Xilinx & Vivado IDE

Speaker: Jia-Ming Lin

Outlines

- Recap. PYNQ Framework
- FPGA Architecture
- FPGA Design Process
- Design Optimization
- Lab1:
 - Creating hardware design by using Vivado HLS
 - Generating Bitstream
 - Validating on FPGA.

Recap. PYNQ Framework



- Interactive computing environment based on Jupyter Notebook
- Using Python to invoke hardware libraries and overlay in Programmable Logic
- Through PYNQ, we can easier to create high performance applications with
 - Parallel hardware execution
 - Hardware acceleration algorithm
 - High frame-rate video processing.

What is an Overlay?

- Hardware library, consisted of one or several IPs
- Extends user application from Processing System(CPU) to Processing Logic(FPGA)
 - Speed-up or customize the hardware platform for a particular application.
- Example: PYNQ base overlay



Case Study: Video input/output

- Two Methods
 - a. OpenCV on ARM A9s
 - b. HDMI_IN and HDMI_OUT on Processing Logic
- Which one is better? In terms of throughput(FPS), energy consumption(Watt)

Method (a)

- USB Webcam on USB port of Zynq PS(CPU)
- Capture image and stored to DRAM
- OpenCV package is used

Source code reference: https://gist.github.com/cathalmccabe/b0ab8 917f748840f0d3959f7eabf0f82



Method (b)

FPGA Architecture

Field Programmable Gate Arrays(FPGAs)

• Recap, Programmable Logic Device

- Contains an array of AND gate & another array of OR gate
- Users can configure these devices in order to implement the Boolean functions
- E.x. some form of sum-of-product.
- Idea behind FPGA,
 - Programmable/configurable logic blocks,
 - Programmable interconnections between logic blocks
 - Configure how the logic blocks are connected
- Two dominant FPGA makers
 - Xilinx and Altera



Configurable Logic Block(CLB) or Slice

• Based on Look-up Tables(LUTs), Flip-Flop and other specialized functions(e.g. Full Adder)





Programmable Interconnections

• Flexible network connections of CLBs(or Slices)



- At each interconnect, there is a switch transistor which is default OFF
- Each switch is controlled by a 1-bit register
- Configuring routing is simply put 0 or 1 into the register

Modern FPGAs



- As the number of transistors on FPGA increasing...
- More "hard" resources
- DSP48
 - Mul, Add, MAC
 - Higher frequence
 - $\circ \quad \text{Not flexible as CLB}$
- Microprocessors
 - ARM or x86
- Block RAM (BRAM)
 - Transfer data
 - Store large data on-chip 11

Comparison of three different on- and off-chip memory storage options

	External		
	Memory	BRAM	\mathbf{FFs}
count	1-4	thousands	millions
size	GBytes	KBytes	Bits
total size	GBytes	MBytes	100s of KBytes
width	8-64	1-16	1
total bandwidth	GBytes/sec	TBytes/sec	100s of TBytes/sec

- **External Memory**: highest density, lowest bandwidth
- **FFs**: highest total bandwidth, limited amount of data storage capability.
- **BRAM**: intermediate value between external memory and FFs

PYNQ-Z2



FPGA Design Process

Block diagram showing a hypothetical embedded FPGA design



- I/O Interface Core: highly constrained by timing, RTL implemented
- Standard Core: Generic, fixed-function processing,
 - Processor core, on-chip memory, interconnections
- Accelerator core:
 - Where we focus on and develop by using HLS.

Design Optimization

Benchmarking Metrics for Hardware Design

- Accuracy
 - Quality of result for a given task
- Throughput
 - Analytics on high volume data
 - Real-time performance (e.g., video at 30 fps)
- Latency
 - For interactive applications (e.g., autonomous navigation)
- Energy and Power
 - Edge and embedded devices have limited battery capacity
 - Data centers have stringent power ceilings due to cooling costs
- Hardware Cost

Reference: <u>https://www.rle.mit.edu/eems/wp-content/uploads/2019/09/2019_icip_tutorial.pdf</u>

Definitions

- Task
 - Function invocation, e.g. int32_t func(int32_t a)
- Task latency
 - Time between task start and when it finishes
- Task interval
 - Time between one task starts and the next starts



- Horizontal Axis: time
- Vertical Axis: Functional units in design
 - E.g. Mul(*), Add(+)
- Four executions of a design, starting a new task every cycle

Area/Throughput Tradeoffs

- Example: FIR filter
 - 1-Dim convolution
 - Key question: what circuit is generated from this code?

#define NUM_TAPS 4

```
void fir(int input, int *output, int taps[NUM_TAPS])
{
    static int delay_line[NUM_TAPS] = {};
```

int result = 0;

for (int $i = NUM_TAPS - 1$; i > 0; i--) { delay_line[i] = delay_line[i - 1]; } delay_line[0] = input;

for (int i = 0; $i < NUM_TAPS$; i++) { result += delay_line[i] * taps[i];

Inner-product



*output = result;

Area/Throughput Tradeoffs

• Architecture generated from previous slide



Design 1

- Advantage: Less resources on multiplication and add operations
- **Disadvantage**: Higher latency
- Disadvantage: more resources on control logics

Area/Throughput Tradeoffs

• With little modifications on HLS(Next week)



Design 2

- Advantage: Higher processing rate, less latency
- Advantage: Simpler architecture
- Disadvantage: Need more resources on Mul and Add

In this course

Chapter		€ CORDIC	LHC 4	G FFT	o SPMV	→ MatMul	∞ Histogram	6 Video	10 Sorting	1 Huffman
Loop Unrolling	X		X	X	Х		Х		Х	
Loop Pipelining	X		X	X	X		X	X	Х	X
Bitwidth Optimization	X	X								X
Function Inlining	X									X
Hierarchy	X			X			X	X	Х	X
Array Optimizations			X	X	X	X	X	X	Х	X
Task Pipelining				X			X	X	Х	X
Testbench					X	X			Х	X
Co-simulation					X					
Streaming						X		X	Х	
Interfacing								X		

Lab 1: Overlay Tutorial

Overview

- Developing a Single IP
 - Intellectual Property(IP) is a reusable IC design that is the intellectual property of one party.
 - Using Vivado HLS
- Generating Bitstream and download to FPGA
- Interact with our own IP by using PYNQ

Open the Vivado HLS 2020.1 IDE



Short path on desktop

Create the Project







Create the Project

Ne	w Vivado HLS Project	0 0
Solution Configuration Create Vivado HLS solution for s	elected technology	E
Solution Name: solution1		
Clock Period: 10	Uncertainty:	
Part Selection Part: xc7vx485tffg1157-1		
O Vitis Bottom Up Flow		7. Select
		Device
< Back		Cancel Finish

chool. Au			
isplay Name: 📶			~
Select P		Family	Vendor
pyng-z2 J. OCICULI	xczu28dr-ffva1517-2-e	zyng	tul.com.tw
Zyng UltraScale+ ZCU106 Evaluation Platform	xczu7ev-ffvc1156-2-e	zynquplus	xilinx.com
Zynq UltraScale+ ZCU104 Evaluation Board xczu7ev-ffvc1156-2-e zynquplus xi			xilinx.com
		Second Second	will may see my
Zynq UltraScale+ ZCU102 Evaluation Board	xczu9eg-ffvb1156-2-e	zynquplus	xiunx.com

Create Project

New Vivado HLS Project	ct 🛛 🕲 🕲
Solution Configuration	
Create Vivado HLS solution for selected technology	
Solution Name: solution1	
Clock	
Period: 10 Uncertainty:	
Part Selection	
Part pynq-z2 (xc7z020clg400-1)	
^{□ viti} f fot Make sure the selected device is PYNQ-Z2	
	12. Click if
	everything ready
< Back	Cancel Finish



Create New Files



Three new files...

Under Source,

 (1) top.h: Define datatype or function, ...
 (2) top.cpp: Hardware design implementation

 Under Test Bench,

(1) **testbench.cpp**: Check hardware/software result is matched.

top.h

#include <iostream>
#include <ap_int.h>

using namespace std;

// 32 bit integer
typedef ap_int<32> data_t;

void add(data_t a, data_t b, data_t &c);

top.cpp

#include	e "to	op.h"
void add	d(dat	ta t a, data t b, data t& c) {
#pragma	HLS	INTERFACE ap ctrl none port=return
#pragma	HLS	INTERFACE s axilite port=a
#pragma	HLS	INTERFACE s axilite port=b
#pragma	HLS	INTERFACE s_axilite port=c
c =	a +	b;

testbench.cpp

#ine	clude "top.h"
int	<pre>main(){ data_t a = 5; data_t b = 6; data_t c;</pre>
	<pre>add(a,b,c); int err = 0;</pre>
	<pre>if (c != (a+b)){ err+=1; cout << "Error, HW c = " << c << endl:</pre>
}	} return err;

- Toolbar: 🗖 ► 🗸 🖶
 - **C Simulation**: Check the design output is same with software implementation
 - ▶ C Synthesis:
 - Generate RTL design(e.g. Verilog) from HLS C code
 - Performance and Resource consumption report
 - C/RTL Cosimulation: Check RTL output is matched with software implementation.
 - **Export RTL**: wrap the design to an IP module
- In this Lab, click "C Simulation" → "C Synthesis" → "Export RTL"

Open Vivado



Create Vivado Project



Create Vivado Project



Create Vivado Project



• Import HLS IP



• Import HLS IP



• Constructing Block Design

The Created Block Design View





• Constructing Block Design



• Generate Bitstream





• Generate Bitstream



• Generate Bitstream

PLOCK DESIGN design 1	
BEOCK DESIGN - design_1	
Sour × Design Signals Board ? _ 🗆 🗅	Address Editor x Diagram x
Q ≚ ≑ + 2 ● 0 ♦	$\textcircled{Q} \hspace{0.1in} Q \hspace{0.1in} \begin{matrix} X \\ X \\ \end{matrix} \\ \begin{matrix} X \\ \end{pmatrix} \hspace{0.1in} \bigcirc \hspace{0.1in} Q \\ \begin{matrix} X \\ X \\ \end{matrix} \\ \begin{matrix} X \\ Y \\ \end{matrix} \\ \begin{matrix} X \\ Y \\ \end{matrix} \\ \begin{matrix} X \\ Y \\ \end{matrix} \\ \begin{matrix} Y \\ Y \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\ \end{matrix} \\ \end{matrix} \\ \end{matrix} \\ \begin{matrix} Z \\ Z \\ \end{matrix} \\$
	add_0 # 4.00_AXI # 9.7k. Bitstream Generation Completed Bitstream Generation successfully completed.
A design_1.bd Cation: /home/jiaming/workplace/	Next © Qpen Implemented Design Wew Reports Open Hardware Manager mb debug systst interconnet presel/001
General Properties	O Generate Memory Configuration File Processor System Reset
Tcl Console × Messages Log Report Q ★ ♦ II ■ ■ ■	S OK Cancel
Exporting to file /home/jiaming/workplace Generated Block Design Tcl file /home/jia Generated Hardware Definition File /home/ [Wed Mar 3 23:38:30 2021] Launched desig Run output will be captured here: design_1_rst_ps7_0_100M 1_synth_1: /home/ design 1_add 0.2_synth 1: /home/iaming/w	/tutorial/add ip_design.add ip_design.srcs/sources_1/bd/design_1/hw_handoff/design_1.hwh monkpthee/tutoria/add_fp_design.arcs/sources_1/bd/design_1/by/http://tutorial/add_fp_design_1.hwh n_1_rst_ps70_100M_1_synth_1, design_1_add_0_2_synth_1, design_1_processing_system70_1_synth_1, design_1_auto_pc_0_synth_1, synth_1. jiaming/workplace/tutorial/add_ip_design/add_ip_design.runs/design_1_rst_ps7_0_100M_1_synth_1/runme.log problace/tutorial/add_ip_design/add_ip_design.runs/design_1_adt_0_xouth_1/runme.log

- Two Design files
 - <Your Design Path>/<Your Project Name>.srcs
 /add_ip_design.srcs/sources_1/bd/design_1/hw_handoff/design_1.hwh
 - o <Your Design Path>/<Your Project Name>.runs/impl_1/design_1_wrapper.bit
- Power on PYNQ-Z2 and create new folder from Jupyter-Notebook



 Transfer two design files to PYNQ-Z2 by using file upload in Jupyter-Notebook



2. Rename "design_1_wrapper.bit" to "design_1.bit"

• Create new Notebook, named "add.ipynb"



Summary

- FPGA Architecture
 - CLB,
 - Programmable Interconnection
 - Modules in Morden FPGA
 - Memory Hierarchy
- Design Optimization
 - o Task
 - Task Latency
 - Task Interval
- Next Week
 - PP4FPGA chapter 2: Hardware Implementation of Finite Impulse Response and Optimizations.