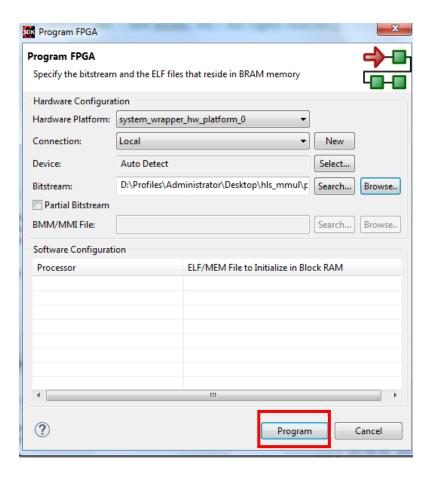




### Step 5 - Verifying HLS by using ARM



 Select bitstream file and click Program. After finish the program, you can see the green light is illumed.

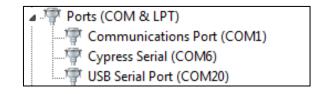


### Step 6 - Verifying HLS by using ARM



- Connect the USB-UART port of ZedBoard (J14) which is labeled UART to a PC using the MicroUSB cable
  - Check COM port at the Device Manager in Windows

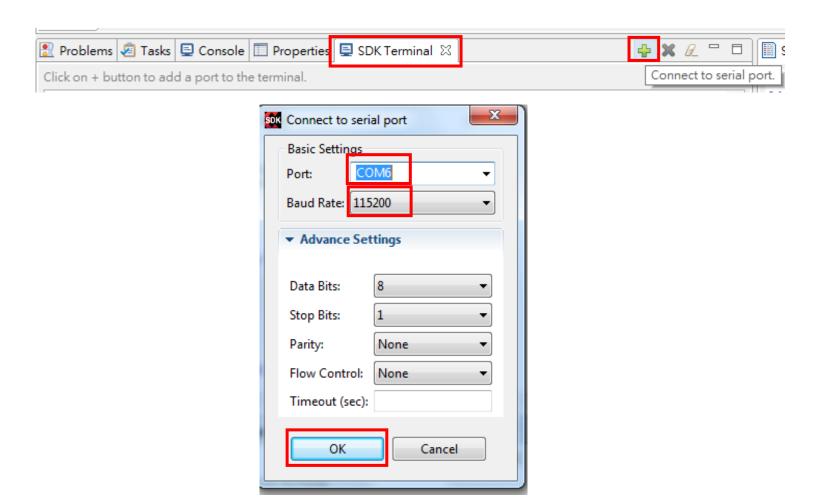




### Step 7 - Verifying HLS by using ARM



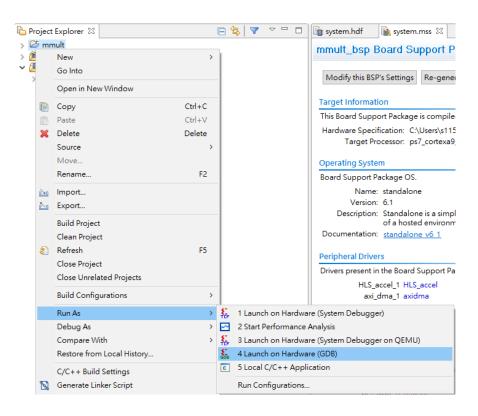
- Set series setting in SDK Terminal by green "+" icon
  - Select the right Port and Baud Rate, and click OK



### **Step 8 - Verifying HLS by using ARM**



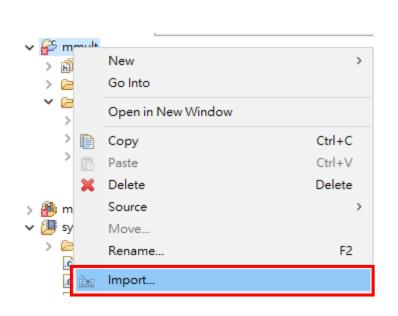
- Test COM port by right click "mmult"->Run As->Launch on Hardware (GDB). Then you can see Hello World in SDK Terminal
  - You should see the "Hello World" in SDK Terminal

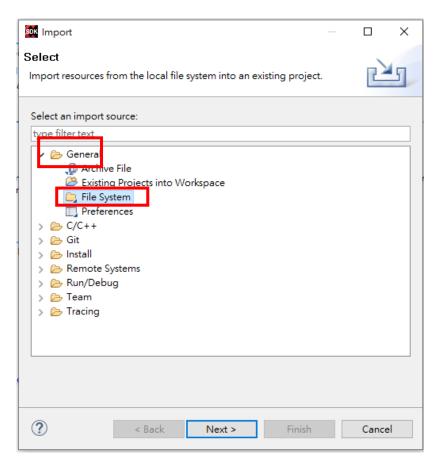


### **Step 9 - Verifying HLS by using ARM**



- Delete the helloworld.c file
- Right click mmult project, click "Import"
  - Select General -> File System

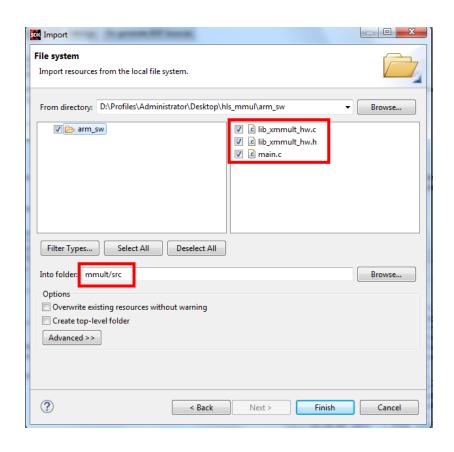




#### Step 10 - Verifying HLS by using ARM



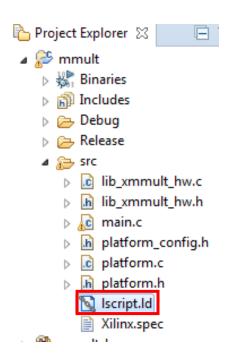
 Import the lib\_xmmult\_hw.h, lib\_xmmult\_hw.c, and main.c into mmult/src folder

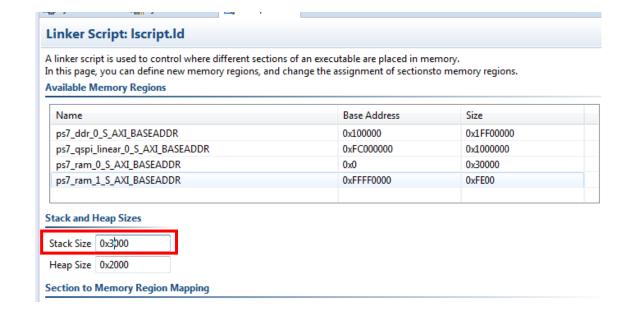


## Step 11 - Verifying HLS by using ARM



- In the Project Explorer mmult/src, double-click on the file Iscript.Id.
  - Change the Stack Size from 0x2000 to 0x3000. Save the file (Ctrl + S)

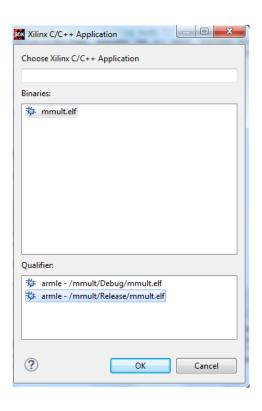




### Step 12 - Verifying HLS by using ARM



- Right click "mmult"
  - Select Build Configuration -> Build All
  - Select Run As -> Launch on Hardware (GDB) and select armlet-/mmult/Release/mmult.elf and click OK



#### Step 13 - Verifying HLS by using ARM



Finally, you can see the SDK Terminal

SDK Log 🚱 Terminal 1 🖂
Serial: (COM6, 115200, 8, 1, None, None - CONNECTED) - Encoding: (ISO-8859-1)
Hello World
*************
DMA Init done
Loop time for 1024 iterations is -2 cycles
Running Matrix Mult in SW
Total run time for SW on Processor is 25880 cycles over 1024 tests.
Cache cleared
Total run time for AXI DMA + HW accelerator is 333830 cycles over 1024 tests
Acce leration factor: 0.77
leration factor: 0.77
SW and HW results match!
Sw did tw resures materi.
*************
************
DMA Init done
Loop time for 1024 iterations is -2 cycles
Running Matrix Mult in SW
Total run time for SW on Processor is 25880 cycles over 1024 tests.
Cache cleared
Total run time for AXI DMA + HW accelerator is 333830 cycles over 1024 tests
Acce
leration factor: 0.77
SW and HW results match!
SW diff IN TESULES MACEIT.

#### Lab08: Your Task



- Task: Develop/optimize the floating point matrix multiplication with HLS
  - Basic: Complete the step-by-step instructions
  - Improve Parallelism:

You just need to modify the mmult\_hw() in mmult.h (From Line 44)

Step 1: Try to add pipeline directive to L1, L2 or L3 Uncomment the "#pragma HLS PIPELINE II=1" to let the L1, L2 or L3 pipeline.

Step 2: Read and analysis the synthesis report to find out what limits parallelism. (Check the "Loop" in the report to see whether the pipeline of the loop works well.)

Set the "#pragma HLS PIPELINE II=1" to L2 is a relatively good choice.

The ideal initiation interval should be 1, however there are some issues make the ideal initiation interval impractical. In the next step you will solve the bottlenecks.

#### Lab08: Your Task



Step 3: Add or modify the directives to further increase the parallelism. Parallelism means that multiple operations occur simultaneously, sometimes the bottleneck is the inability to read the data for calculation in time, sometimes the bottleneck is that the pipeline stops because of the in/out loop.

- 1) Uncomment the "#pragma HLS array\_partition variable=<ARRAY> block factor=<FACTOR> dim=<DIM>" then set the factor and dim to let the array be partitioned.
- 2) Remove the "#pragma HLS loop\_flatten off", which prevents the HLS tool from automatically merging outer loops of the "pipelining" into one larger loop.

#### Submission Rule:

- Submit your
  - 1) Revised mmult.h
  - 2) Screen shot of "Performance Estimates" for each HLS design (as shown in P.12)
  - 3) Screen shot of SDK verification for **each** HLS design (as shown in P.40)

In total: 3 + 3 screen shots

Deadline: 12:30 on 26 March 2025 (Late submission is NOT acceptable)

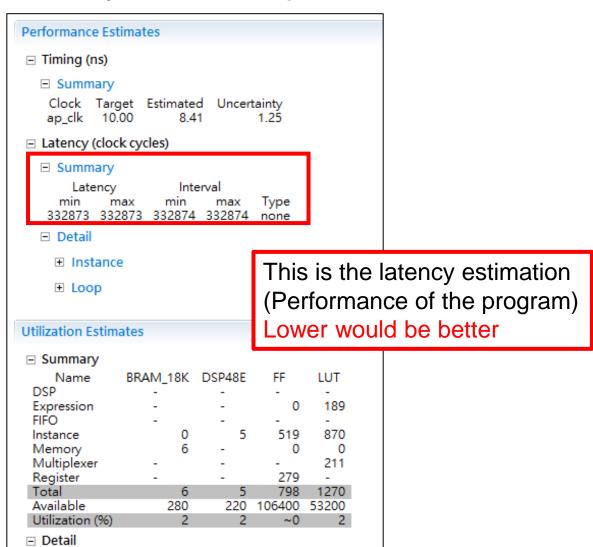
## Hints: Workflow (p.3)



- You just need to modify the mmult\_hw() in mmult.h
  (From Line 44)
- Run C Synthesis
- Until you got a decent result in the performance report, then you export RTL, move on to the block design and SDK verification steps

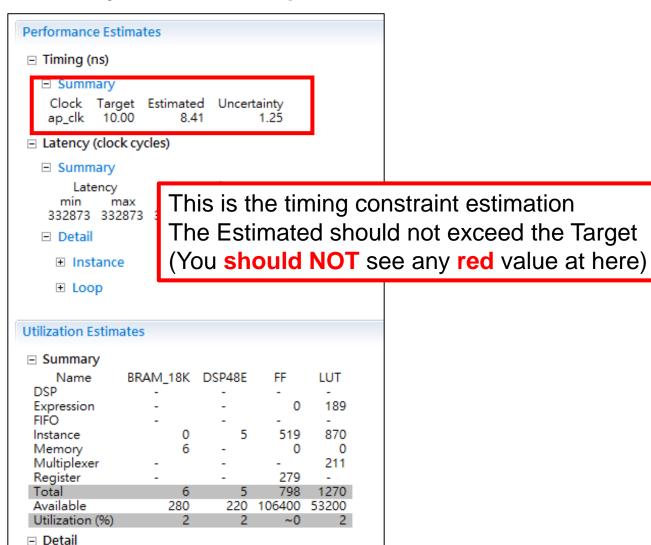
#### **Hints: Performance of the solution**





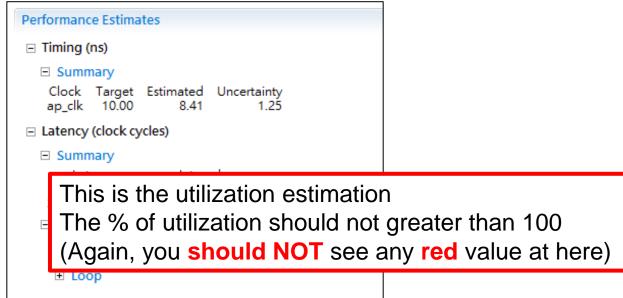
#### **Hints: Comply with time constraints**

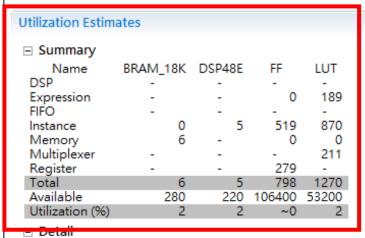




#### Hints: Use resources wisely



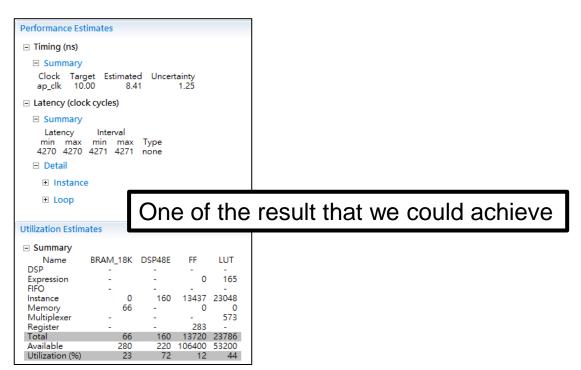




#### **Hints: Final performance**



- For Loop Optimization, you need to choose the which loop do you want to optimize? which pragma do you use?
- For Array Partition Optimization, you need to decide the parameters, such as dimension and factor.





## 香港中文大學

The Chinese University of Hong Kong

# Thank you

