



# The Evolution of Accelerators upon Deep Learning Algorithms

Song Yao, Shuang Liang, Junbin Wang, Zhongmin Chen, Shaoxia Fang, Lingzhi Sui, Qian Yu, Dongliang Xie, Xiaoming Sun, Song Han, Yi shan, and Yu Wang

August 21<sup>st</sup>, 2018, Cupertino

\*The authors represent the combined efforts from many talented and hard-working engineers in DeePhi

www.deephi.com



### 

What's New About DeePhi Tech

**02** The Evolution of Algorithms

**03** The Evolution of DeePhi's DPU

**04** DL Software: New Changes

#### Summary

01

05



# Part 1 What's New About DeePhi Tech





Collaboration with Xilinx University Program

Deep learning acceleration Time series analysis Stereo vision Development of products on Xilinx FPGA platform since inception of DeePhi

Face recognition Video analysis Speech recognition acceleration

•••••

Co-Marketing and Co-Sales with Xilinx Team

> Data Center Automotive Video surveillance

> > . . . . . .







DEEPHi 深鉴科技

© 2018 DeePhi Tech. All Rights Reserved.





https://aws.amazon.com/marketplace/pp/B079N2J42R?from=timeline&isappinstalled=0 https://market.huaweicloud.com/product/00301-110982-0--0



# Part 2

### **Evolution of Algorithms**



#### LeNet-5: 1998



#### AlexNet: 2012



VGG-Net: 2014



GoogLeNet: 2014



ResNet: 2015

DEEPHI Now XILINX.



#### ReLU is not the only widely used activation function



Group Conv: Not every input channel is connected to every output channel



#### Channels can have different weights



Depth-wise Conv and point-wise Conv are widely used



DenseNet: Layers are not necessarily serial



Dilated Conv: Conv kernels are not necessarily dense



Deformable Conv: Conv kernels are not necessarily rectangle

DEEPHI Now XILINX.



## Part 3

### **Evolution of DeePhi's DPU**

Architecture (single core)	DPU-V1	DPU-V2	DPU-V2.5	
Peak performance (GOPS)	120	1350	1367	
On-chip Memory (KB)	300	1123	1123	
LUT	30k	75k	37k	
FF	35k	146k	80k	
Frequency (MHz)	214	333	333	
Typical FPGA Platform	Zynq 7020 (28nm)	ZU9 (16nm)	ZU9 (16nm)	
Total Cores	1	2	3	
Total Peak Perf (GOPS)	120	2700	4100	

#### **Direction of improvements:**

Better scheduling strategy, higher resource utilization, more supported features, and more flexibility

#### From DPU-V1 to DPU-V2.5



#### DPU-V1



#### 1. Convolution

(a) Kernel size=3\*3(b) Stride=1(c) Arbitrary padding size

#### 2. Pooling

(a) Kernel size=2\*2(b) Stride=2

### **3. Activation function** (a) ReLU



DPU-V2

**1. Feature map** (a) Elementwise

# 2. Convolution (a) Arbitrary kernel size (b) Arbitrary stride (c) Arbitrary padding size

#### 3. Pooling

(a) avg/max pooling
(b) Size=2\*2,3\*3
(c) Stride=1,2
(d) Arbitrary padding size
4. Activation function: (a) ReLU
5. FC(INT8)
6. Multi-Batch





Supported network example: LeNet/ResNet50/InceptionV1/InceptionV2/MobileNet\*/YOLO\*/SegNet\*/FPN\* note: the networks with \* is only supported on DPU\_V2.5

#### **DPU-V2.5**

#### **1. Feature map**

- (a) Elementwise
- (b) Split
- (c) Concat
- (d) Resize
- (e) Batch Normalization

#### **2.** Convolution

- (a) arbitrary kernel size
- (b) arbitrary stride
- (c) arbitrary padding size
- (d) Deconv
- (e) Dilated Conv
- (f) Depthwise Conv

#### **3.Pooling**

- (a) avg/max pooling
- (b) arbitrary size
- (c) arbitrary stride
- (d) arbitrary padding
- (e) ROI pooling

#### **4.**Activation function

- (a) ReLU
- (b) PReLU
- (c) LeakyReLU
- (d) Sigmoid

#### 5.FC (INT8/FP32) 6.Multi-Batch



#### Depth-Wise: Number of groups = Number of channels

#### Efficient Networks May Not Be Friendly to Hardware



#### Times of improvement of MobileNet-v1/v2 over ResNet-50 on Tesla P100 GPU

Model	# Parameters	# Computations	Runtime (ms)	Reference Top-1 Accuracy
ResNet-50	25.5M	7.72 GOPS	4.20	76.1%
MobileNet-v1	4.2M	1.14 GOPS	1.37	70.6%
MobileNet-v2	3.5M	0.60 GOPS	1.49	72.0%

DEEPHI Now XILINX.

#### Efficient Networks May Not Be Friendly to Hardware



#### Times of improvement of MobileNet-v1 over three benchmark networks on **DeePhi's DPU on ZU9 FPGA**

Model	Computations	Reference Top-1 Accuracy	Performance	FPS	Utilization Rate	Year
VGG-16	30.7 GOPS	71.9%	2.36 TOPS	76.9	87.30%	2013
GoogleNet	3.89 GOPS	71.0%	0.99 TOPS	254	36.80%	2014
ResNet-50	7.72 GOPS	76.1%	1.06 TOPS	137	39.30%	2015
MobileNet-v1	1.14 GOPS	70.6%	0.77 TOPS	675	28.62%	2017

DEEPHI Now XILINX.





CTC Ratio of MobileNet-V1

#### CTC Ratio of VGG16

- The Communication/Computation (CTC) Ratio of depth-wise is very high
- Point-wise (1x1, PW) conv is memory-bound

#### **Our Strategy Resolving the Problem**





(a) Fused-layer Convolution for MobileNet



(c) Computing Architecture



Ratio of workload (p.w./d.w.)

(b) The workload ratio between adjacent PW-Conv and DW-Conv



#### **Results on ZU9 FPGA**

DEEPHI Now XILINX.

- Dual cores on ZU9 FPGA
- Each core uses 1024KB BRAM
- In total 3648KB BRAM on ZU9 FPGA



### Simulation results with doubled BRAM

Each core uses
 2048KB BRAM



### Part 4

### **DL Software: New Changes**



#### More Optimizations Should be Considered in Compiler





Node fusion/decomposition & data stream optimization



Memory allocation/scheduling/reuse to improve performance



### Part 5 Summary

Discover the Philosophy behind Deep Learning Computing

© 2018 DeePhi Tech. All Rights Reserved.



- Algorithms are evolving at an increasingly faster rate
- Modern neural networks like MobileNet are not necessarily friendly for hardware acceleration on existing ASICs
- We propose a fusing-layer strategy together with compiling optimization to better accelerate Depth-wise/Point-wise convolutions
- More optimization strategies for new types of networks should be considered
- Pure hardware evolves slowly due to the long period in designing and manufacturing chips
- FPGA can be benefited from latest DL techniques in both hardware and software side







# **THANK YOU!**

Friday, August 17, 2018