

FPGA

Chen Zhang,et al

Motivation Background Method Exploration Implementat

Experiment and Result

Conclusion

Optimizing FPGA-Based Accelerator Design for Deep Convolutional Neural Network

Chen Zhang, Peng Li, Guangyu Sun, et al

Presented by Zhuangwei Zhuang

South China University of Technology

August 14, 2017

1/33



Content

FPGA

Chen Zhang,et al

Motivation Background Method Exploration

Experiment and Result

Conclusion

1 Motivation

2 Background

3 Method

4 Exploration

5 Implementation

6 Experiment and Result

7 Conclusion



Motivation Acceleration Platforms

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploration

Implementation

Experiment and Result

Conclusion

Convolutional neural network (CNN) has been widely employed for image recognition for its ability to achieve high accuracy.

Recently, various FPGA-based accelerators for deep CNN have been proposed.

Platforms to Accelerate CNN

- FPGA Field Programmable Gate Array
- GPU Graphics Processing Unit
- ASIC Application Specific Integrated Circuit



Motivation Compare FPGA with GPU

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploration

Implementation

Experiment and Result

Conclusion

Table: Projected ImageNet 1K Performance and Power¹

Platform	ImageNet 1K	Max Device	Efficiency
	(Image/s)	Power (W)	(Image/Sec/W)
Catapult Server + Stratix V D5	134	25	5.36
Catapult Server + Arria 10 GX1150	233	25	9.32
Caffe + cuDNN on Tesla K20	376	235	1.6
Caffe + cuDNN on Tesla K40	500~824	235	2.13~3.51

Table: Projected AlexNet Performance and Power²

CNN Classification Platform	Performance(Image/s)	Power(W)	Efficiency (Image/Sec/W)
E52699 Dual Xeon CPU (18 core per Xeon)	1320	321	4.11
PCIe w/Dual Arria 10 1150	1200	130	9.27

 $^{^1 \}rm Kalin$ et al.Accelerating Deep Convolutional Neural Networks Using Specialized Hardware

 $^{^2 {\}sf Xilinx}^{\textcircled{B}}. {\sf Efficient Implementation of Neural Network Systems Built on {\tt FPGAs, Programmed with OpenCL}}$



Motivation FPGA Structure



Chen Zhang,et a

Motivation

Background

Method

Exploratior

Implementation

Experiment and Result

Conclusion



Figure: FPGA Structure



Motivation FPGA-based Accelerator for Deep CNN

FPGA

Chen Zhang,et al

Motivation

- Background
- Method
- Exploratior
- Implementation
- Experiment and Result
- Conclusion

Advantages

- High performance
- High energy efficiency
- Capability of Reconfiguration
- Fast development round

Disadvantages

- Limited computation resource
- Limited memory bandwidth



Motivation Problem

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploratio

Implementation

Experiment and Result

Conclusion

There would be as much as **90%** performance difference between two different solutions with the same logic resource utilization.

・ロン ・四 と ・ ヨ と ・ ヨ と

7/33



Background FPGA Design in Roofline Model

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploration

Implementation

Experiment and Result

Conclusion



FPGA-based Accelerator Structure



Background FPGA Design in Roofline Model



Chen Zhang,et a

Motivation

Background

Method

Exploratio

Implementation

Experiment and Result

Conclusion





Background FPGA Design in Roofline Model

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploratio

Implementation

Experiment and Result

Conclusion



Computation To Communication (CTC) Ratio



Background Attainable Performance vs. Computation performance

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploratio

Implementatior

Experiment and Result

Conclusion



Computation To Communication (CTC) Ratio



Background Attainable Performance vs. Computation performance

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploratio

Implementation

Experiment and Result

Conclusion



Computation To Communication (CTC) Ratio



Background Optimization

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploration

Implementation

Experiment and Result

Conclusion

Optimizing Target:

Find a design with maximized attainable performance

Optimizing Methods

- Loop Unroll
- Pipeline
- Loop Tiling
- Data Reuse



Background Convolutional Neural Network

FPGA

Chen Zhang,et a

Motivation

Background Method

Exploration

Implementation

Experiment and Result

Conclusion





Background Convolutional Neural Network

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploration

Implementation

Experiment and Result

Conclusion

Convolutional layers account for over **90%** computation^{3,4}

 $^{^3\}text{A}.$ Krizhevsky, etc. Imagenet classification with deep convolutional neural networks. NIPS 2012.

⁴ J. Cong and B. Xiao. Minimizing computation in convolutional neural networks. ICANN 2014. 🗐 🗠 🔿 🗬



Method Loop Unroll

FPGA

Chen Zhang,et a

Motivation Background Method

Exploratio

Implementation

Experiment and Result

Conclusion



 $Y[m][r][c] = \sum_{n=0}^{N-1} \sum_{i=0}^{K-1} \sum_{j=0}^{K-1} W[m][n][i][j] \times X[n][S \times r+i][S \times c+j]$ (1)



Method Loop Unroll

4 5 6

}}}}

FPGA

Method

Pseudo code of a convolutional layer

$$\begin{aligned} & \textbf{for}(j=0;j<\textbf{K};j++) \\ & \text{output}_fm[to][row][col] += \\ & \text{weights}[to][i][j]*input_fm[ti][S*row+i][S*col+j]; \end{aligned}$$

<ロ> <同> <同> < 回> < 回> э 17/33



Ν

Method Loop Tiling

FPGA	Pseudo code of a tiled convolutional layer		
Chen Zhang,et al	External data transfer (Tile Loop)		
otivation ackground ethod xploration	$ \begin{array}{ll} 1 \mbox{ for}(row=0;row$		
periment d Result	On-chip data computation (Point Loop)		
nclusion		C	
	1111	18 / 33	



Method Computation Optimization

FPGA

Chen Zhang,et al

Motivation

Background

Method

Exploration

Implementation

Experiment and Result

Conclusior

Different unrolling way will affect the complexity of generated hardware, and even affect the hardware operation frequency.

On-chip data computation (Point Loop)
<pre>3 for(i=0;i<k;i++){ #pragma="" +="</td" for(i="0;i<K;i++){" for(tcc="col;tcc<min(col+Tc,C);tcc++){" for(ti="t;ti<min(ti+Tn,N);tii++){" for(too="tc;too<min(To+Tm,M);too++){" for(trr="row;trr<min(row+Tr,R);trr++){" hls="" if="" l:="" output_fm[too][trr][tcc]="" pipeline="" unroll=""></k;i++){></pre>

Before optimized

After optimized



Method Computation Optimization

FPGA

Chen Zhang,et a

- Motivation
- Background

Method

- Exploration
- Implementation
- Experiment and Result
- Conclusion



On-chip data computation (Point Loop)

8 for(i=0;i<K;i++){ 9 for(i=0;i<K;i++)10 for(trr=row;trr<min(row+Tr,R);trr++){ for(tcc=col;tcc<min(col+Tc,C);tcc++){</pre> 11 #pragma HLS pipeline 12 for(too=to;too<min(To+Tm,M);too++){</pre> #pragma HLS UNROLL 13 for(tii=ti;ti < min(ti+Tn,N);tii++)#pragma HLS UNROLL 14 L: output_fm[too][trr][tcc] += weight[too][tii][i][j]*
input_fm[tii][S*trr+i][S*tcc+j];

}}}}}



Method Computational Roof

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploration

Implementation

Experiment and Result

Conclusion

Computational roof $= \frac{total number of operations}{number of execution cycles}$ $= \frac{2 \times R \times C \times M \times N \times K \times K}{\left[\frac{M}{T_m}\right] \times \left[\frac{N}{T_n}\right] \times \left[\frac{R}{T_r}\right] \times \left[\frac{C}{T_c}\right] \times (T_r \times T_c \times K \times K + P)}$ $\approx \frac{2 \times R \times C \times M \times N \times K \times K}{\left[\frac{M}{T_m}\right] \times \left[\frac{N}{T_n}\right] \times R \times C \times K \times K}$ $\approx 2 \times T_m \times T_n$ (2)

where

$$\begin{cases} P = pipeline \ depth - 1 \\ 0 < T_m \le M \\ 0 < T_n \le N \\ 0 < T_r \le R \\ 0 < T_c \le C \end{cases} \xrightarrow{(\Box)} (\Box) < (\Box$$



Method Memory Access Optimization

FPGA

Method

External data transfer (Tile Loop)	External data transfer (Tile Loop)
$ 1 \ for(row=0;row$	$ \begin{array}{llllllllllllllllllllllllllllllllllll$

Before local memory promotion

After local memory promotion



Method Computation to Communication Ratio

=

=

where

FPGA

Chen Zhang,et al

Motivation

Background

Method

Exploration

Implementation

Experiment and Result

Conclusion

Computation to Communication Ratio	
total number of operations	
total amount of external data access	
$2\times R\times C\times M\times N\times K\times K$	(2)
$\overline{\alpha_{in} \times B_{in} + \alpha_{weight} \times B_{weight} + \alpha_{out} \times B_{out}}$	(3)

$B_{in} = T_n (ST_r + K - S)(ST_c + K - S)$ (4)

$$B_{wight} = T_m T_n K^2 \tag{5}$$

$$B_{out} = T_m T_r T_c \tag{6}$$

$$0 < B_{in} + B_{wight} + B_{out} \le BRAM_{capacity} \tag{7}$$

$$\alpha_{in} = \alpha_{weight} = \frac{M}{T_m} \times \frac{N}{T_n} \times \frac{R}{T_r} \times \frac{C}{T_c}$$
(8)

$$\alpha_{out} = \frac{M}{T_m} \times \frac{R}{T_r} \times \frac{C}{T_c}$$
(9)

◆□ → ◆□ → ◆ ■ → ◆ ■ → ○ Q (?)
23 / 33



FPGA

Explorat

Exploration Design Space

n	Computation Engine:	Legal Solutions of Tm&Tn:
et al	Constrains for CNN configurations:	
on	$T_m \in (Integer, 1 < T_m < M) \ M = 128$	
und	$T_n \times T_m \in (Integer, 1 < Tn < N)N = 129$	0007
	Constrains for FPGA resource:	2097
on	$T_n \times T_m \in (Integer, 1 < T_m \times T_n < \#of \ PE) \\ \#of \ PE = 450$	
ntation		
ent Ilt	Communication:	Legal Solutions:
on		
	# of memory access methods	3

Total Legal Solutions: 6291



Exploration Design Space Exploration

FPGA

Chen Zhang,et a

Motivation Background

Method

Exploration

Implementation

Experiment and Result

Conclusion





Exploration Multi-Layer CNN Accelerator Design

		0	
E	Р	с а.	А

Chen Zhang,et a

Motivation

Buchgroui

Method

Exploration

Implementation

Experiment and Result

Conclusion

	Optimal Unroll	Execution
	Factor $< T_m, T_n >$	Cycles
Layer 1	< 48, 3 >	366025
Layer 2	< 20, 24 >	237185
Layer 3	< 96, 5 >	160264
Layer 4	< 95, 5 >	120198
Layer 5	< 32, 15 >	80132
Total	-	963804
Cross-Laver Optimization	< 64.7 >	1008246

With unified unroll factors, the degradation is within 5% compared to the total execution cycles of each optimized convolutional layer.



Implementation System Overview





Implementation Accelerator structure





Implementation Resource Utilization

FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploration

Implementation

Experiment and Result

Conclusion



FPGA Resource Utilization

≣ ∽ 29/33



Experiment and Result $_{\text{Vs. CPU}}$

FPGA

Chen Zhang,et a

Motivation Background Method Exploration Implementation Experiment

and Result



CPU	Xeon E5-2430 (32nm)	16 cores	2.2GHz	gcc 4.7.2 -O3 OpenMP 3.0
EDCA	Virtex7-485t	118 PEc	100MH-	Vivado 2013.4
TFGA	(28nm)	440 F LS	100101112	Vivado HLS 2013.4



Experiment and Result

Vs. Other Implementations

FPGA

Chen Zhang,et al

- Motivation Background Method Exploration
- Implementatio

Experiment and Result

Conclusion









Conclusion

FPGA

Chen Zhang,et al

Motivation

Background

Method

Exploration

Implementation

Experiment and Result

Conclusion

An accelerator for convolutional neural network

Contribution:

- Accurate analytical model for computation & communication
- Find the best solution with roofline model

Result:

- On-board run implementation
- 3.5x better performance over other FPGA implementations
- 80% performance/area improvement



FPGA

Chen Zhang,et a

Motivation

Background

Method

Exploratio

Implementation

Experiment and Result

Conclusion

THANK YOU Q & A?