# Report – Lab 2

R13943009 鄭至盛
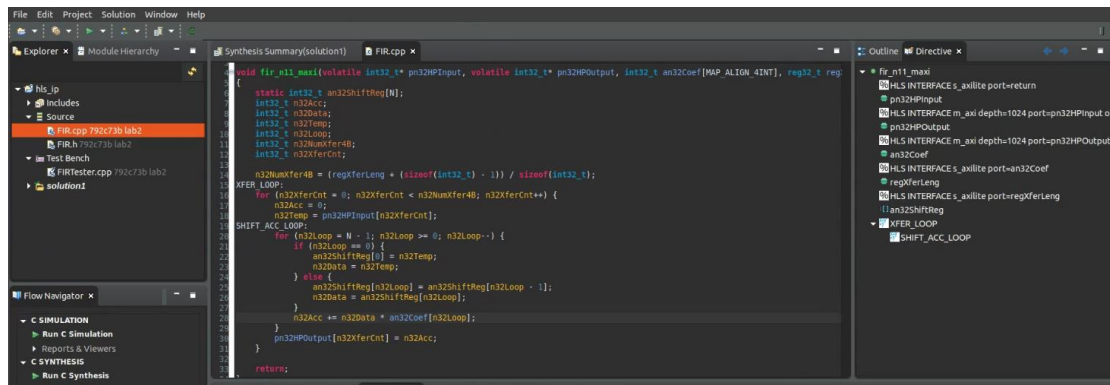
## Part 1 – AXI-Master

## Overview of the HLS Design:

### The C code

The .h file simply defines our top function (multip_2num) & data type (int32_t).
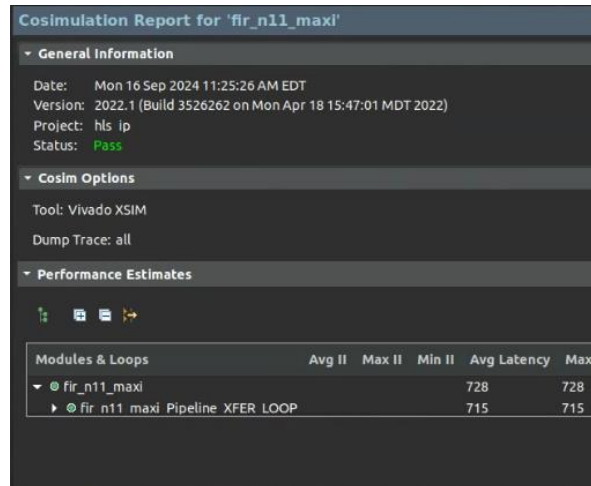
The .cpp file defines a shift register to store the temporary filtered value & the static for-loop that describes the behavior ($A = A + B \times C$) & boundary condition of the FIR filter. The directives are set up with a .tcl file, containing AXI-Lite for the configurations & AXI-Master for the I/O samples.
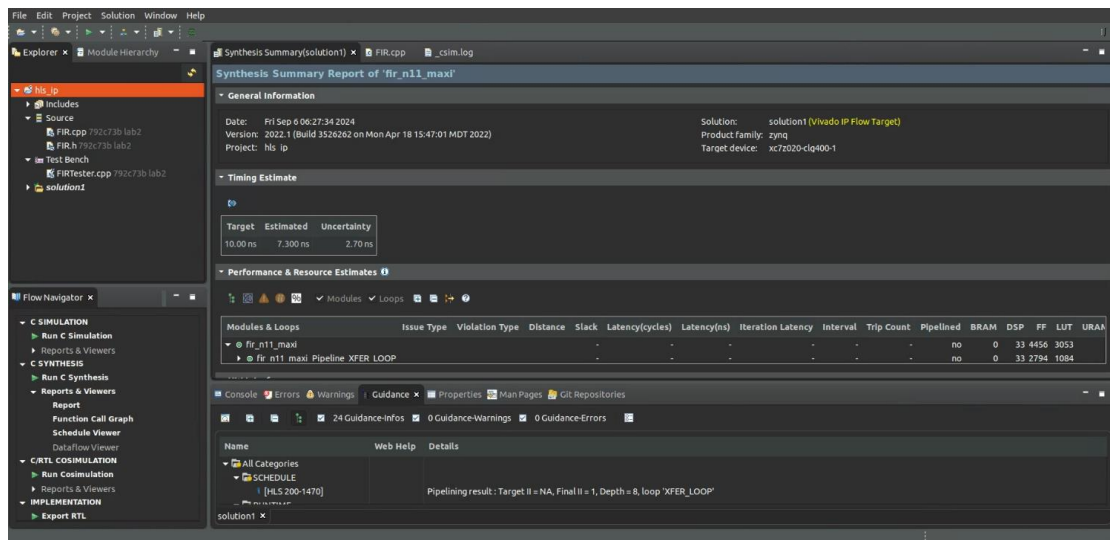


The tester file first set up the taps & I/O array for the HLS IP. Then, it would verify the NUM_SAMPLES of the filtered output & the golden data. By calling the diff command with system() API, we can determine whether the IP behaves correctly.

### C-simulation (left) & Co-simulation (right)

The co-sim error would occur if the DEPTH of the testbench & that of the directive is mismatched.

**Cosimulation Report for 'fir_n11_maxi'**

▾ **General Information**

Date: Mon 16 Sep 2024 11:25:26 AM EDT
Version: 2022.1 (Build 3526262 on Mon Apr 18 15:47:01 MDT 2022)
Project: hls ip
Status: Pass

▾ **Cosim Options**

Tool: Vivado XSIM

Dump Trace: all

▾ **Performance Estimates**

| Modules & Loops | Avg II | Max II | Min II | Avg Latency | Max |
|---|---|---|---|---|---|
| ▾ ⊙ fir_n11_maxi | | | | 728 | 728 |
| ▸ ⊙ fir n11 maxi Pipeline XFER LOOP | | | | 715 | 715 |

---

Synthesis Summary(solution1)    FIR.cpp    _csim.log ×

```
 1 INFO: [SIM 2] *************** CSIM start ***************
 2 INFO: [SIM 4] CSIM will launch GCC as the compiler.
 3   Compiling ../../../../hls_FIRN11MAXI/FIRTester.cpp in debug mode
 4   Generating csim.exe
 5 >> Start test!
 6 >> Comparing against output data...
 7 /home/ubuntu/course-lab_2/hls_ip/solution1/csim/build
 8 >> Test passed!
 9 --------------------
10 INFO: [SIM 1] CSim done with 0 errors.
11 INFO: [SIM 3] *************** CSIM finish ***************
12
```

## C-synthesis

The design is synthesized with clock period=10ns

---

File  Edit  Project  Solution  Window  Help

**Synthesis Summary Report of 'fir_n11_maxi'**

▾ **General Information**

Date: Fri Sep 6 06:27:34 2024
Version: 2022.1 (Build 3526262 on Mon Apr 18 15:47:01 MDT 2022)
Project: hls ip

Solution: solution1 (Vivado IP Flow Target)
Product family: zynq
Target device: xc7z020-clg400-1

▾ **Timing Estimate**

| Target | Estimated | Uncertainty |
|---|---|---|
| 10.00 ns | 7.300 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 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ▾ ⊙ fir_n11_maxi | | | | | - | - | - | - | | no | 0 | 33 | 4456 | 3053 | |
| ▸ ⊙ fir n11 maxi Pipeline XFER LOOP | | | - | - | - | - | - | - | - | no | 0 | 33 | 2794 | 1084 | |

▾ Console  Errors  Warnings  Guidance ×  Properties  Man Pages  Git Repositories

24 Guidance-Infos  0 Guidance-Warnings  0 Guidance-Errors

| Name | Web Help | Details |
|---|---|---|
| ▾ All Categories | | |
| ▾ SCHEDULE | | |
| [HLS 200-1470] | | Pipelining result : Target II = NA, Final II = 1, Depth = 8, loop 'XFER_LOOP' |

solution1 ×

## Generated Waveform

The aqua/purple part is the input/output signal of the RTL part.

The aqua/purple part & the marker reveal the data/address signal of the RTL part,
including addresses of I/O buffers, the coefficient of each tap, and the number of
samples.

# Overview of the Hardware Wrapper:

## Block Diagram

In the block diagram, we can see 2 AXI peripherals are connected between the generated HLS IP (fir_n11_maxi_0) & ZYNQ7 Processing System. With the protocol of AXI-Lite, the configuration such as the tap coefficient & buffer address can be sent to HLS IP. On the other hand, the AXI-Master protocol can handle the transfer of data before/after filtering.



## Verilog Code

The generated Verilog code can be divided into 2 parts. The computation unit can be further divided into a pipelined loop (containing pipelined MAC structure) & multiplier unit. On the other hand, the controlling unit deals with the AXI-Lite & AXI-Master protocol respectively.

## Hardware Utilization

From the aspect of DSP, the 11-tap multipliers require 33 DSP48E1. On the other hand, the HLS IP/AXI-Lite/AXI-Master requires ~3000/~400/~600 FFs respectively.



Under 100MHz clock signals, both setup/hold slack are met.



In the driver files of MISC, we can find the address of different control signals which PYNQ APIs can further exploit.

# Overview of the PS-side (PYNQ)

## Python code

We can see the Python code simply loads the bitstream file, selects IP, set up configurations, and writes & reads until there's no output according to the address defined in the driver files.

```python
from __future__ import print_function

import sys, os
import numpy as np
from time import time
import matplotlib.pyplot as plt

sys.path.append('/home/xilinx')
os.environ['XILINX_XRT'] = '/usr'
from pynq import Overlay
from pynq import allocate

if __name__ == "__main__":
    print("Entry:", sys.argv[0])
    print("System argument(s):", len(sys.argv))

    print("Start of \"" + sys.argv[0] + "\"")

    ol = Overlay("/home/xilinx/jupyter_notebooks/FIRN11MAXI.bit")
    ipFIRN11 = ol.fir_n11_maxi_0

    fiSamples = open("samples_triangular_wave.txt", "r+")
    numSamples = 0
    line = fiSamples.readline()
    while line:
        numSamples = numSamples + 1
        line = fiSamples.readline()

    inBuffer0 = allocate(shape=(numSamples,), dtype=np.int32)
    outBuffer0 = allocate(shape=(numSamples,), dtype=np.int32)
    fiSamples.seek(0)
    for i in range(numSamples):
        line = fiSamples.readline()
        inBuffer0[i] = int(line)
    fiSamples.close()

    numTaps = 11
    n32Taps = [0, -10, -9, 23, 56, 63, 56, 23, -9, -10, 0]
    #n32Taps = [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]
    n32DCGain = 0
    timeKernelStart = time()
    for i in range(numTaps):
        n32DCGain = n32DCGain + n32Taps[i]
        ipFIRN11.write(0x40 + i * 4, n32Taps[i])
    if n32DCGain < 0:
        n32DCGain = 0 - n32DCGain
    ipFIRN11.write(0x28, len(inBuffer0) * 4)
    ipFIRN11.write(0x10, inBuffer0.device_address)
    ipFIRN11.write(0x1C, outBuffer0.device_address)
    ipFIRN11.write(0x00, 0x01)
    while (ipFIRN11.read(0x00) & 0x4) == 0x0:
        continue
    timeKernelEnd = time()
    print("Kernel execution time: " + str(timeKernelEnd - timeKernelStart) + " s")

    plt.title("FIR Response")
    plt.xlabel("Sample Point")
    plt.ylabel("Magnitude")
    xSeq = range(len(inBuffer0))
    if n32DCGain == 0:
        plt.plot(xSeq, inBuffer0, 'b.', xSeq, outBuffer0, 'r.')
    else:
        plt.plot(xSeq, inBuffer0, 'b.', xSeq, outBuffer0 / n32DCGain, 'r.')
    plt.grid(True)
    plt.show() # In Jupyter, press Tab + Shift keys to show plot then redo run

    print("============================")
    print("Exit process")
```
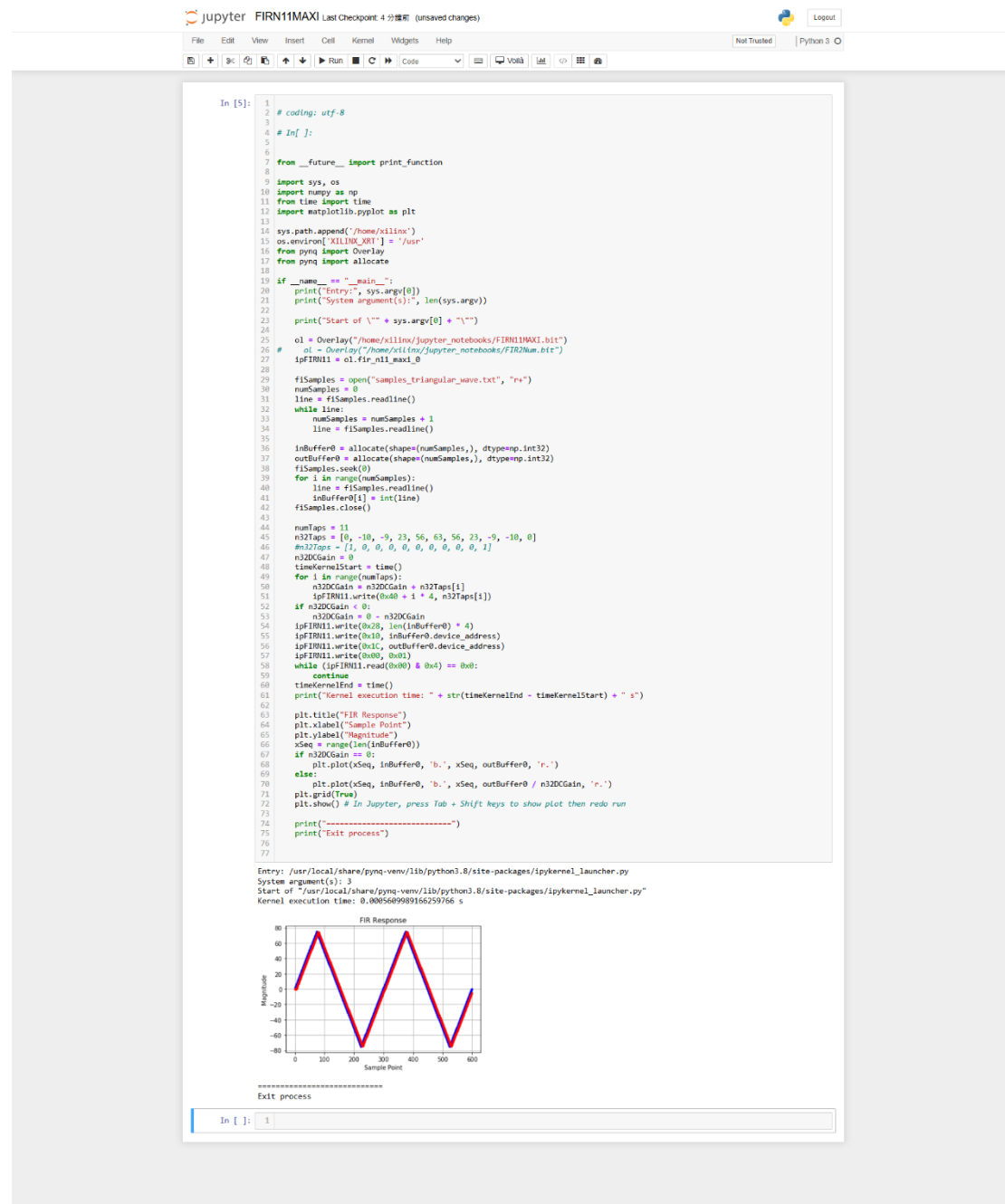
## Final results

As shown in the notebook, the computed results from HLS IP are correct & it requires 5.6e-4 seconds for 600 samples.

# Part 2 – AXI-Stream

## Overview of the HLS Design:

### The C code
Compared to the last examples, the code simply replaces the AXI-Master interface with AXI-Stream one. We also need to add read/write functions into .cpp file.





For the tester file, we need to add control of last signals & manually pull up

### C-simulation (left) & Co-simulation (right)

## C-synthesis

The design is synthesized with clock period=10ns



## Generated Waveform

We can see that the output is terminated when the TLAST signals is pulled up.

# Overview of the Hardware Wrapper:

## Block Diagram

In the block diagram, unlike the last example, we need to manually add DMA IP for both input/output of HLS IP. The DMA IP facilitates high-bandwidth data transfers without involving the CPU, while AXI peripherals typically only handle memory-mapped transactions.



## Verilog Code

The generated Verilog code can be divided into 2 parts. The computation unit can be further divided into a pipelined loop (containing pipelined MAC structure) & multiplier unit. On the other hand, the controlling unit deals with the AXI-Lite & AXI-Stream protocols respectively.

## Hardware Utilization

From the aspect of DSP, the 11-tap multipliers also require 33 DSP48E1. On the other hand, the HLS IP/AXI-Lite/AXI-Master requires ~1400/~600/~1100 FFs respectively. The DMA in/out IP also introduces ~750/~1400 additional FFs.



Under 100MHz clock signals, both setup/hold slack are met.



In the driver files of MISC, we can find the address of different control signals which PYNQ APIs can further exploit.

# Overview of the PS-side (PYNQ)

## Python code

We can see the Python code simply loads the bitstream file, selects IP, and writes to & reads from the address defined in the driver files. In addition, the ipDMAOut needs to be manually modified to "ol.axi_dma_out_0". We can also see that the order of transfer & start signals is different between AXI-Master & AXI-Stream.

```python
from __future__ import print_function

import sys, os
import numpy as np
from time import time
import matplotlib.pyplot as plt

sys.path.append('/home/xilinx')
os.environ['XILINX_XRT'] = '/usr'
from pynq import Overlay
from pynq import allocate

if __name__ == "__main__":
    print("Entry:", sys.argv[0])
    print("System argument(s):", len(sys.argv))

    print("Start of \"" + sys.argv[0] + "\"")

    ol = Overlay("/home/xilinx/jupyter_notebooks/FIRN11Stream.bit")
    ipFIRN11 = ol.fir_n11_strm_0
    ipDMAIn = ol.axi_dma_in_0
    ipDMAOut = ol.axi_dma_out_1

    fiSamples = open("samples_triangular_wave.txt", "r+")
    numSamples = 0
    line = fiSamples.readline()
    while line:
        numSamples = numSamples + 1
        line = fiSamples.readline()

    inBuffer0 = allocate(shape=(numSamples,), dtype=np.int32)
    outBuffer0 = allocate(shape=(numSamples,), dtype=np.int32)
    fiSamples.seek(0)
    for i in range(numSamples):
        line = fiSamples.readline()
        inBuffer0[i] = int(line)
    fiSamples.close()

    numTaps = 11
    n32Taps = [0, -10, -9, 23, 56, 63, 56, 23, -9, -10, 0]
    #n32Taps = [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]
    n32DCGain = 0
    timeKernelStart = time()
    for i in range(numTaps):
        n32DCGain = n32DCGain + n32Taps[i]
        ipFIRN11.write(0x40 + i * 4, n32Taps[i])
    if n32DCGain < 0:
        n32DCGain = 0 - n32DCGain
    ipFIRN11.write(0x10, len(inBuffer0) * 4)
    ipFIRN11.write(0x00, 0x01)
    ipDMAIn.sendchannel.transfer(inBuffer0)
    ipDMAOut.recvchannel.transfer(outBuffer0)
    ipDMAIn.sendchannel.wait()
    ipDMAOut.recvchannel.wait()
    timeKernelEnd = time()
    print("Kernel execution time: " + str(timeKernelEnd - timeKernelStart) + " s")

    plt.title("FIR Response")
    plt.xlabel("Sample Point")
    plt.ylabel("Magnitude")
    xSeq = range(len(inBuffer0))
    if n32DCGain == 0:
        plt.plot(xSeq, inBuffer0, 'b.', xSeq, outBuffer0, 'r.')
    else:
        plt.plot(xSeq, inBuffer0, 'b.', xSeq, outBuffer0 / n32DCGain, 'r.')
    plt.grid(True)
    plt.show() # In Jupyter, press Tab + Shift keys to show plot then redo run

    print("==============================")
    print("Exit process")
```

## Final results

As shown in the notebook, the computed results from HLS IP are correct & it requires 1.6e-3 seconds for 600 samples.

In [2]:

```
# coding: utf-8

# In[3]:


from __future__ import print_function

import sys, os
import numpy as np
from time import time
import matplotlib.pyplot as plt

sys.path.append('/home/xilinx')
os.environ['XILINX_XRT'] = '/usr'
from pynq import Overlay
from pynq import allocate

if __name__ == "__main__":
    print("Entry:", sys.argv[0])
    print("System argument(s):", len(sys.argv))

    print("Start of \"" + sys.argv[0] + "\"")

    ol = Overlay("/home/xilinx/jupyter_notebooks/FIRN11Stream.bit")
    ipFIRN11 = ol.fir_n11_strm_0
    ipDMAIn = ol.axi_dma_in_0
    ipDMAOut = ol.axi_dma_out_0

    fiSamples = open("samples_triangular_wave.txt", "r+")
    numSamples = 0
    line = fiSamples.readline()
    while line:
        numSamples = numSamples + 1
        line = fiSamples.readline()

    inBuffer0 = allocate(shape=(numSamples,), dtype=np.int32)
    outBuffer0 = allocate(shape=(numSamples,), dtype=np.int32)
    fiSamples.seek(0)
    for i in range(numSamples):
        line = fiSamples.readline()
        inBuffer0[i] = int(line)
    fiSamples.close()

    numTaps = 11
    n32Taps = [0, -10, -9, 23, 56, 63, 56, 23, -9, -10, 0]
    #n32Taps = [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]
    n32DCGain = 0
    timeKernelStart = time()
    for i in range(numTaps):
        n32DCGain = n32DCGain + n32Taps[i]
        ipFIRN11.write(0x40 + i * 4, n32Taps[i])
    if n32DCGain < 0:
        n32DCGain = 0 - n32DCGain
    ipFIRN11.write(0x10, len(inBuffer0) * 4)
    ipFIRN11.write(0x00, 0x01)
    ipDMAIn.sendchannel.transfer(inBuffer0)
    ipDMAOut.recvchannel.transfer(outBuffer0)
    ipDMAIn.sendchannel.wait()
    ipDMAOut.recvchannel.wait()
    timeKernelEnd = time()
    print("Kernel execution time: " + str(timeKernelEnd - timeKernelStart) + " s")

    plt.title("FIR Response")
    plt.xlabel("Sample Point")
    plt.ylabel("Magnitude")
    xSeq = range(len(inBuffer0))
    if n32DCGain == 0:
        plt.plot(xSeq, inBuffer0, 'b.', xSeq, outBuffer0, 'r.')
    else:
        plt.plot(xSeq, inBuffer0, 'b.', xSeq, outBuffer0 / n32DCGain, 'r.')
    plt.grid(True)
    plt.show() # In Jupyter, press Tab + Shift keys to show plot then redo run

    print("=============================")
    print("Exit process")
```
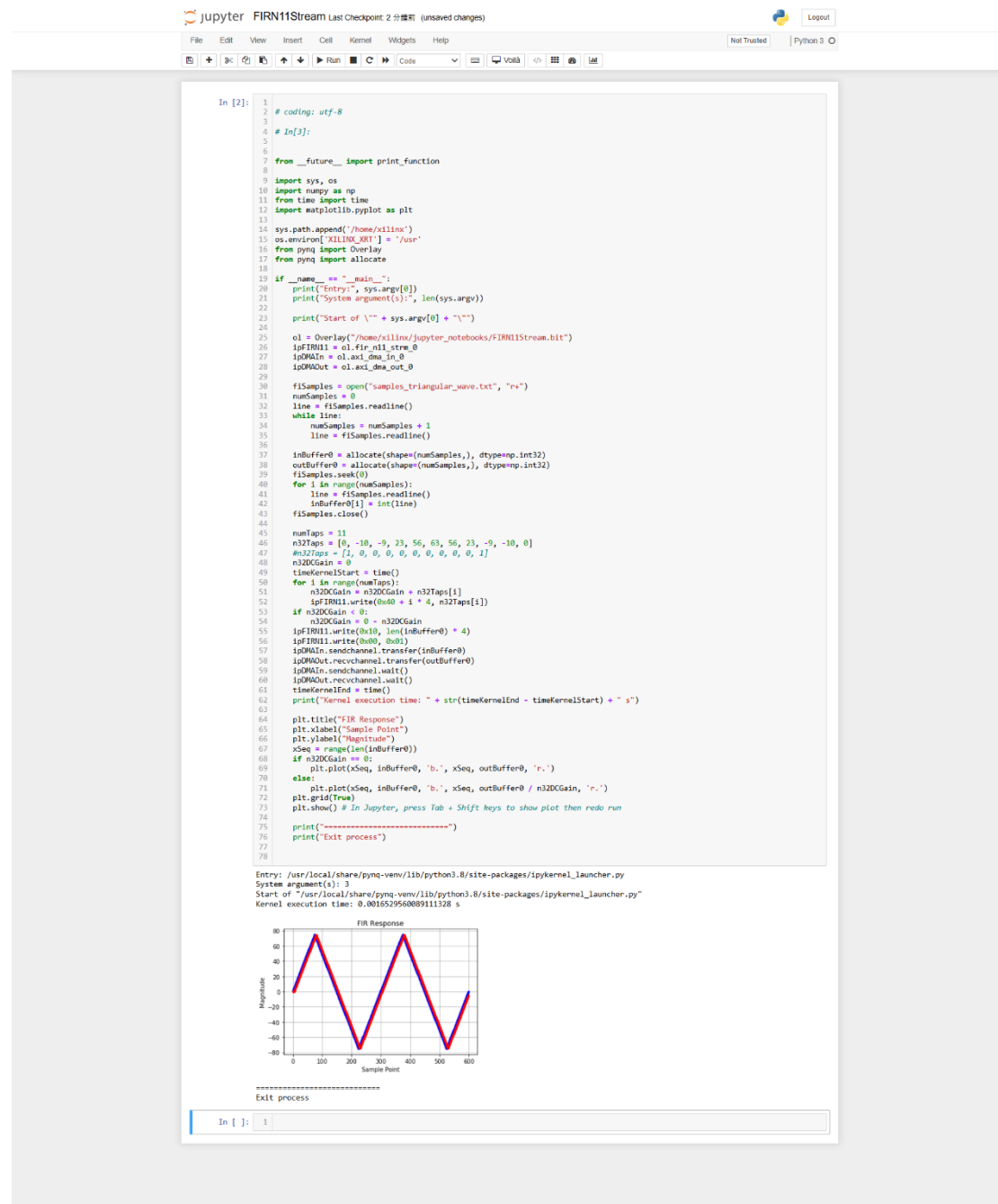
```
Entry: /usr/local/share/pynq-venv/lib/python3.8/site-packages/ipykernel_launcher.py
System argument(s): 3
Start of "/usr/local/share/pynq-venv/lib/python3.8/site-packages/ipykernel_launcher.py"
Kernel execution time: 0.0016529560089111328 s
```

FIR Response

```
=============================
Exit process
```

In [ ]:

# Summary

The 2 examples of lab 2 simply walk us through 2 different AXI protocols: AXI-Master & AXI-Stream. We can also discover some differences between the behavior & performance of them.

First, we can discover that the simulation time of AXI-Stream is much longer than that of AXI-Master from the co-simulation result. It might be caused by the additional control signals in the input for-loop of the .cpp tester.

Second, we can discover that the simulation time of AXI-Stream is also still longer than that of AXI-Master from the PYNQ on-board test. Although I cannot fully investigate the black-box behavior inside both PYNQ API & the board, when using DMA with AXI-Stream interfaces, there may be additional software overhead involved in managing the data flow compared to using an AXI-Master interface directly with memory-mapped operations. This overhead might further slow down performance when the amount of data is low.

In summary, while AXI-Stream is suitable for specific applications like video or audio streaming where continuous data flow is essential, it may not match the efficiency of AXI-Master interfaces in scenarios requiring high throughput and low latency due to its inherent design limitations and the complexities involved in managing backpressure and flow control.