## Speech Recognition

Tsang-Yung Wu Chen-Chia Chang Cynthia Y Liu

Team06



### Outline

- Project Idea 0
- About CNN
- Workflow
- Structure of System
- Software Design
- Calculating Implement 0
- DEMO
- **Future Work**



# Project Idea

### "There is only 21 hours a day, and the remaining 3 hours are for sleeping."

### **Speech Recognition**

- Machine learning is widely used
- Often using GPU parallel operation
- Take advantage of FPGA hardware acceleration
- Implement a speech recognition system on the FPGA
- CNN







### About CNN

"I just want to live in the lab, or else?"





### 

- Convolution Neural Network
- Can reduce the amount of weight compared to DNN
- We use it to implement speech recognition:
  - Convolution \* 2
  - Max pooling
  - Mean pooling
  - Fully connected



## Python predict

### **DCNN**

Current model : mfcc-based with relu precision: 8 bit, clip at 7.9375 Take 200 test data as the benchmark::



	precision	time
General Version	81.5%	35.2s
Weight truncate	82%	34.5s
Clip	80.5%	34s
Clip + Weight truncate	81.5%	32s





### Workflow

"You're too strong."





### Workflow



Transmit pre-trained weight and spectrum

### **RS-232**



![](_page_9_Picture_0.jpeg)

# Structure of System

![](_page_9_Picture_2.jpeg)

# Structure

![](_page_10_Figure_1.jpeg)

![](_page_10_Picture_2.jpeg)

### Structure of System

- In order to use NIOS as the master, put all slaves into QSYS
- PLL synthesizes 75MHz clock, while SDRAM requires 3 ns faster
- Set the written MVM (hardware calculation) to Avalon slave and connect to NIOS.

![](_page_11_Picture_4.jpeg)

![](_page_11_Picture_5.jpeg)

![](_page_11_Picture_6.jpeg)

![](_page_12_Picture_0.jpeg)

### Software Design

![](_page_12_Picture_2.jpeg)

![](_page_12_Picture_3.jpeg)

### Software design NIOS II :

- C code Implementing CNN's flow control
- Memory using SDRAM (128MB)
- Use RS232 to transmit pre-trained weight for storage in SDRAM

![](_page_13_Picture_4.jpeg)

![](_page_13_Picture_5.jpeg)

![](_page_13_Picture_8.jpeg)

## oftware

- **CNN Flow : desiden** 
  - Flatten matrix
- **Convolution two times** 
  - matrix block
  - multiplication(use hardware)
  - reshape and truncate
- Max pooling
- Mean pooling
  - Fully connected

![](_page_14_Figure_12.jpeg)

![](_page_14_Figure_13.jpeg)

![](_page_14_Figure_14.jpeg)

![](_page_14_Picture_15.jpeg)

![](_page_14_Picture_16.jpeg)

MVM

![](_page_14_Picture_17.jpeg)

![](_page_14_Picture_18.jpeg)

![](_page_15_Picture_0.jpeg)

### **Calculating Implement**

![](_page_15_Picture_2.jpeg)

## Controller

- Implement fractional multiplication
- Reduce the required arithmetic unit
- Stochastic number generator

![](_page_16_Figure_4.jpeg)

![](_page_16_Figure_7.jpeg)

# Controller

- Implement fractional multiplication
- Reduce the required arithmetic unit
- Stochastic number generator

![](_page_17_Figure_4.jpeg)

![](_page_17_Figure_6.jpeg)

![](_page_17_Figure_7.jpeg)

## Multiplexer

- The way to implement matrix multipliers
- Synchronous operation during transmission (reduced waiting transmission time) -> Hardware acceleration

![](_page_18_Figure_3.jpeg)

![](_page_18_Figure_4.jpeg)

![](_page_18_Figure_5.jpeg)

![](_page_18_Picture_6.jpeg)

![](_page_19_Picture_0.jpeg)

### 

![](_page_19_Picture_2.jpeg)

![](_page_20_Picture_0.jpeg)

### Future Work "I can finally go back to take a shower."

![](_page_20_Picture_2.jpeg)

### Future work

- Parallel operation using more computing resources
- Change the model representation to reduce storage space
- Actual use of recognized voice messages (ex: robot dog)
- If you have a higher level FPGA, you can try a more complicated model.

![](_page_21_Picture_5.jpeg)

![](_page_21_Picture_7.jpeg)

![](_page_21_Picture_8.jpeg)

### Reference

- Chen, Yu-Hsin, et al. "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks." IEEE Journal of Solid-State Circuits 52.1 (2017): 127-138.
- Sim, Hyeonuk, and Jongeun Lee. "A new stochastic computing multiplier with application to deep convolutional neural networks." Proceedings of the 54th Annual Design Automation Conference 2017. ACM, 2017.
- Kim, Daewoo, et al. "FPGA implementation of convolutional neural network based on stochastic computing." Field Programmable Technology (ICFPT), 2017 International Conference on. IEEE, 2017.
  Alaghi, Armin, and John P. Hayes. "Survey of stochastic computing." ACM Transactions on Embedded
- Alaghi, Armin, and John P. Hayes. "Survey of st computing systems (TECS) 12.2s (2013): 92.

![](_page_22_Picture_5.jpeg)

Any questions?

![](_page_23_Picture_2.jpeg)

![](_page_23_Picture_3.jpeg)

![](_page_23_Picture_4.jpeg)