

fpgaConvNet: A Framework for Mapping Convolutional Neural Networks on FPGAs

Stylianos I. Venieris, Christos-Savvas Bouganis stylianos.venieris10@imperial.ac.uk

FCCM 2016, Washington DC 2 May 2016

Deep Learning and Al



Deep Learning Success Stories - ConvNets



Deep Learning on FPGAs

Hand-tuned implementations

Hardware Accelerated Convolutional Neural Networks for Synthetic Vision Systems

Clément Farabet^{1,2}, Berin Martini², Polina Akselrod², Selçuk Talay², Yann LeCun¹ and Eugenio Culurciello² ¹ The Courant Institute of Mathematical Sciences and Center for Neural Science, New York University, USA ² Electrical Engineering Department, Yale University, New Haven, USA

Going Deeper with Embedded FPGA Platform for Convolutional Neural Network

Jiantao Qiu^{1,2}, Jie Wang¹, Song Yao^{1,2}, Kaiyuan Guo^{1,2}, Boxun Li^{1,2},Erjin Zhou¹, Jincheng Yu^{1,2}, Tianqi Tang^{1,2}, Ningyi Xu³, Sen Song^{2,4}, Yu Wang^{1,2}, and Huazhong Yang^{1,2} Memory I/O Optimisation

Memory Access Optimized Routing Scheme for Deep Networks on a Mobile Coprocessor

Aysegul Dundar*, Jonghoon Jin[†], Vinayak Gokhale[†], Berin Martini* and Eugenio Culurciello*[†] *Weldon School of Biomedical Engineering, Purdue University

Memory-Centric Accelerator Design for Convolutional Neural Networks

plications Vs and mo able suppo

Maurice Peemen, Arnaud A. A. Setio, Bart Mesman and Henk Corporaal Department of Electrical Engineering, Eindhoven University of Technology, Eindhoven, the Netherlands Email: m.c.j.peemen@tue.nl, arnaud.arindra.adiyoso@student.tue.nl, b.mesman@tue.nl, h.corporaal@tue.nl

Abstract—In the near future, cameras will be used everywhere as flexible sensors for numerous applications. For mobility and privacy reasons, the required image processing should be local on embedded computer platforms with performance requirements



Design Space Exploration

Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks

Chen Zhang¹ chen.ceca@pku.edu.cn Peng Li² pengli@cs.ucla.edu Guangyu Sun^{1,3} gsun@pku.edu.cn

Yijin Guan¹ guanyijin@pku.edu.cn Bingjun Xiao² xiao@cs.ucla.edu Jason Cong^{2,3,1,*} cong@cs.ucla.edu

Throughput-Optimized OpenCL-based FPGA Accelerator for Large-Scale Convolutional Neural Networks

Naveen Suda, Vikas Chandra[‡], Ganesh Dasika^{*}, Abinash Mohanty, Yufei Ma, Sarma Vrudhula[†], Jae-sun Seo, Yu Cao



Our approach - fpgaConvNet



Convolutional Neural Networks (ConvNets)



fpgaConvNet – ConvNet Modelling Framework



- Each layer mapped to a tunable set of *hardware building blocks*

fpgaConvNet - Modelling ConvNets with SDF

ConvNet Hardware SDF Graph



fpgaConvNet - Design Space Perspective



Bottlenecks:

- Limited compute resources
- Limited off-chip memory bandwidth
- Limited *on-chip memory* for model parameters

Define a set of *actions* to move around the design space

Action 1: Coarse-grained Folding



Action 1: Coarse-grained Folding



Action 3: Partitioning through Reconfiguration



fpgaConvNet – SDF Analytical Power



Evaluation - Experimental Setup

- fpgaConvNet
 - Xilinx Zynq-7000 XC7Z020 SoC with 220 DSPs at 100 MHz
 - Q8.8 fixed-point precision to match existing work (also supports floating-point)
 - Current toolflow supports the Vivado HLS toolchain

Performance Model Accuracy

Performance Model Accuracy Scene Labelling ConvNet Sign Recognition ConvNet CNP Error between 1.73% and 11.76% MPCNN LeNet-5 CFF 12 2 10 0 4 6 8 14

Measured Performance (GOps/s) Predicted Performance (GOps/s)

fpgaConvNet vs. Existing FPGA Work

Performance Density Comparison (GOps/s/Slice)



[1] C. Farabet et al., "CNP: An FPGA-Based Processor for Convolutional Networks", in FPL, IEEE, 2009.

[2] M. Peemen et al., "Memory-centric accelerator design for Convolutional Neural Networks", in ICCD, IEEE, 2013.

fpgaConvNet vs. Existing Embedded GPU Work





- Existing Work (GOps/s/Watt)
- fpgaConvNet (GOps/s/Watt)

Hand-tuned Embedded GPU

- Tegra K1 at 800 MHz
- Memory Bandwidth: 12 GB/s

fpgaConvNet

- Zynq-7000 XC7Z020 at 100 MHz
- Memory Bandwidth: 4.26 GB/s

[3] L. Cavigelli et al., "Accelerating real-time embedded scene labeling with convolutional networks", in DAC, ACM/EDAC/IEEE, 2015.

