ALL PROGRAMMABLE





A Framework for Reduced Precision Neural Networks on FPGAs

Kees Vissers Xilinx Research

- Background
- CNNs and their challenges
- Quantized neural networks
- > FINN: Framework for fast exploitation of Quantized Neural Networks

MPSoC 2017

© Copyright 2017 Xilinx

- > Experimental results
- Conclusion

What are FPGAs?

Programmable devices that contain:

- > MAC (DSP48) for floating point, 16 bit integer, 8 bit integer
- Logic Lookup tables (LUTS) for any function at bit-precision, including
 2 bit and 1 bit MAC (xnor, popcount), and compression, security, etc.
- > Large number of flexible memory blocks, with high internal bandwidth
- > High-speed I/O, good external memory interfaces
- > ARM cores
- > Family of devices

Customizable hardware architectures with fine-grain programmability



Challenge 1: Diverse Applications with Diverse Design Targets



Challenge 2: Neural Networks Will Continue to Change



Customized ML Processor Datapath



Challenge 3: Highly Compute and Memory Intensive

> The predominant CNN computation is linear algebra

- Demands lots of (simple) computation and lots of parameters (memory)
 - AlexNet: 244MB & 1.5GOPS, VGG16: 552MB & 30.8GOPS; GoogleNet: 41.9MB & 3.0GOPS for ImageNet





Challenge 3:

billions of multiply-accumulate ops & tens of megabytes of parameter data

Increasingly Reduced Precision Networks

> Floating point (FP) CNNs contain a lot of redundancy

- Reducing precision is shown to work down to 1b with minimal loss of accuracy –
 - ICLR 2017 with ternary weight networks on par with FP for AlexNet top-1 and top-5, ResNet20,32,44,56
 - Accuracy gap is closing
- > Reducing precision brings numerous advantageous

MPSoC 2017

© Copyright 2017 Xilinx

- Power
- Performance
- Memory requirements
- Not just for FPGAs

		Re	lative	Energ	y Cost	
Operation:	Energy (pJ)] .				
8b Add	0.03					
16b Add	0.05					
32b Add	0.1					
16b FP Add	0.4					
32b FP Add	0.9					
8b Mult	0.2					
32b Mult	3.1					
16b FP Mult	1.1					
32b FP Mult	3.7					
32b SRAM Read (8KB)	5					
32b DRAM Read	640					
		1 1	0	100	1000	10000

Source: Bill Dally (Stanford), Cadence Embedded Neural Network Summit, February 1, 2017

XILINX ➤ ALL PROGRAMMABLE.

Accuracy of Quantized Neural Networks (QNNs) Improving Published Results for FP CNNs, QNNs and binarized NNs (BNNs)



 Accuracy results are improving rapidly through for example new training techniques, topological changes and other methods

MPSoC 2017 © Copyright 2017 Xilinx

Potential of Reduced Precision on FPGAs

- > Cost per operation is greatly reduced
 - For example, for BNN: FP multiply accumulate becomes XNOR with bit counts
- > Memory cost is greatly reduced
 - Large networks can fit entirely into on-chip memory (OCM) (UltraRAM, BRAM)

UT DSP 100ks LUTs Ks DSPs

- > Today's FPGAs have a much higher peak performance for reduced precision operations
 - FPGA performance is anti-proportional to the cost per operation when applications are sufficiently parallel
 - Lower cost per op & massively parallel = more ops every cycle

Precision	Cost per Op LUT	Cost per Op DSP	MB needed (AlexNet)	TOps/s (KU115)*	TOps/s (VU9P)**	TOps/s (ZU19EG)*
1b	2.5	0	7.6	~46	~100	~66
4b	16	0	30.5	~11	~15	~16
8b	45	0	61	~3	~6	~4
16b	15	0.5	122	~1	~4	~1
32b	178	2	244	~0.5	-1	~0.3
<u>Assumptions</u> : Applica	ation can fill device to 70 cation can fill device to 70	% (fully parallelizable) 250 0% (fully parallelizable) 30 MPSoC 2017	0MHZ 00MHZ 7 © Copyright 20		amazon webservices	

© Copyright 2017 Xilinx

Potential of QNNs on FPGAs (ZU19EG)



Exploitation of Quantized Neural Networks through

FINN: A Framework for Fast, Scalable Neural Network Inference



https://arxiv.org/abs/1612.07119 http://arxiv.org/abs/1701.03400

FINN Design Principles

Custom-tailored hardware

- -Customized data types
- Customized dataflow architecture to match network topology
- Keep all parameters on-chip, if possible

>C++ design entry

To support portability, scalability & rapid exploration



Customized Dataflow Architecture

Work Flow for Exploration of NNs of FPGAs



Experimental Results

- Embedded platforms (Zyng Z7045 & 7020): ZC706, PYNQ open source platform
- Server class accelerator: ADM_PCIE_8K5 & TUL Accel. kit in OpenPOWER (& x86 with SDAccel)

© Copyright 2017 Xilinx









EXILINX > ALL PROGRAMMABLE.

Input Data





Page 16



numbers

MPSoC 2017



 German road signs

© Copyright 2017 Xilinx



EXILINX > ALL PROGRAMMABLE.





Test Networks

- > Multilayer Perceptron
 - Input images: 28x28 pixels, black-white (MNIST)
 - Up to 5.8MOPS/frame

> VGG-16 derivative

- Input images: 32x32 pixels, RGB image (SVHN, CIFAR-10, traffic signs, playing cards)
- Up to 1.2GOPS/frame

> DorefaNet

- Input images: 226x226 pixels, RGB (ImageNet)
- Up to
- > YoloV2, TinyYolo
 - Input images: 448x448, RBG (VOC, COCO)
 - 35 and 7GOPS/frame





QNN Results - Latency, Performance/Power



PYNQ FINN Open Source Release

- > FINN is open sourced and available at
 - -<u>https://github.com/Xilinx/BNN-PYNQ</u>
 - New features are continuously rolled out
- Supported on low cost open source platform Pynq
 - Visit <u>www.pynq.io</u>
- Easy to use with precooked overlays & examples – BNNs
- > 1000x faster than Raspberry Pi3

jupyter

Z7020 ARM FPS	Raspberry Pi3 FPS
17.3	44.3
1.2	2.3

EXILINX > ALL PROGRAMMABLE,

LFC 168k 974 112 30.6K 102 <2.5 (80%) (57.6%) CNV 3.04k 341 140 28.5K 1580 <2.5 (100%) (53.5%)	Z7020	FPS (FPGA)	GOPS/s	BRAM	LUT	Latency [us]	Power [W]
CNV 3.04k 341 140 28.5K 1580 <2.5 (100%) (53.5%)	LFC	168k	974	112 (80%)	30.6K (57.6%)	102	<2.5
	CNV	3.04k	341	140 (100%)	28.5K (53.5%)	1580	<2.5

MPSoC 2017

© Copyright 2017 Xilinx

Summary

> Benefits of FPGA implementations:

- -Extreme performance with reduced precision
- -Low latency through dataflow no batching needed
- -Flexibility
- Low power total solution: keep data on chip, compress data, compute at reduced precision (good for memory bandwidth too)

MPSoC 2017

© Copyright 2017 Xilinx





> Get started with FINN & Pynq

-www.pynq.io

-https://github.com/Xilinx/BNN-PYNQ

The Name of the Game: Designing Hardware-Optimal CNNs *It's a trade-off...*



hardware cost/ performance/ power

Example: ImageNet Classification Published Results & Xilinx Research internal Experiments



Floating point is too expensive

Below 10% uses ensembles and cost likely prohibitively high Pareto optimal: 1b – 8b provide good compromises

Inference Accelerators – Accuracy vs Hardware Cost for a fixed topology

- Just reducing precision, reduce hardware cost & increases error
- Recuperate accuracy by retraining & increasing network size
- 1b, 2b and 4b provide pareto optimal solutions





Conclusions: We're only at the start...

- > We presented a framework for exploring Neural Networks at any precision ranging from 32bit Floating Point to 1bit for weight and activation.
- > We have shown that for a number of cases optimal networks can use compute and storage below 8bit!
- Lots of scope for research in exploring the design space between accuracy, performance, cost, power etc.

> Very Exciting times for Neural Networks on heterogeneous platforms.

XII INX > ALL PROGRAMMABLE.



Thanks to a large team at Xilinx including Xilinx Research Ireland











© Copyright 2017 Xilinx

Technical Details on Finn architectures



Architecture of a Matrix-Vector Threshold Unit (MVTU)

- Fully connected layers & convolutional layers are mapped on matrix-vector multiply threshold units (MVTUs)
- > MVTUs support OFM (neuron) and folding over weights (synaptic)
- > Weight and output stationary (weights and popcounts are retained locally)
- > Max pool units are optionally placed behind MVTUs



XILINX > ALL PROGRAMMABLE.

Synthesizable C++ Network Description



MVTU



Architecture of Infrastructure on Zynq SOC

MPSoC 2017



© Copyright 2017 Xilinx

XILINX > ALL PROGRAMMABLE.