

# Quantisation

Efficient implementation of convolutional neural networks

Philip Leong

Computer Engineering Lab  
The University of Sydney

July 2018 / PAPAA Workshop

# Australia



## *Australia and Europe Area size comparison*

Darwin to Perth 4396km • Perth to Adelaide 2707km • Adelaide to Melbourne 726km  
Melbourne to Sydney 887km • Sydney to Brisbane 972km • Brisbane to Cairns 1748km



Population: 24M (2016)  
Europe: 741M (2016)  
Hong Kong: 7M (2016)  
Area 1/25<sup>th</sup> Tasmania

# Outline

## 1 Introduction

Number Systems

Convolutional Neural Networks

Integer Quantisation

SYQ: Low Precision DNN Training

FINN: A Binarised Neural Network

## 2 Tutorial

# Outline

## 1 Introduction

### Number Systems

Convolutional Neural Networks

Integer Quantisation

SYQ: Low Precision DNN Training

FINN: A Binarised Neural Network

## 2 Tutorial

# Introduction

- There are several degrees of freedom to explore when optimising DNNs
  - NN architecture (SqueezeNet, MobileNet)
  - Compression (SVD, Deep Compression, Circulant)
  - Quantization (FP16, TF-Lite, FINN, DoReFa-Net)
- This talk: quantisation

# Unsigned Numbers

$$\begin{aligned} U &= (u_{W-1}u_{W-2}\dots u_0), u_i \in \{0, 1\} \\ &= \sum_{i=0}^{W-1} u_i 2^i \end{aligned}$$

- $U$  is a  $W$ -bit unsigned integer
- Range  $[0, 2^W)$

# Two's Complement Numbers

$$\begin{aligned} X &= (x_{W-1}x_{W-2}\dots x_0), x_i \in \{0, 1\} \\ &= -x_{W-1}2^{W-1} + \sum_{i=0}^{W-2} x_i2^i \end{aligned}$$

- $X$  is a  $W$ -bit signed integer
- Range  $[-2^{W-1}, 2^{W-1})$

## Two's Complement Fractions

$$\begin{aligned} Y &= \left( \overbrace{y_{W-1} \dots y_F}^{\text{I-bit integer}} \overbrace{y_{F-1} \dots x_0}^{\text{F-bit fraction}} \right), y_i \in \{0, 1\} \\ &= 2^{-F} \times \left( -x_{W-1} 2^{W-1} + \sum_{i=0}^{W-2} x_i 2^i \right) \end{aligned}$$

- Y is a W-bit signed fraction with F-bit fraction
- Are two's complement numbers scaled by  $2^{-F}$
- Notation used: (I,F) (with  $I + F = W$ )
  - (W,0) same as two's complement integers
  - (1,W-1) has range [-1,1) and multiplication never overflows

## Dynamic Fixed Point [CBD14]

$$D = (-1)^S \cdot 2^{-F} \sum_{i=0}^{W-2} x_i 2^i$$

- $D$  is dynamic fixed point number with sign bit  $S$ , fractional length  $F$ ,  $W$  is word length
- Sign-magnitude fraction with  $F$  being shared within a group
- Allows number format to be adapted to different network segments e.g. layer inputs, weights and outputs can have different  $F$

# Operations on Two's Complement Fractions

- Addition and subtraction same as two's complement
- Multiplication
  - An (I,F) multiplication gives a (2I,2F) result, need to discard F bits
  - For (1,3)

$$\begin{aligned}0.75 \times 0.75 &= 0.110 \times 0.110 \\ &= 00.100100 \quad \text{in (2I,2F) format} \\ &\approx 0.100 \quad \text{in (I,F) format (truncated)}\end{aligned}$$

- Integer part controls range
- Fractional part controls spacing between numbers

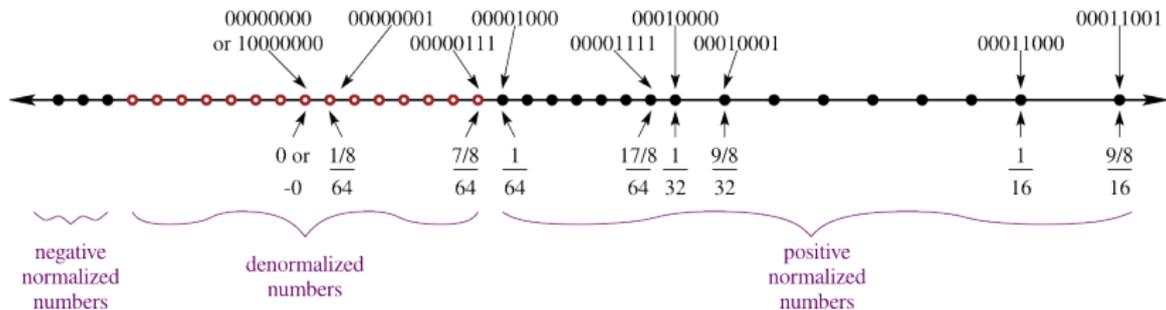
# Floating Point 1

$$Z = (\overbrace{a_0}^A \overbrace{b_{J-1} \dots b_0}^B \overbrace{c_{F-1} \dots c_0}^C), (a_i, b_i, c_i) \in \{0, 1\}$$

- Treating A, B and C as unsigned integers
  - The sign bit is  $S = \begin{cases} +1 & \text{if } a_0 = 0 \\ -1, & \text{otherwise} \end{cases}$
  - The exponent is stored in a biased representation with  $E = B - (2^{J-1} - 1)$
  - For normalised numbers,  $B \neq 0$ , and  $M$  is a positive (1,F) two's complement fraction  $M = 1 + C2^{-F}$
  - For denormalised numbers  $B = 0$  and there is no implicit 1 in the positive (0,F) two's complement fraction  $M = C2^{-F}$

## Floating Point 2

$$Z = \begin{cases} S \times 2^E \times M & \text{if } (0 < B < 2^J - 1) \\ S \times 2^E \times (M - 1) & \text{if } (B = 0) \\ S \times \infty & \text{if } (B = 2^J - 1 \text{ and } C = 0) \\ \text{NaN} & \text{if } (B = 2^J - 1 \text{ and } C \neq 0) \end{cases}$$



# Operations on Floating Point Numbers

- Much larger resource utilisation
- Longer latency
- We will focus on fixed point

# Outline

## 1 Introduction

Number Systems

**Convolutional Neural Networks**

Integer Quantisation

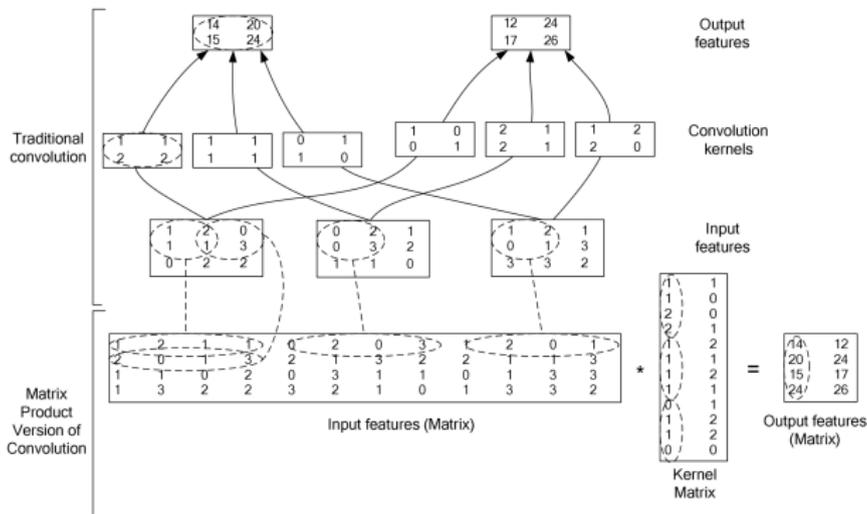
SYQ: Low Precision DNN Training

FINN: A Binarised Neural Network

## 2 Tutorial

# Convolution Layer as MM

- Convolution layers converted to GEMM [CPS06]
- Efficient BLAS libraries can be exploited



## DNN Computation

Computational problem in DNNs is to compute a number of dot products

$$h = g(\mathbf{w}^T \mathbf{x}) \quad (1)$$

where

- $g$  is an element-wise nonlinear activation function
- $\mathbf{x} \in \mathbb{R}^{i.w.h}$  is the input vector
- $\mathbf{w} \in \mathbb{R}^{i.w.h}$  is the weight vector

# Arithmetic Intensity

- Computation of a DNN layer is MV multiplication
- For MV multiply is  $O(1)$ , for MM is  $O(b)$  where  $b$  is block size
- Efficient CPU/GPU implementations use batch size  $\gg 1$  (process a number of inputs together)
- For latency-critical applications (e.g. object detection for self-driving car), we want a batch size of 1
- **Make sure comparisons are at the same batch size!**

# Outline

## 1 Introduction

Number Systems

Convolutional Neural Networks

**Integer Quantisation**

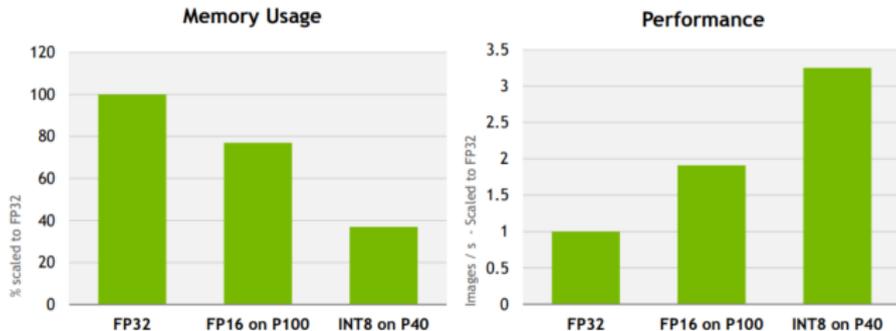
SYQ: Low Precision DNN Training

FINN: A Binarised Neural Network

## 2 Tutorial

# Role of Wordlength on Performance

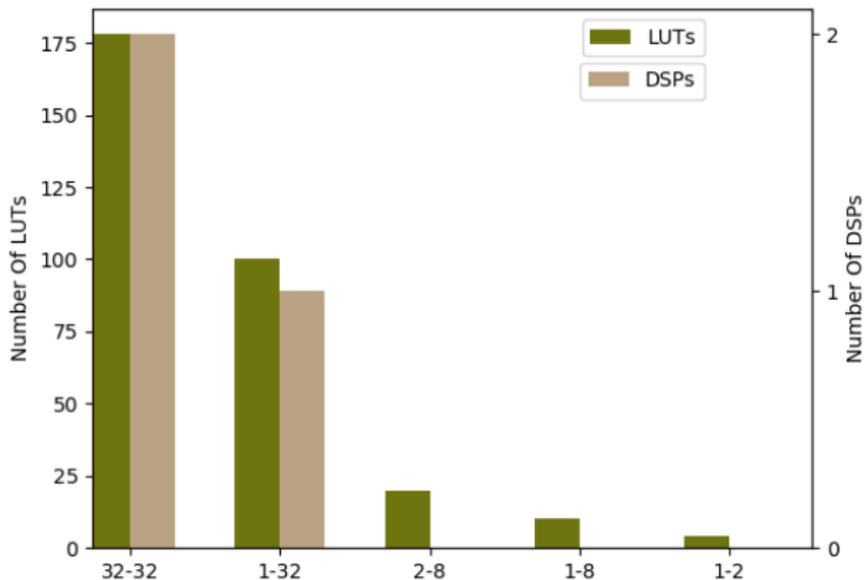
- CPU/GPU
  - Floating point performance comparable to fixed
  - Integer data types usually vectorisable hence faster
  - Nvidia offers FP64, FP32 and FP16 (> Tegra X1 and Pascal)
- FPGA
  - Datapath is flexible
  - No floating point unit so fixed point normally preferred



ResNet50 Model, Batch Size = 128, TensorRT 2.1 RC prerelease

## Role of Wordlength on Resources

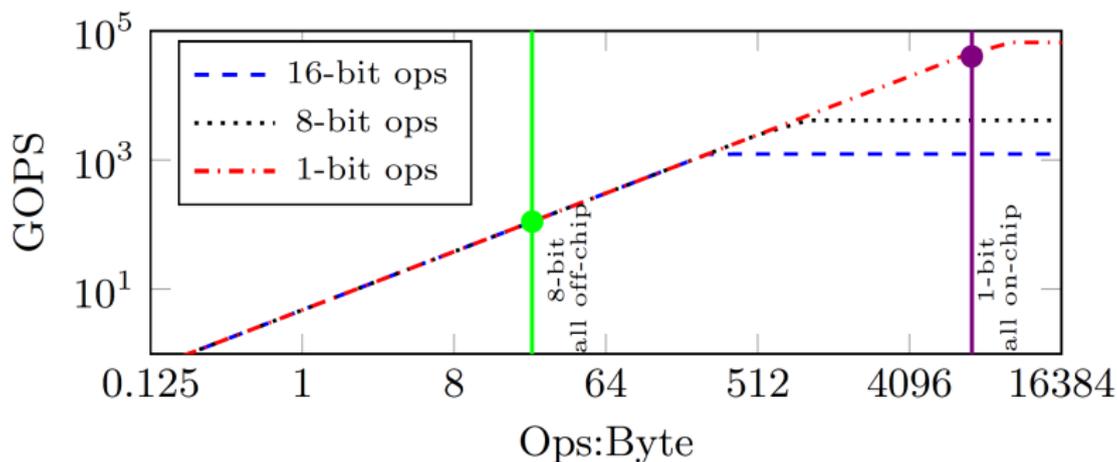
- X axis is bitwidth (weight-activation) and Y axis Number of LUTs/DSPs for MAC
- For  $k$ -bits, area is  $\mathcal{O}(k^2)$



## Roofline Model

### Roofline model for Xilinx ZU19EG

- X axis is computational intensity (ops to perform / byte fetch), Y axis is performance
- Diagonal parts show memory-bandwidth limited space
- Horizontal parts show computation limited space
- Actually this is a better metric to optimise than say GOPs/s
- **Low precision extremely advantageous for performance**



## Integer Quantization [Jac+18]

A way to map numbers  $r \in \mathbb{R}$  to unsigned integers  $q \in \mathbb{U}_+$  is via an affine transformation

$$r = S(q - Z) \quad (2)$$

- $\mathbb{U}_+$  is the set of unsigned  $W$ -bit integers
- $S, Z$  are the quantisation parameters
  - $S \in \mathbb{R}_+$  represents a scaling constant
  - $Z \in \mathbb{U}_+$  represents a zero-point

## Integer MM [Jac+18]

- $N \times N$  MM defined as

$$r_3^{(i,k)} = \sum_{j=1}^N r_1^{(i,j)} r_2^{(j,k)}, \quad (3)$$

substituting  $r = S(q - Z)$  (2) and rewriting we get

$$q_3^{(i,k)} = Z_3 + M \left( NZ_1 Z_2 - Z_1 a_2^{(k)} - Z_2 \bar{a}_1^{(i)} + \sum_{j=1}^N q_1^{(i,j)} q_2^{(j,k)} \right) \quad (4)$$

- Multiplication with  $M = \frac{S_1 S_2}{S_3}$  is implemented in (high-precision) two's complement fixed point
- $a_2^{(k)}$  and  $\bar{a}_1^{(i)}$  together only take  $2N^2$  additions
- Sum in (4) takes  $2N^3$  and is a standard integer MAC
- CPU implementation uses uint8, accumulated as int32

## Quantisation Range [Jac+18]

For each layer, quantisation parameterised by  $(a,b,n)$ :

$$\text{clamp}(r; a, b) = \min(\max(x, a), b)$$

$$s(a, b, n) = \frac{b - a}{n - 1}$$

$$q(r; a, b, n) = \text{rnd}\left(\frac{\text{clamp}(r; a, b) - a}{s(a, b, n)}\right)s(a, b, n) + a \quad (5)$$

where  $r \in \mathbb{R}$  is number to be quantised,  $[a,b]$  is quantisation range,  $n$  is number of quantisation levels and  $\text{rnd}()$  rounds to nearest integer

Figure from [Jac+18] (with permission)

# Training Algorithm [Jac+18]

- 1 Create training graph of the floating-point model
- 2 Insert quantisation operations for integer computation in inference path using (5)
- 3 Train with quantised inference but floating-point backpropagation until convergence
- 4 Use weights thus obtained for low-precision inference

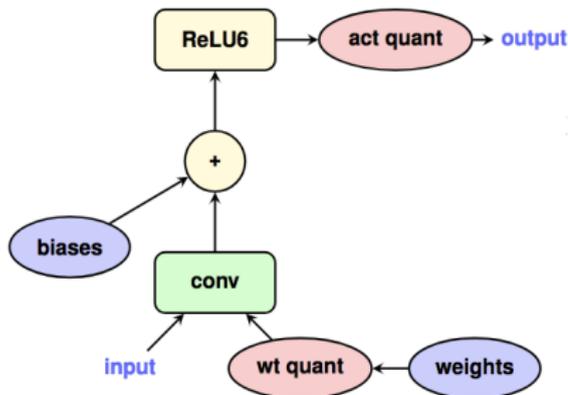


Figure from [Jac+18] (with permission)

## Accuracy vs Precision [Jac+18]

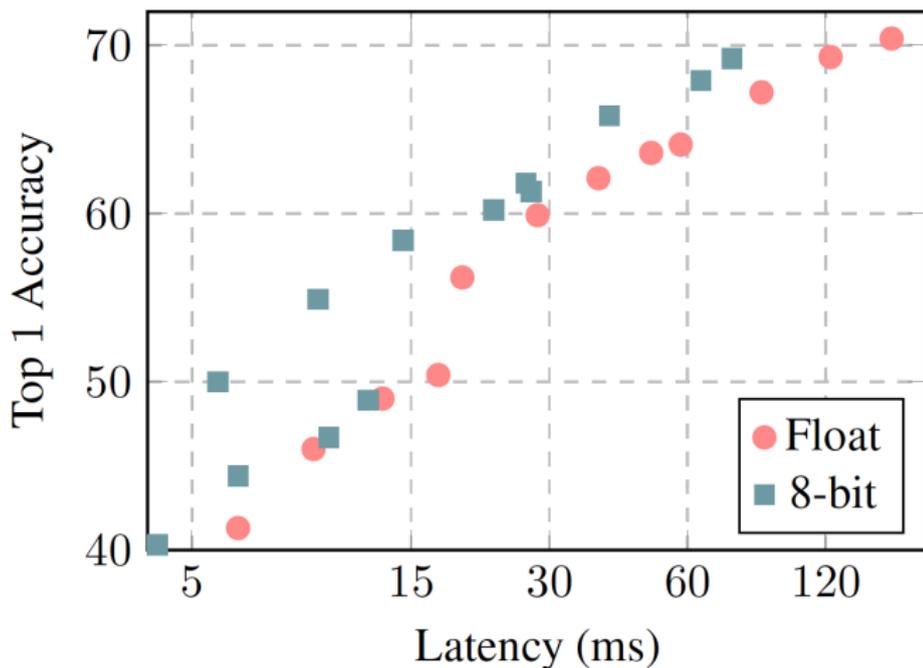
ResNet50 on ImageNet, comparison with other approaches

Scheme	BWN	TWN	INQ	FGQ	Ours
Weight bits	1	2	5	2	8
Activation bits	float32	float32	float32	8	8
Accuracy	68.7%	72.5%	74.8%	70.8%	74.9%

Table 4.2: ResNet on ImageNet: Accuracy under various quantization schemes, including binary weight networks (BWN [21, 15]), ternary weight networks (TWN [21, 22]), incremental network quantization (INQ [33]) and fine-grained quantization (FGQ [26])

## Accuracy vs Latency [Jac+18]

ImageNet classifier on Google Pixel 2 (Qualcomm Snapdragon 835 big cores)



# Outline

## 1 Introduction

Number Systems

Convolutional Neural Networks

Integer Quantisation

**SYQ: Low Precision DNN Training**

FINN: A Binarised Neural Network

## 2 Tutorial

# Symmetric Quantisation (SYQ) [Far+18]

- To compute quantised weights from FP weights

$$\mathbf{Q}_l = \text{sign}(\mathbf{W}_l) \odot \mathbf{M}_l \quad (6)$$

with,

$$M_{l,i,j} = \begin{cases} 1 & \text{if } |W_{l,i,j}| \geq \eta \\ 0 & \text{if } -\eta < W_{l,i,j} < \eta \end{cases} \quad (7)$$

$$\text{sign}(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{otherwise} \end{cases} \quad (8)$$

where  $\mathbf{M}$  represents a masking matrix,  $\eta$  is the quantization threshold hyperparameter (0 for binarised)

# Symmetric Quantisation (SYQ) [Far+18]

- Make approximation  $W_I \approx \alpha_I Q_I$ ,  $Q_I \in \mathcal{C}$
- $\mathcal{C}$  is the codebook,  $\mathcal{C} \in \{\mathcal{C}_1, \mathcal{C}_2, \dots\}$  e.g.  $\mathcal{C} = \{-1, +1\}$  for binary,  $\mathcal{C} = \{-1, 0, +1\}$  for ternary
- A diagonal matrix  $\alpha_I$  is defined by the vector  $\alpha_I = [\alpha_I^1, \dots, \alpha_I^m]$ :

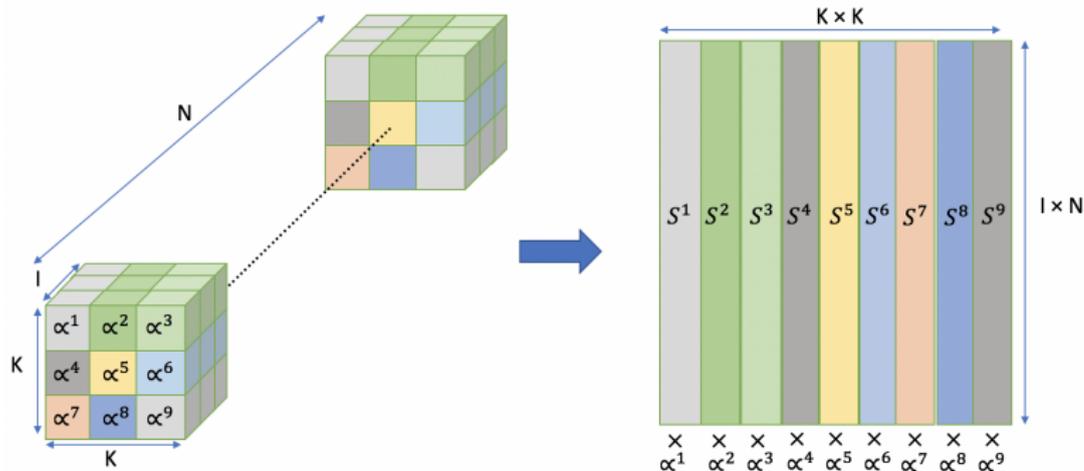
$$\alpha = \text{diag}(\alpha) := \begin{bmatrix} \alpha^1 & 0 & .. & 0 & 0 \\ 0 & \alpha^2 & .. & \vdots & 0 \\ \vdots & \vdots & .. & \alpha^{m-1} & \vdots \\ 0 & 0 & .. & 0 & \alpha^m \end{bmatrix}$$

- Train by solving

$$\alpha_I^* = \underset{\alpha}{\text{argmin}} E(\alpha, \mathbf{Q}) \quad \text{s.t.} \quad \alpha \geq 0, Q_{I,i,j} \in \mathcal{C}$$

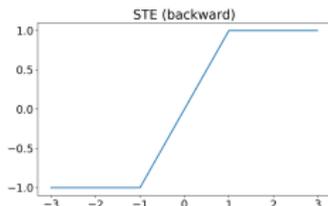
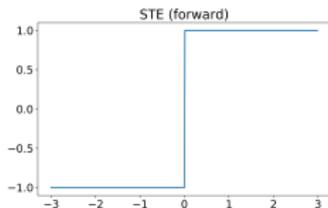
# Subgroups

- Finer-grained quantisation improves weight approximation
- Pixel-wise shown, layer-wise has similar accuracy



# Dealing with Non-differentiable Functions

- Recall (6)  $\mathbf{Q}_l = \text{sign}(\mathbf{W}_l) \odot \mathbf{M}_l$
- This step function has a derivative which is zero everywhere: *vanishing gradients* problem
- Address via a straight through estimator (STE)
- Consider  $q = \text{sign}(r)$  and  $g_r \approx \frac{\partial C}{\partial q}$  then  $\frac{\partial C}{\partial r} \approx g_q \mathbf{1}_{|r| \leq 1}$



## Results for 8-bit activations

<b>Model</b>		<b>Bin</b>	<b>Tern</b>	<b>FP32</b>	<b>Reference</b>
AlexNet	Top-1	<b>56.6</b>	<b>58.1</b>	56.6	57.1
	Top-5	<b>79.4</b>	<b>80.8</b>	80.2	80.2
VGG	Top-1	<b>66.2</b>	<b>68.7</b>	69.4	-
	Top-5	<b>87.0</b>	<b>88.5</b>	89.1	-
ResNet-18	Top-1	<b>62.9</b>	<b>67.7</b>	69.1	69.6
	Top-5	<b>84.6</b>	<b>87.8</b>	89.0	89.2
ResNet-34	Top-1	<b>67.0</b>	<b>70.8</b>	71.3	73.3
	Top-5	<b>87.6</b>	<b>89.8</b>	89.1	91.3
ResNet-50	Top-1	<b>70.6</b>	<b>72.3</b>	76.0	76.0
	Top-5	<b>89.6</b>	<b>90.9</b>	93.0	93.0

- Our ResNet and AlexNet reference results are obtained from <https://github.com/facebook/fb.resnet.torch> and <https://github.com/BVLC/caffe>

## Alexnet Comparison

Model	Weights	Act.	Top-1	Top-5
DoReFa-Net [Zho+16]	1	2	49.8	-
QNN [Hub+16]	1	2	51.0	73.7
HWGQ [Cai+17]	1	2	52.7	76.3
<b>SYQ</b>	<b>1</b>	<b>2</b>	<b>55.2</b>	<b>78.4</b>
DoReFa-Net [Zho+16]	1	4	53.0	-
<b>SYQ</b>	<b>1</b>	<b>4</b>	<b>56.2</b>	<b>79.4</b>
BWN [Ras+16]	1	32	56.8	79.4
<b>SYQ</b>	<b>1</b>	<b>8</b>	<b>56.6</b>	<b>79.4</b>
<b>SYQ</b>	<b>2</b>	<b>2</b>	<b>55.7</b>	<b>79.1</b>
FGQ [Mel+17]	2	8	49.04	-
TTQ [Zhu+16]	2	32	57.5	79.7
<b>SYQ</b>	<b>2</b>	<b>8</b>	<b>58.1</b>	<b>80.8</b>

# Outline

## 1 Introduction

Number Systems

Convolutional Neural Networks

Integer Quantisation

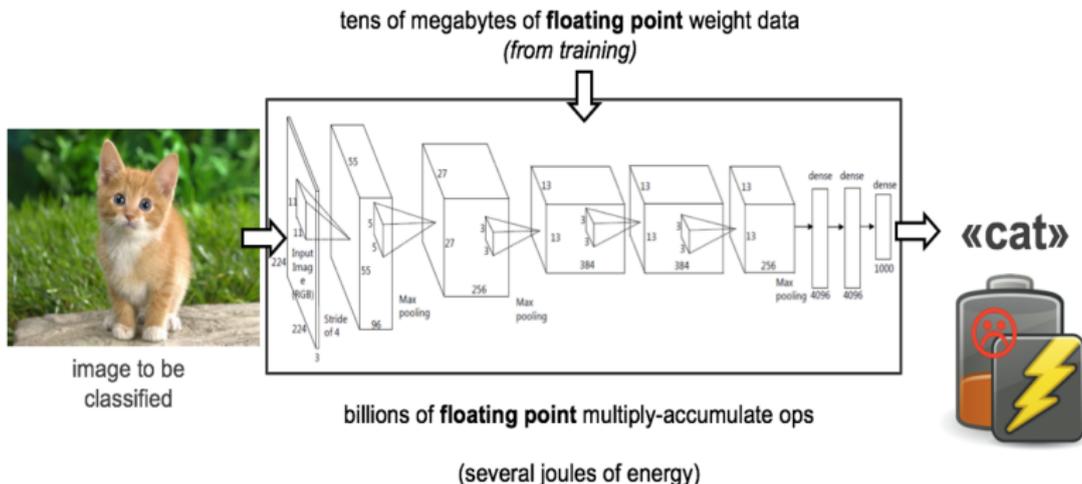
SYQ: Low Precision DNN Training

**FINN: A Binarised Neural Network**

## 2 Tutorial

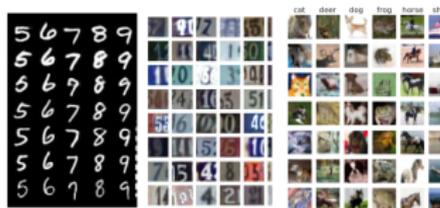
# Inference with Convolutional Neural Networks

Slides from Yaman Umuroglu et. al., "FINN: A framework for fast, scalable binarized neural network inference," FPGA'17



# Binarized Neural Networks

- › The extreme case of quantization
  - Permit only two values: +1 and -1
  - Binary weights, binary activations
  - Trained from scratch, not truncated FP



- › Courbariaux and Hubara et al. (NIPS 2016)

- Competitive results on three smaller benchmarks
- Open source training flow
- Standard “deep learning” layers
  - Convolutions, max pooling, batch norm, fully connected...

	MNIST	SVHN	CIFAR-10
Binary weights & activations	0.96%	2.53%	10.15%
FP weights & activations	0.94%	1.69%	7.62%
BNN accuracy loss	-0.2%	-0.84%	-2.53%

*% classification error (lower is better)*

# Advantages of BNNs

## Vivado HLS estimates on Xilinx UltraScale+ MPSoC ZU19EG

### › *Much* smaller datapaths

- Multiply becomes XNOR, addition becomes popcount
- No DSPs needed, everything in LUTs
- Lower cost per op = more ops every cycle

### › *Much* smaller weights

- Large networks can fit entirely into on-chip memory (OCM)
- More bandwidth, less energy compared to off-chip

Precision	Peak TOPS	On-chip weights
1b	~66	~70 M
8b	~4	~10 M
16b	~1	~5 M
32b	~0.3	~2 M

↑ 200x (Peak TOPS)

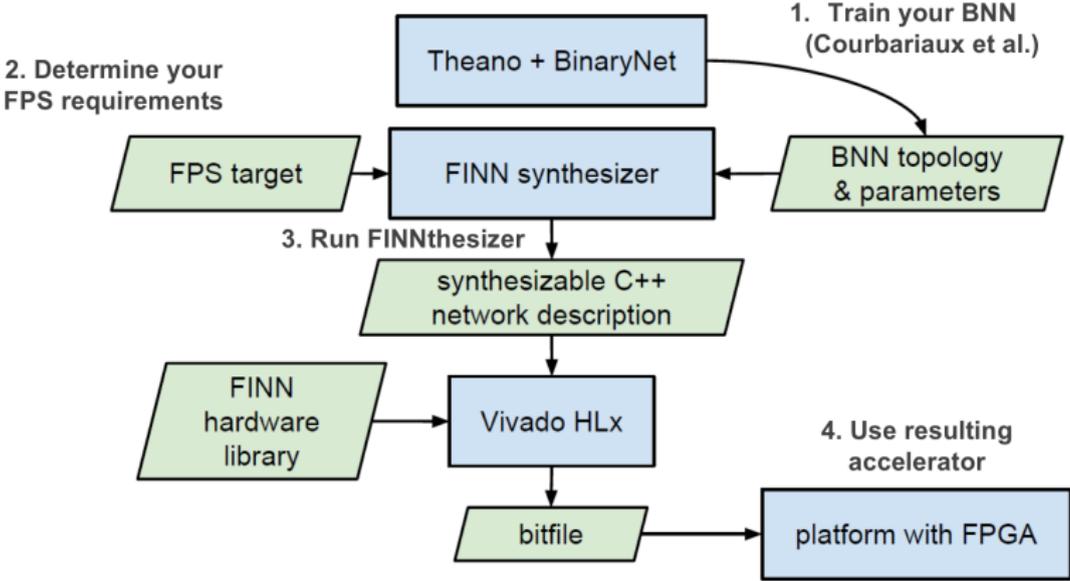
↑ 30x (On-chip weights)

### › **fast** inference with **large BNNs**

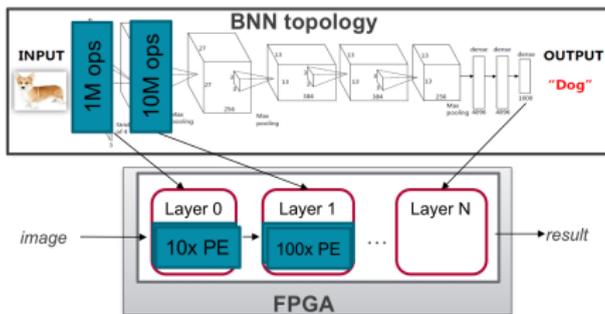
## Design Flow

- One size does not fit all - Generate tailored hardware for network and use-case
- Stay on-chip - Higher energy efficiency and bandwidth
- Support portability and rapid exploration - Vivado HLS (High-Level Synthesis)
- Simplify with BNN-specific optimizations - Exploit compile time optimizations to simplify hardware, e.g. batchnorm and activation => thresholding

# Design Flow



# Heterogeneous Streaming Architecture



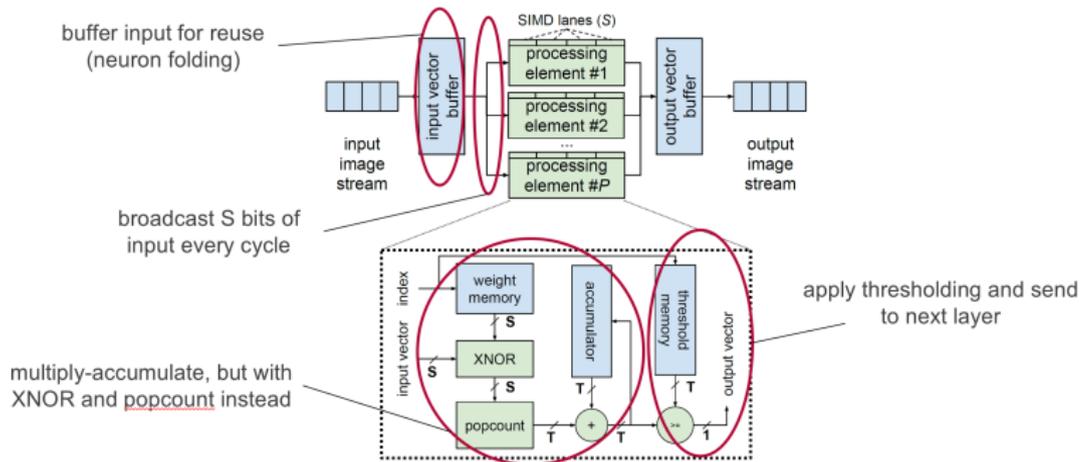
**1x FPS**  
**10x FPS**

- **One hardware layer per BNN layer, parameters built into bitstream**
  - Both inter- and intra-layer parallelism
- **Heterogeneous: Avoid “one-size-fits-all” penalties**
  - Allocate compute resources according to FPS and network requirements
- **Streaming: Maximize throughput, minimize latency**
  - Overlapping computation and communication, batch size = 1

Slide © Copyright 2016 Xilinx

# Matrix-Vector Threshold Unit (MVTU)

- Core computational element of FINN, tiled matrix-vector multiply
- Computes a (P) row x (S) column chunk of matrix every cycle, per-layer configurable tile size



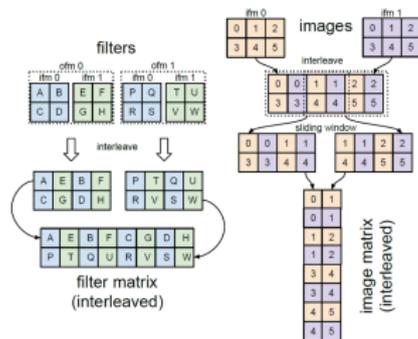
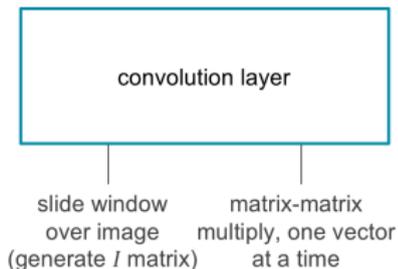
Slide © Copyright 2016 Xilinx

# Convolutional Layers

➤ Lower convolutions to matrix-matrix multiplication,  $W \cdot I$

- $W$  : filter matrix (generated offline)
- $I$  : image matrix (generated on-the-fly)

➤ Two components:



# Performance

	Accuracy	FPS	Power (chip)	Power (wall)	kFPS / Watt (chip)	kFPS / Watt (wall)	Precision	
FINN	MNIST, SFC-max	95.8%	12.3 M	7.3 W	21.2 W	1693	583	1
	MNIST, LFC-max	98.4%	1.5 M	8.8 W	22.6 W	177	269	1
	CIFAR-10, CNV-max	80.1%	21.9 k	3.6 W	11.7 W	6	2	1
	SVHN, CNV-max	94.9%	21.9 k	3.6 W	11.7 W	6	2	1
Prior Work	MNIST, <a href="#">Alemdar et al.</a>	97.8%	255.1 k	0.3 W	-	806	-	2
	CIFAR-10, <a href="#">TrueNorth</a>	83.4%	1.2 k	0.2 W	-	6	-	1
	SVHN, <a href="#">TrueNorth</a>	96.7%	2.5 k	0.3 W	-	10	-	1

Max accuracy loss: ~3%

10 – 100x better performance

CIFAR-10/SVHN energy efficiency comparable to TrueNorth ASIC

# Summary

- Reducing precision
  - Significantly reduce computational costs in DNNs
  - Data may now fit entirely on chip, avoiding external memory accesses
  - Computations greatly simplified
  - Key dimension for optimisation in CPU/GPU/FPGA implementations
- Convolutional layer can be computed as a MM
- Still an active research topic

# Tutorial Question 1

- 1 Download VM (quantisation\_usyd.ova) from <https://bluemountain.eee.hku.hk/papaa2018/>
- 2 Import to Virtualbox, and inside VM do  

```
git clone https://gitlab.com/phwl/syq-cifar10.git
```
- 3 Derive Equation (4) from (3) and (2)

## Tutorial Question 2

- 1 Cifar10 is a *very small* neural network benchmark<sup>1</sup>. Test precision with SYQ using:

```
cd syq-cifar10/src  
python cifar10_eval.py
```

(this can run during training)

- 2 The code provided performs binary quantisation. Modify the code to determine precision for binary, ternary and floating-point (use the checkpoint files provided to initialise your training).

```
python cifar10_train.py
```

---

<sup>1</sup> [https://www.tensorflow.org/tutorials/images/deep\\_cnn](https://www.tensorflow.org/tutorials/images/deep_cnn)

## References I

- [CBD14] Matthieu Courbariaux et al. “Low precision arithmetic for deep learning”. In: *CoRR* abs/1412.7024 (2014). arXiv: 1412.7024. URL: <http://arxiv.org/abs/1412.7024>.
- [CPS06] Kumar Chellapilla et al. “High Performance Convolutional Neural Networks for Document Processing”. In: *Tenth International Workshop on Frontiers in Handwriting Recognition*. Ed. by Guy Lorette. <http://www.suvisoft.com>. Université de Rennes 1. La Baule (France): Suvisoft, Oct. 2006. URL: <https://hal.inria.fr/inria-00112631>.

## References II

- [Jac+18] Benoit Jacob et al. “Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference”. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2018.
- [Far+18] Julian Faraone et al. “SYQ: Learning Symmetric Quantization For Efficient Deep Neural Networks”. In: *Proc. Computer Vision and Pattern Recognition (CVPR)*. Utah, US, June 2018.
- [Zho+16] Shuchang Zhou et al. “DoReFa-Net: Training Low Bitwidth Convolutional Neural Networks with Low Bitwidth Gradients”. In: *CoRR abs/1606.06160* (2016). URL: <http://arxiv.org/abs/1606.06160>.

## References III

- [Hub+16] Itay Hubara et al. “Quantized Neural Networks: Training Neural Networks with Low Precision Weights and Activations”. In: *CoRR* abs/1609.07061 (2016). arXiv: 1609.07061. URL: <http://arxiv.org/abs/1609.07061>.
- [Cai+17] Zhaowei Cai et al. “Deep Learning with Low Precision by Half-wave Gaussian Quantization”. In: *CoRR* abs/1702.00953 (2017). arXiv: 1702.00953. URL: <http://arxiv.org/abs/1702.00953>.
- [Ras+16] Mohammad Rastegari et al. “XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks”. In: *CoRR* abs/1603.05279 (2016). URL: <http://arxiv.org/abs/1603.05279>.

## References IV

- [Mel+17] Naveen Mellempudi et al. “Ternary Neural Networks with Fine-Grained Quantization”. In: *CoRR* abs/1705.01462 (2017). URL: <http://arxiv.org/abs/1705.01462>.
- [Zhu+16] Chenzhuo Zhu et al. “Trained Ternary Quantization”. In: *CoRR* abs/1612.01064 (2016). URL: <http://arxiv.org/abs/1612.01064>.