## Ming Hsieh Department of Electrical Engineering

## **Convolutional Neural Network on FPGA**

Chi Zhang **FPGA/Parallel Computing Lab** 

# USCViterbi

School of Engineering

#### Motivation and Problem Definitions

- Convolutional Neural Network (CNN) achieves the state-of-art performance in image recognition, natural language processing and bioinformatics.
- High computation complexity in both inference and training, which needs specific hardware to accelerate.
- FPGA plays an important role due to its reconfigurability and massive parallelism.
- Automatic Generation Tool: How to design a system that apply to a wide range of CNN architectures remain challenging.

### Approach

- Analyze the basic computational structure of CNN and identify fast algorithm for acceleration.
- Build highly-optimized and parameterized hardware accelerator module for specific algorithm including Matrix operations, Fast Fourier Transform, Winograd algorithm, etc.
- Given CNN architecture, utilize the hardware modules from library to build a full system.
- Explore the system level optimization including memory sub-system, resources distribution and scheduling algorithm.
- Prototype on specific FPGA device.

## Algorithm Optimization Hardware Module Implementation System Level Optimization

Prototype on FPGA

### Algorithm and Hardware

#### CPU + FPGA Mapping

### System Level Optimization







- Increase on-chip data reuse to reduce FPGA-memory bandwidth requirement.
- Exploit task parallelism if there are available resources and bandwidth.



#### **Experiments and Results**



Figure 7: Floating Point Operations required in various state-of-art CNNs

Platform Clock (MHz) Data Precision Bandwidth (GB/s) CNN Model	[21] Virtex7 VX485t 100 32-bit float 12.8 AlexNet	[14] Zynq XC7Z045 150 16-bit fixed 4.2 VGG16-SVD	Our W Intel Quick QPI FP 200 32-bit ff 5.0 AlexNet	ork Assist GA loat VGG16	<ul> <li>Experimenta</li> <li>Intel QuickAs</li> <li>Shared memory</li> <li>6 GB/s FPGA</li> <li>Experimented</li> </ul>	Table V) <b>Perfo</b> • 1.1x I • With eNet arch	e 9: GoogLeNet Convolutional Layer Performance <b>prmance Comparison</b> Ix delay with 3x less resources consumption thout loss of accuracy for any modern CNN chitectures.							
Memory	4.5 MB	2.13 MB	4.0 MB	4.0 MB			AlexNet Execution Time (ms)							
Multipliers	747	780	224	224	Layer (Group)	FPGA	FPGA	CDU	Sequential	Concurrent	FPGA	CDU	Sequential	Congument
(CFLOPS)	61.62	187.80	83.00	123.48		(Theoretical)	(Actual)				(Actual)	CFU	Sequential	Concurrent
			CPU:17.17		CONV1	30.96	31.53	7.76	39.29	32.74	-	17.17	17.17	17.17
Delay (CONV) (ms)	21.61	163.42	FPGA:23.64	263.27	CONV2	44.36	46.01	4.18	50.19	46.48	7.86	0.09	7.95	7.94
Delay×Multipliers	16142	127467	9141/5295	58972	CONV3	81.92	82.27	3.54	85.81	82.75	4.42	0.12	4.54	4.50
Classification Accuracy	Lossless	Lossy	Lossless		CONV4	81.92	82.77	1.37	84.14	82.90	6.64	0.15	6.79	6.71
Flexibility	Any CNN	Limited	Any CNN		CONV5	17.69	18.36	0.27	18.63	18.40	4.42	0.11	4.53	4.49
Table 7: Performance Comparison with the State-of-Art					CONV Total	256.85	262.94	17.12	280.06	263.27	23.34	17.64	40.98	40.81
CNN Implementations on FPGA					Table 6: Execution Time for VGG16 and AlexNet									

**Asymptotic Analysis** 

**GFLOP** Reduction

Space convolution:  $O(N^2F^2)$ 

OaA convolution:  $O(N^2 \log F)$ 

CaffeNet (2012): 48.82%

GoogLeNet (2014): 39.43%

VGG16 (2014): 54.10%

#### GoogLeNet (convolutional layer only) Execution Time (ms) CPU CPU (1 thd) Layer (16 thds)+ FPGA CONV1 38.5312.70CONV2 117.13 17.14 91.05 Inception3 16.1629.16 173.21 Inception4 45.458.48 Inception5 465.37 Total 83.64 Throughput 17.3696.60 (GFLOPS)

Ming Hsieh Institute

Ming Hsieh Department of Electrical Engineering

### Discussions and Future Work

Energy Efficiency. A CPU-FPGA based design will consume more power than FPGA-only based design. However, the CPU adds more flexibility to the design. Moreover, since most of the computations are inside convolutional layer, the CPU simply performs adding and data rearrangement and the energy consumption will not scale up quickly if we increase the CNN size. Automatic Code Generation. Our framework provides complete solution to accelerate CNN on FPGA including inter-layer data rearranging. Modern CNNs' convolutional layers are mainly consist of small kernels. Thus, by zero-padding various kernel sizes to fit a chosen FFT size, and using FPGA to accelerate it by exploiting massive parallelism, we can achieve considerable performance improvement for various CNN models. We can use our framework to develop an automatic code generation tool so high-level users can specify CNN models and generate the design automatically.

Fixed Point vs. Floating Point. Many previous approaches use fixed point instead of floating point for computations. The advantage is less resources consumption and the power efficiency. However, it penalizes the classification accuracy. Some approaches claim that the classification accuracy is tolerable according to experiments. However, it is hard to generalize to an arbitrary CNN model.

fpga.usc.edu