# Low-precision CNN Inference and Training for Edge Applications on FPGAs

Philip Leong | Computer Engineering Laboratory
School of Electrical and Information Engineering,
The University of Sydney

THE UNIVERSITY OF
SYDNEY

› Focuses on how to use parallelism to solve demanding problems

- Novel architectures, applications and design techniques using VLSI, FPGA and parallel computing technology

› Research

- Reconfigurable computing

- Machine learning

- Signal processing

› Collaborations

- Xilinx, Exablaze (now Cisco)

- Defence and DSTG

- clustertech.com

› Implementation of *small* CNNs

- Unrolling Ternary Networks

- Training deep neural networks in low-precision with high accuracy using FPGAs

# Unrolling Ternary Networks

*Stephen Tridgell, Martin Kumm, Martin Hardieck, David Boland, Duncan Moss, Peter Zipf, Philip H.W. Leong*
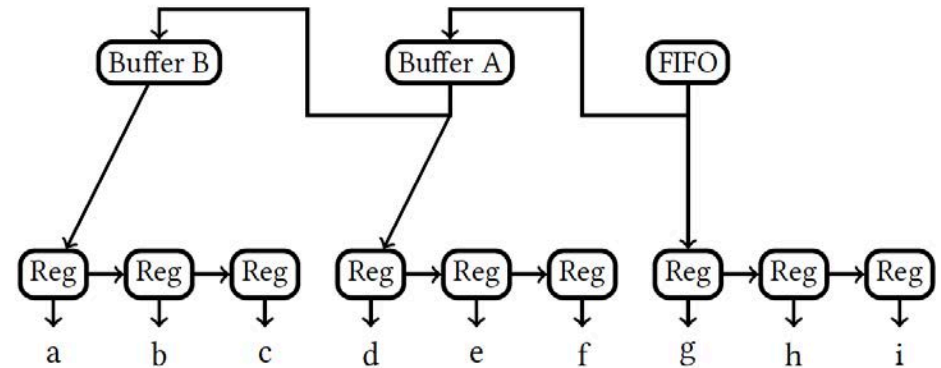
THE UNIVERSITY OF
SYDNEY

› **A fully pipelined DNN** implementation with ternary coefficients

› Difficult to make fully parallel implementations of a NN on contemporary FPGA due to size

› Fit entire DNN on FPGA by exploiting unstructured sparsity and the following techniques:

1. Buffering of streaming inputs in a pipelined manner

2. Ternary weights implemented as pruned adder trees

3. Common subexpression merging

4. 16-bit bit serial arithmetic to minimize accuracy loss with low area

5. Sparsity control

› VGG-7 network

› Ternary weights

› 16-bit activations

› Accept a single pixel every cycle (p=1)

- W*W image takes W*W cycles

| Operation | Image Size In | Channel In | Channel Out |
|---|---|---|---|
| Buffer | 32x32 | 3 | 3 |
| Conv | 32x32 | 3 | 64 |
| Scale and Shift | 32x32 | 64 | 64 |
| Buffer | 32x32 | 64 | 64 |
| Conv | 32x32 | 64 | 64 |
| Scale and Shift | 32x32 | 64 | 64 |
| Buffer | 32x32 | 64 | 64 |
| Max Pool | 32x32 | 64 | 64 |
| Buffer | 16x16 | 64 | 64 |
| Conv | 16x16 | 64 | 128 |
| Scale and Shift | 16x16 | 128 | 128 |
| **Buffer** | 16x16 | 128 | 128 |
| **Conv** | 16x16 | 128 | 128 |
| Scale and Shift | 16x16 | 128 | 128 |
| Buffer | 16x16 | 128 | 128 |
| Max Pool | 16x16 | 128 | 128 |
| Buffer | 8x8 | 128 | 128 |
| Conv | 8x8 | 128 | 256 |
| Scale and Shift | 8x8 | 256 | 256 |
| Buffer | 8x8 | 256 | 256 |
| Conv | 8x8 | 256 | 256 |
| Scale and Shift | 8x8 | 256 | 256 |
| Buffer | 8x8 | 256 | 256 |
| Max Pool | 8x8 | 256 | 256 |
| FIFO | 4x4 | 256 | 256 |
| MuxLayer | 4x4 | 256 | 4096 |
| Dense | 1x1 | 4096 | 128 |
| Scale and Shift | 1x1 | 128 | 128 |
| MuxLayer | 1x1 | 128 | 128 |
| Dense | 1x1 | 128 | 10 |

## Implement Pipelined 3x3 Convolution



Input FIFO outputs the pixel each cycle to both Buffer A and the first stage of a shift register.
Buffer A and Buffer B delay the output by the image width

› Weights are ternary

- So multiplication with $\pm 1$ is either addition or subtraction
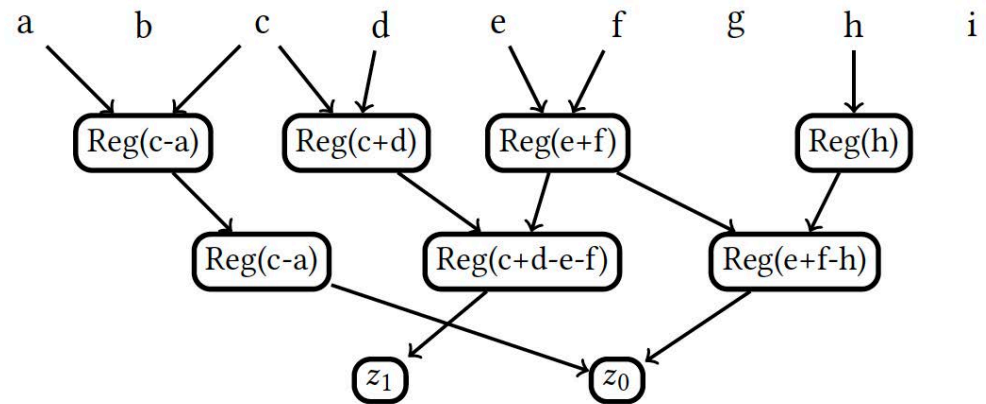
- Multiplication with 0 makes matrix sparse

$$a \times (-1) \qquad b \times 0 \qquad c \times 1$$
$$d \times 0 \qquad e \times 1 \qquad f \times 1$$
$$g \times 0 \qquad h \times (-1) \qquad i \times 0$$

› Weights are ternary

- Reduces convolution to constructing adder tree

- Subexpression merged to reduce implementation

$$
\begin{array}{lll}
a \times (-1) & b \times 0 & c \times 1 \\
d \times 0 & e \times 1 & f \times 1 \\
g \times 0 & h \times (-1) & i \times 0
\end{array}
$$



Computing $z_0 = c + e + f - (a + h)$ and $z_1 = c + d - e - f$

› None – doesn't fit

› Reduced Pipelined Adder Graph (RPAG) – too slow

› Top down CSE (TD-CSE) – doesn't find good solutions

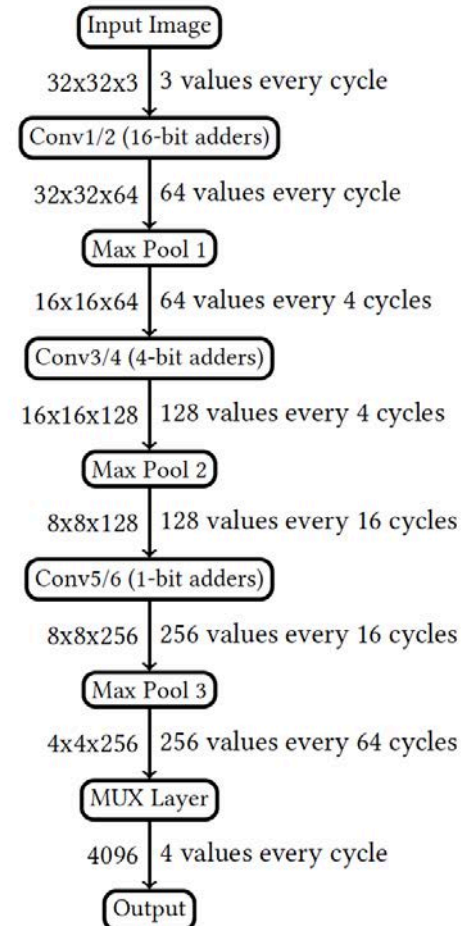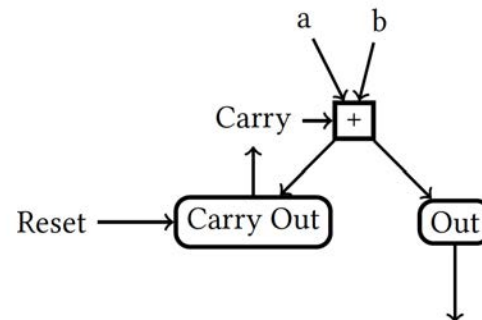› Bottom up CSE (BU-CSE) – worst of both worlds, runs out of memory

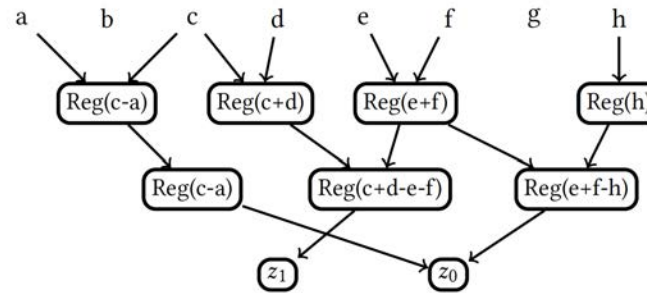› RPAG too computationally expensive for layers 2-6

› BU-CSE a bit better than TD-CSE

| Layer | Method | Adds+Regs | Time(s) | Mem(GB) | P&R(hrs) |
|-------|--------|-----------|---------|---------|----------|
| 1 | None | 868 | - | - | 0.5 |
|   | RPAG | **482** | 64 | 0.008 | 0.48 |
|   | TD-CSE | 599 | 0.4 | 0.029 | - |
|   | BU-CSE | 616 | 0.5 | 0.03 | 0.45 |
| 2 | None | 8681 | - | - | 1.08 |
|   | TD-CSE | 5299 | 24 | 0.1 | - |
|   | BU-CSE | **4544** | 64 | 0.17 | 0.93 |
| 3 | None | 17972 | - | - | 1.9 |
|   | TD-CSE | 10765 | 89 | 0.18 | - |
|   | BU-CSE | **10370** | 545 | 1.1 | 1.13 |
| 4 | None | 36741 | - | - | 4.25 |
|   | TD-CSE | 21357 | 873 | 0.63 | - |
|   | BU-CSE | **20365** | 2937 | 6.6 | 2.68 |
| 5 | None | 72248 | - | - | 3.86 |
|   | TD-CSE | 39659 | 3088 | 1.2 | - |
|   | BU-CSE | **39135** | 25634 | 44 | 1.72 |
| 6 | None | 146083 | - | - | 11.15 |
|   | TD-CSE | 76505 | 26720 | 4.8 | - |
|   | BU-CSE | **73935** | 147390 | 191.0 | 3.08 |

| Layer | % decrease in Adds+Regs | % decrease in CLBs |
|-------|-------------------------|---------------------|
| 1 | -29.0 | -41.4 |
| 2 | -47.7 | -32.6 |
| 3 | -42.3 | -42.1 |
| 4 | -44.6 | -44.6 |
| 5 | -45.8 | -58.8 |
| 6 | -49.4 | -60.8 |

› Used 16-bit fixed point

› Each layer followed by batch normalization with floating point scaling factor

› Suppose that for a given layer, $p$ pixels arrive at the same time

- For $p \geq 1$ have $p$ adder trees in parallel

- For $p < 1$ word or bit-serial adders can match input rate with hardware resources

- 4-bit digit serial has 1/4 area

- 1-bit bit serial has 1/16 area

› Avoids idle adders





| | |
|---|---|
| Input Image | |
| 32x32x3 | 3 values every cycle |
| Conv1/2 (16-bit adders) | |
| 32x32x64 | 64 values every cycle |
| Max Pool 1 | |
| 16x16x64 | 64 values every 4 cycles |
| Conv3/4 (4-bit adders) | |
| 16x16x128 | 128 values every 4 cycles |
| Max Pool 2 | |
| 8x8x128 | 128 values every 16 cycles |
| Conv5/6 (1-bit adders) | |
| 8x8x256 | 256 values every 16 cycles |
| Max Pool 3 | |
| 4x4x256 | 256 values every 64 cycles |
| MUX Layer | |
| 4096 | 4 values every cycle |
| Output | |

› CIFAR10 dataset

› Image padded with 4 pixels each side and randomly cropped back to 32x32

› Weights are compared with threshold $\Delta^* \approx \epsilon \cdot E(|W|)$

 - 0 if less than threshold, $s(\pm 1)$ otherwise (s is a scaling factor)

› We introduce the idea of changing $\epsilon$ to control sparsity

| TNN Type | $\epsilon$ | Sparsity ( % ) | Accuracy |
|---|---|---|---|
| Graham [Graham 2014] (Floating Point) | - | - | 96.53% |
| Li et al. [Li et al. 2016], full-size | 0.7 | $\approx 48$ | 93.1% |
| Half-size | 0.7 | $\approx 47$ | 91.4% |
| Half-size | 0.8 | $\approx 52$ | 91.9% |
| Half-size | 1.0 | $\approx 61$ | 91.7% |
| Half-size | 1.2 | $\approx 69$ | 91.9% |
| Half-size | 1.4 | $\approx 76$ | 90.9% |
| Half-size | 1.6 | $\approx 82$ | 90.3% |
| Half-size | 1.8 | $\approx 87$ | 90.6% |

| Layer Type | Input Image Size | Num Filters | $\epsilon$ | Sparsity |
|---|---|---|---|---|
| Conv2D | $32 \times 32 \times 3$ | 64 | 0.7 | 54.7% |
| Conv2D | $32 \times 32 \times 64$ | 64 | 1.4 | 76.9% |
| Max Pool | $32 \times 32 \times 64$ | 64 | - | - |
| Conv2D | $16 \times 16 \times 64$ | 128 | 1.4 | 76.1% |
| Conv2D | $16 \times 16 \times 128$ | 128 | 1.4 | 75.3% |
| Max Pool | $16 \times 16 \times 128$ | 128 | - | - |
| Conv2D | $8 \times 8 \times 128$ | 256 | 1.4 | 75.8% |
| Conv2D | $8 \times 8 \times 256$ | 256 | 1.4 | 75.4% |
| Max Pool | $8 \times 8 \times 256$ | 256 | - | - |
| Dense | 4096 | 128 | 1.0 | 76.2% |
| Softmax | 128 | 10 | 1.0 | 58.4% |

› System implemented on Ultrascale+ VU9P @ 125 MHz

› Open Source Verilog generator

- https://github.com/da-steve101/binary_connect_cifar

› Generated code using in AWS F1 implementation

- https://github.com/da-steve101/aws-fpga

| Block | LUTs/1182240 | FFs/2364480 |
|---|---|---|
| Conv1 | 3764 ( 0.3% ) | 10047 ( 0.4% ) |
| Conv2 | 40608 ( 3.4% ) | 71827 ( 3.0% ) |
| Conv3 | 55341 ( 4.7% ) | 56040 ( 2.4% ) |
| Conv4 | 111675 ( 9.4% ) | 110021 ( 4.7% ) |
| Conv5 | 73337 ( 6.2% ) | 79233 ( 3.4% ) |
| Conv6 | 127932 ( 10.8% ) | 139433 ( 5.9% ) |
| All Conv | 535023 ( 45.3% ) | 631672 ( 26.7% ) |
| Dense | 12433 ( 1.1% ) | 19295 ( 0.8% ) |
| SM | 500 ( 0.04% ) | 442 ( 0.02% ) |
| Whole CNN | 549358 ( 46.5% ) | 659252 ( 27.9% ) |
| Whole design | 787545 ( 66.6% ) | 984443 ( 41.6% ) |

# Summary of Sparsity and CSE Improvement

| Layer | Num Mults | Num Mults | With Sparsity | With CSE |
|---|---|---|---|---|
| Conv1 | 32*32*3*3*3*64 | 1769472 | 716800 | 630784 |
| Conv2 | 32*32*3*3*64*64 | 37748736 | 8637440 | 4653056 |
| Conv3 | 16*16*3*3*64*128 | 18874368 | 4559616 | 2654720 |
| Conv4 | 16*16*3*3*128*128 | 37748736 | 9396480 | 5213440 |
| Conv5 | 8*8*3*3*128*256 | 18874368 | 4656768 | 2504640 |
| Conv6 | 8*8*3*3*256*256 | 37748736 | 9356736 | 4731840 |
| Dense | 4096*128 | 524228 | 524228 | 1048456[1] |
| SM | 128*10 | 1280 | 1280 | 2560[1] |
| Total | 153289924 | 153 MMACs/Image | 38 MMACs/Image | 21 MOps/Image |

[1] Obtained by converting one MACs to two Ops

## Comparison with ASIC and FPGA implementations (CIFAR10)

| Reference | Hardware ($mm^2$,nm,LE[5]/LC[5] $\times 10^6$) | Precision (wghts, actv) | Freq. [MHz] | Latency | TOps/sec A/L/E[6] | FPS | Accuracy |
|---|---|---|---|---|---|---|---|
| [Venkatesh et al. 2017] | ASIC(1.09,14,–) | $(2,16^2)$ | 500 | – | 2.5/2.5/2.5 | – | 91.6%[3] |
| [Andri et al. 2017] | ASIC(1.9,65,–) | (1,12) | 480 | – | 1.5/1.5/1.5 | 434 | – |
| [Jouppi et al. 2017] | ASIC(331,28,–) | (8,8) | 700 | ≈10 ms | 86/86/86[4] | – | – |
| [Baskin et al. 2018] | 5SGSD8(1600,28,0.7) | (1,2) | 105 | – | – | 1.2 k[3] | 84.2% |
| [Li et al. 2017] | XC7VX690(1806.25,28,0.7) | $(1^1, 1)$ | 90 | – | 7.7/3.9/7.7 | 6.2 k | 87.8% |
| [Liang et al. 2018] | 5SGSD8(1600,28,0.7) | (1,1) | 150 | – | 9.4/4.7/9.4 | 7.6 k[3] | 86.31% |
| [Prost-Boucle et al. 2017] | VC709(1806.25,28,0.7) | (2,2) | 250 | – | 8.4/4.2/8.4 | 27 k | 86.7% |
| [Umuroglu et al. 2017] | ZC706(961,28,0.35) | (1,1) | 200 | 283 µs | 2.4/1.2/2.4 | 21.9 k | 80.1% |
| [Fraser et al. 2017] | KU115(1600,20,1.45) | (1,1) | 125 | 671 µs | 14.8/7.4/14.8 | 12 k | 88.7% |
| This work | VU9P(2256.25,20,2.6) | (2,16) | 125 | 29 µs | 2.5/2.5/37.3 | 122k | 90.9% |

[1]First layer is fixed point, [2]floating point, [3]estimated, [4] 92 TOps/sec peak, [5] LE and LC are from Xilinx or Altera documentation of the FPGAs, [6] Actual/Logical/Equivalent

› Presented method to unroll convolution with ternary weights and make parallel implementation

- Exploits unstructured sparsity with no overhead

- Uses CSE, sparsity control and digit serial adders to further reduce area

- Limited amount of buffering and only loosely dependent on image size

› As larger FPGAs become available this technique may become more favourable

# Training deep neural networks in low-precision with high accuracy using FPGAs

*Sean Fox, Julian Faraone, David Boland, Kees Vissers, and Philip H.W. Leong*

› How do we get deep learning to the edge?

- Train on GPUs

- Move trained model to the edge

(and make use of fast, low power FPGA inference engines?)

- .

= 1 billion devices

| <100M servers | 3B phones | 12B IoT | 150B embedded devices |

source: https://heartbeat.fritz.ai/

› How do we get deep learning to the edge?

- Train on GPUs

- Move trained model to the edge

What about changing environments?

We need to be able to train at the edge



= 1 billion devices

To the edge!

<100M servers

3B phones

12B IoT

150B embedded devices

source: https://heartbeat.fritz.ai/

- Forward path (Inference)



- Each layer computes:
  - Linear transform (Matrix operation)
  - Non-linear activation function

$$f(\boldsymbol{x}) = g(\boldsymbol{W}^T \boldsymbol{x} + \boldsymbol{b})$$

- Backward path:

$$\nabla L(.) = \frac{\partial Loss}{\partial y_{pred}}$$

- Step 1: Weight gradient

$$\nabla W_2 = \frac{\partial Loss}{\partial W_2} = X_2^T . \nabla L(.)$$

- Step 2: Activation gradient

$$\nabla X_2 = \frac{\partial Loss}{\partial X_2} = W_2 . \nabla L(.)$$

- Backward path



- Two additional Matrix Mults at each layer:
  - Weight gradients
  - Activation gradients (or loss back-propagated)

$$\nabla W_l = \frac{\partial Loss}{\partial W_l} = X_l^T . \nabla X_{l+1}$$

$$\nabla X_l = \frac{\partial Loss}{\partial X_l} = W_l . \nabla X_{l+1}$$

› DNN training is dominated by matrix multiplication

› GPUs good at floating-point matrix operations

   - Large dynamic range and avoids excessive quantisation errors critical for DNN training

**Algorithm 1: Convolution Layer**

**Define:** layer $l$; time $t$; 8-bit weights $\bar{W}_l^t$; input activations $x_l^t$; deltas $\nabla x_l^t$; weight updates $\nabla W_l^t$; quantisation functions $Q_w, Q_a, Q_e$; quantisation scaling coefficients $qw, qa, qe$; gemm inputs $A, B$; gemm output $C$; batch size $K$;

**1. Forward:**

    Software:

1        $\bar{x}_l^t, qa = Q_a(x_l^t); B = im2col(\bar{x}_l^t);$

    Hardware:

2        $A = (\bar{W}_l^t)^T;$

3        $C = tofloat(gemm(A, B), qw, qa);$

    Software:

4        $x_{l+1}^t = C;$

**2. Backward:**

    Software:

5        $\nabla \bar{x}_{l+1}^t, qe = Q_e(\nabla x_l^t); tmp = im2col(\bar{x}_l^{t\,T});$

6        for $i = 1, 2, \ldots, K$ do

    Hardware:

7        $A = \nabla \bar{x}_l^t(i)^T; B = tmp(i);$

8        $C = tofloat(gemm(A, B), qe, qa);$

    Software:

9        $\nabla W_l^t += C;$

10      end

    Hardware:

11      $A = \bar{W}_l^t; B = \nabla \bar{x}_{l+1}^t;$

12      $C = tofloat(gemm(A, B), qw, qe);$

    Software:

13      $\nabla x_l^t = col2im(C);$

**FPGA**

- Low-Precision (8-bit)
  - All matrix multiplications
  - >95% of DNN operations

**ARM**

- High-Precision
  - Everything else!
  - Of particular importance is the **weight update and gradient accumulator**

- Suits a Zynq platform
  - Fast DDR, shared between PL and floating-point

- Weight Update

$$W^{t+1} = W^t - \alpha \nabla W^t$$

If learning rate $\alpha$ and weight gradient $\nabla W^t$ are small, then the update will be rounded to 0

- Weight Gradient Accumulator

$$\nabla W = \sum_{i=1}^{batch\_size} x^T . \nabla x$$

If accumulator is low precision, then small numbers are rounded to 0

- Weight Update

$$W^{t+1} = W^t - \alpha \, \nabla W^t$$

If learning rate $\alpha$ and weight gradient $\nabla W^t$ are small, then the update will be rounded to 0

- Weight Gradient Accumulator

$$\nabla W = \sum_{i=1}^{batch\_size} x^T . \nabla x$$

High-precision accumulator can avoid underflow

- A group of numbers that share the same exponent

- Outcome: Dynamic Range
        (at a cost of some accuracy degradation)

**Float**   **Fixed**   **BFP**

We use 1 block for  weights, b for activations and deltas (where b is batch size)

› SWALP - averaging SGD iterates of the weights with a higher learning rate can recover quantisation errors

Representable points

● Representable points    ● LP-SGD Trajectory

● Representable points          ● LP-SGD Trajectory

● Representable points    ● LP-SGD Trajectory

Representable points    LP-SGD Trajectory

Representable points     LP-SGD Trajectory

● Representable points    ● LP-SGD Trajectory    ★ SWA Solution

- Heterogenous computing platform:
    1. Processing System (ARM processor)
    2. Programmable Logic (8-bit GEMM)
    3. Shared DRAM

- Well suited to standalone embedded applications

- Convolution must be transformed into Matrix Multiplication

- Convolution must be transformed into Matrix Multiplication

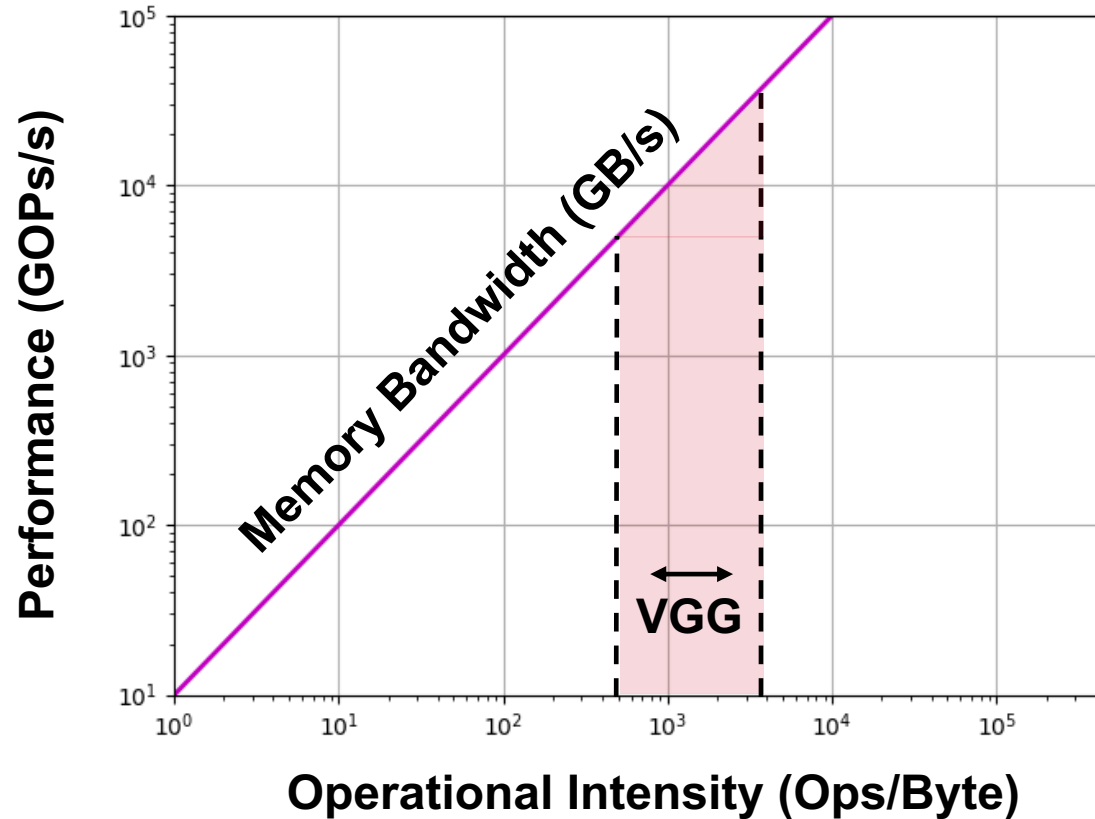- Convolution must be transformed into Matrix Multiplication

- 8-bit GEMM ($C = A.B$ or $C = A^T.B$)
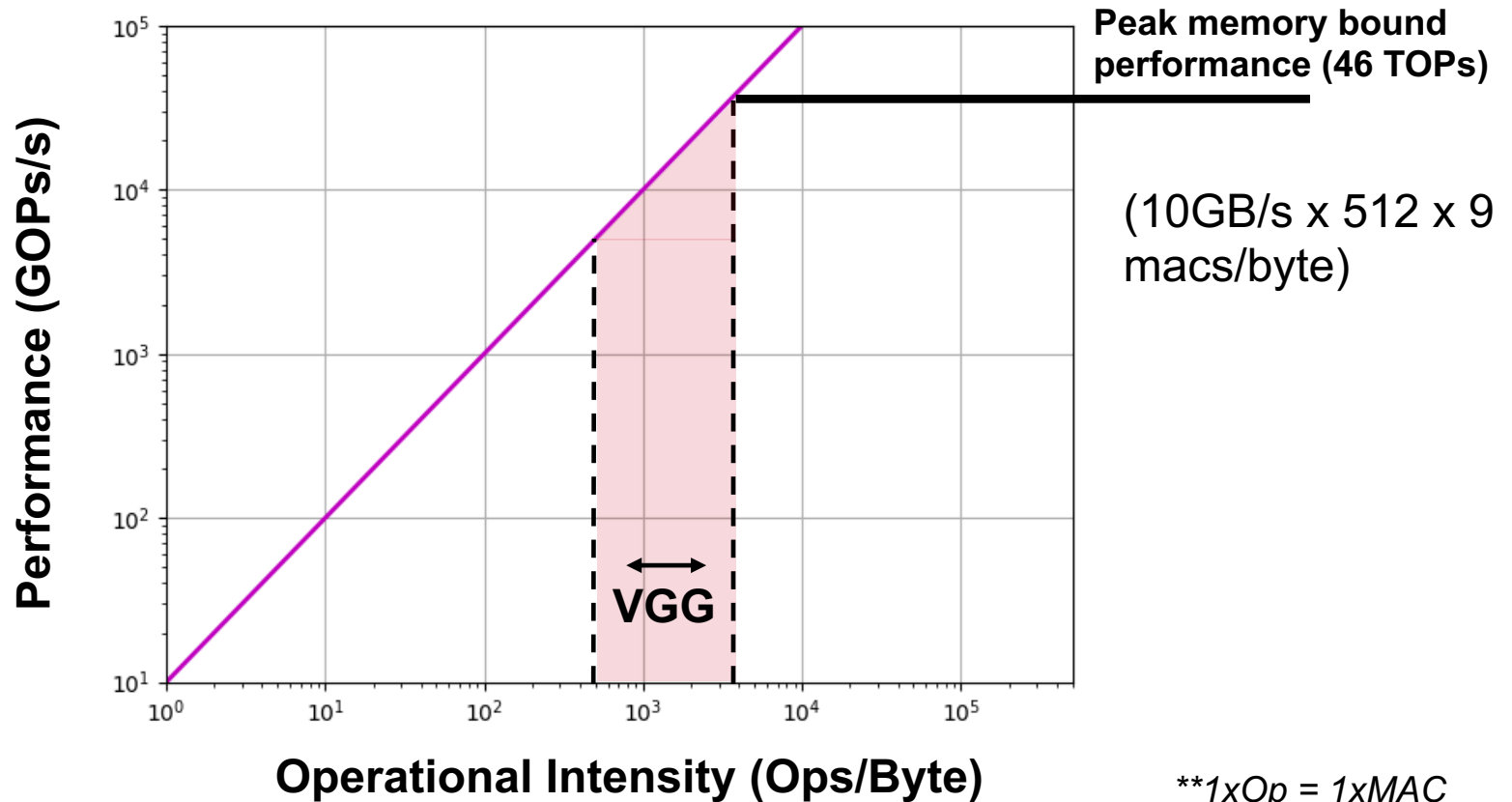
- Block Floating Point (BFP) rescaling
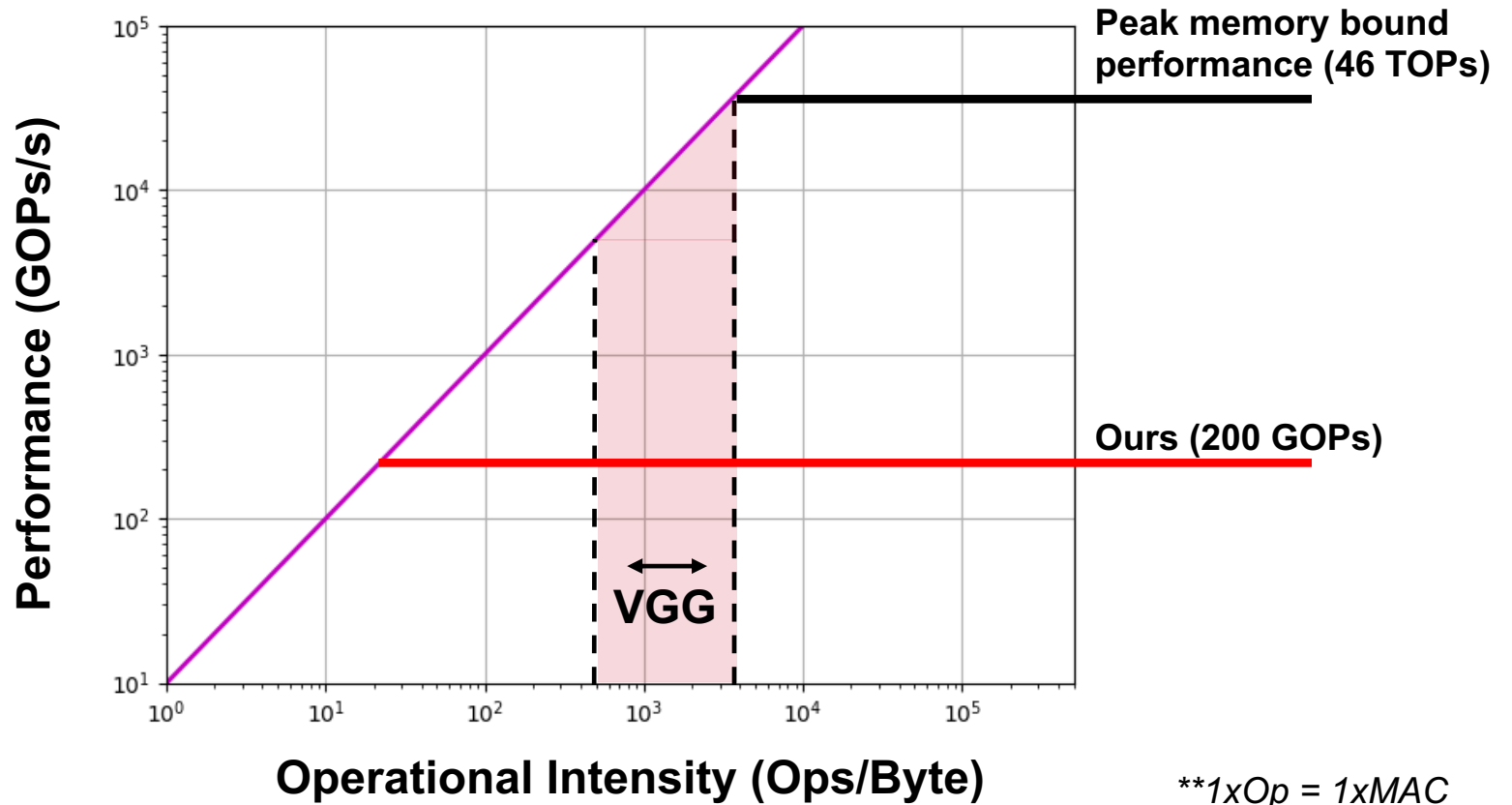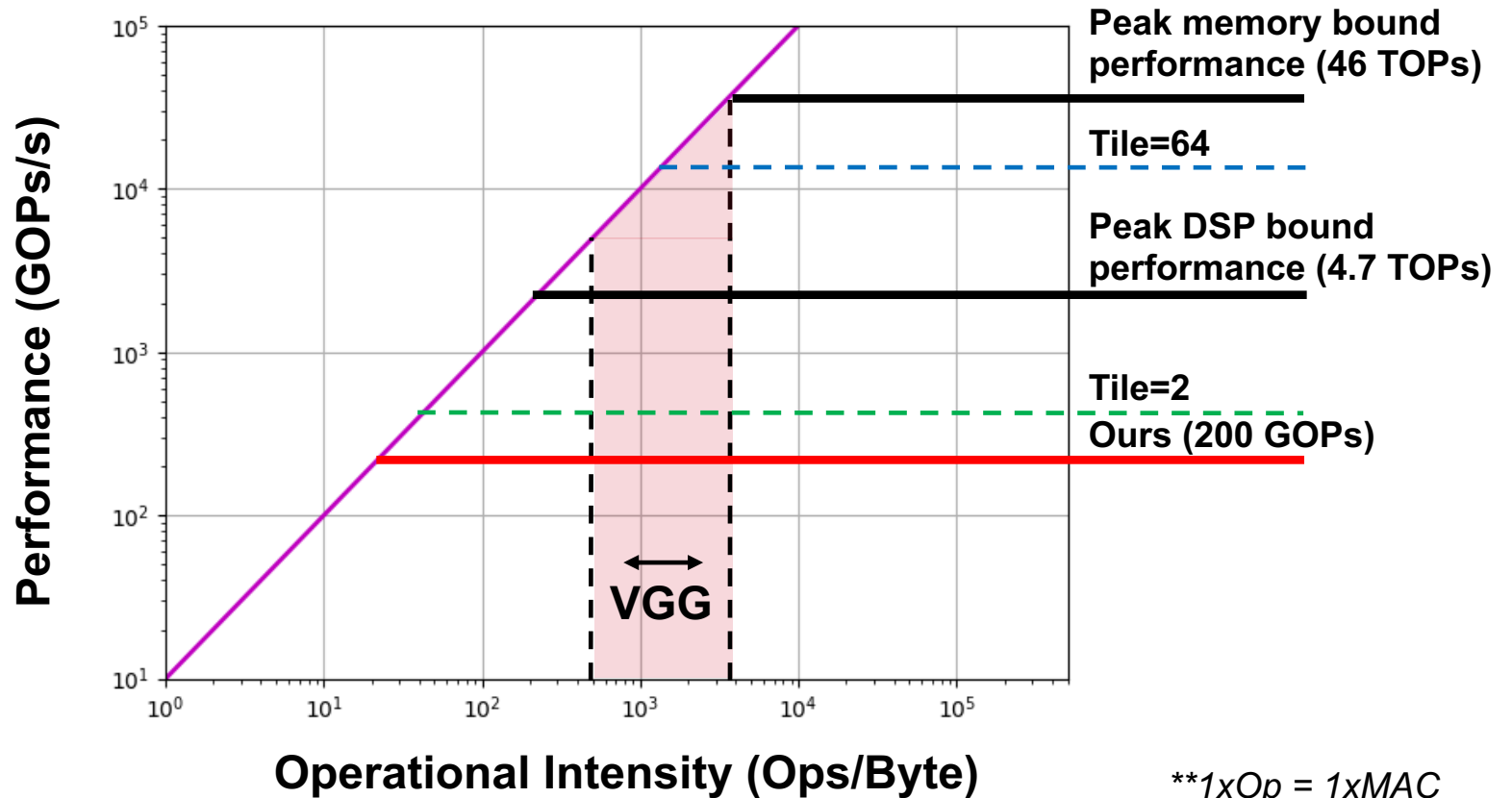
- Red line: 32x32 Array @ 200Mhz on ZCU111



**Operational Intensity (Ops/Byte)**

*Performance (GOPs/s)* — y-axis

*Memory Bandwidth (GB/s)*

VGG

**1xOp = 1xMAC*

- Red line: 32x32 Array @ 200Mhz on ZCU111



Peak memory bound performance (46 TOPs)

(10GB/s x 512 x 9 macs/byte)

**1xOp = 1xMAC*

- Red line: 32x32 Array @ 200Mhz on ZCU111



**Peak memory bound performance (46 TOPs)**

**Ours (200 GOPs)**

**VGG**

Performance (GOPs/s)

Operational Intensity (Ops/Byte)

*\*\*1xOp = 1xMAC*

- Red line: 32x32 Array @ 200Mhz on ZCU111



**Operational Intensity (Ops/Byte)**

*\*\*1xOp = 1xMAC*

- Software integration with Darknet (Yolo backend)
- Available: https://github.com/sfox14/darknet-zynq

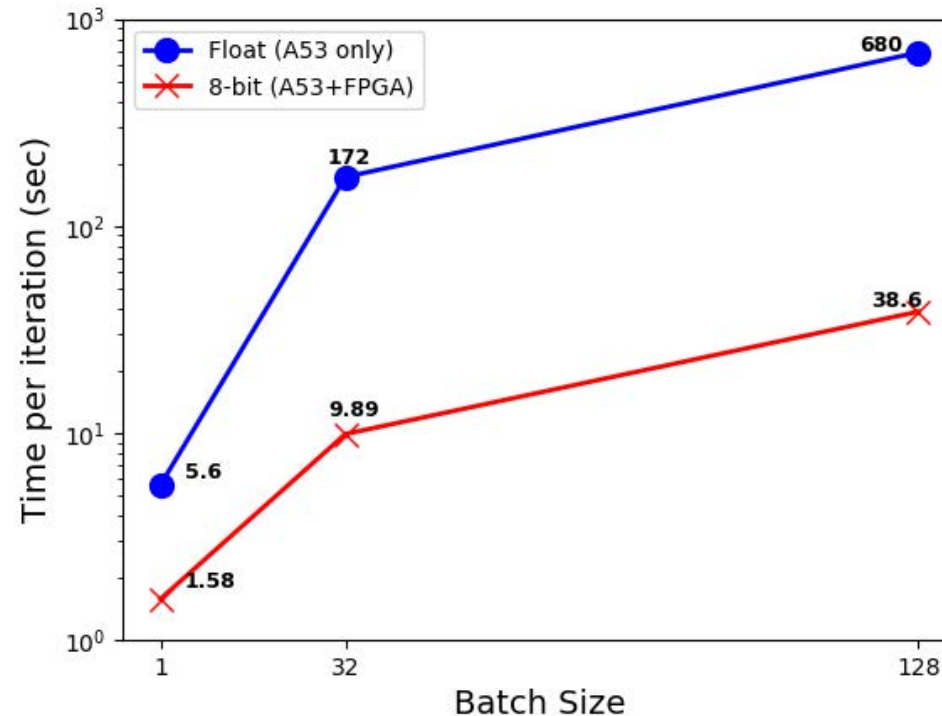Bitstream (.bit) and C/C++ runtime library (.so)

Open source C library for DNN training and inference

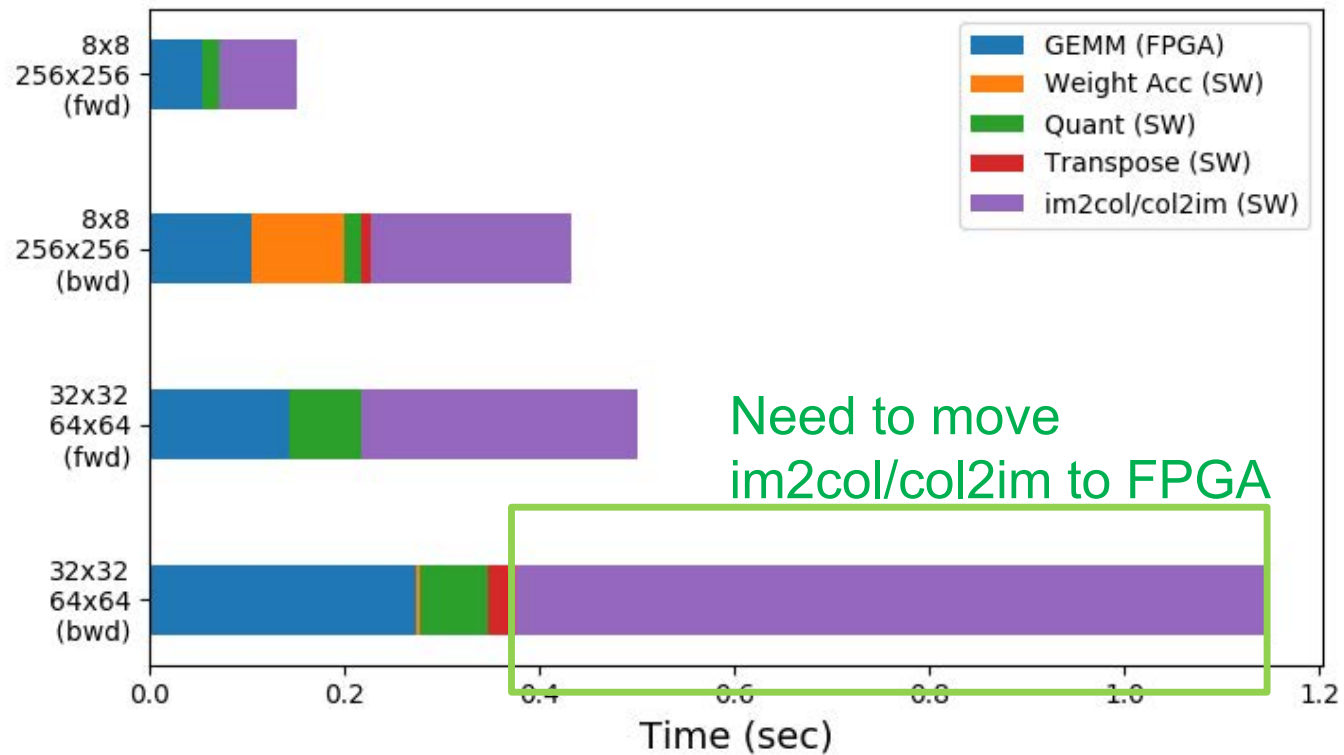Our code can be compiled and run on ZYNQ boards with PYNQ v2.3 image

| Dataset | Model | Float | SWALP | Ours (8-bit) |
|---------|-------|-------|-------|--------------|
| CIFAR10 | VGG16 | 93.02 | 92.47 | 92.93 |
|  | PreResNet-20 | 93.50 | 93.29 | 93.29 |
| MNIST | Logistic Regression | 92.60 | 92.06 | 92.70 |

❑ High accuracy can be sustained with mostly 8-bits
  • only requires a few more epochs than float
❑ Batch=128, no bias, no batchnorm
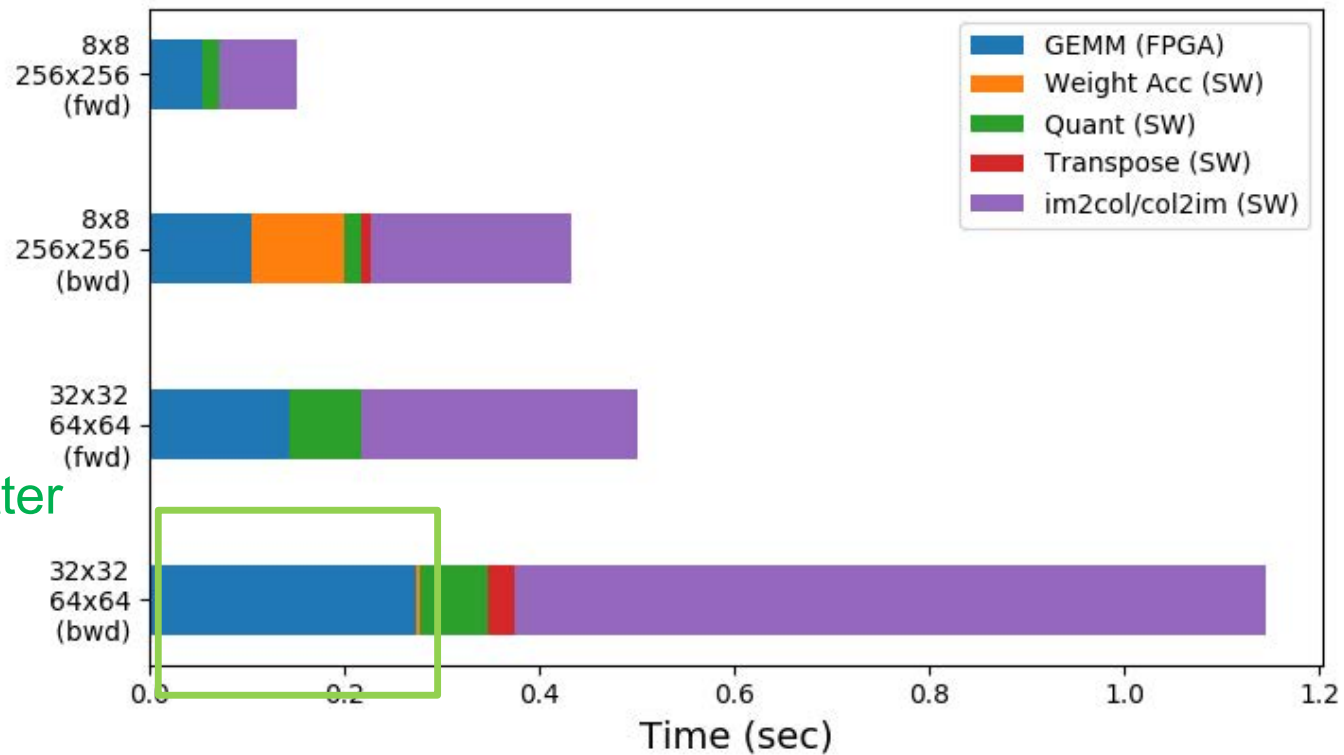❑ Lr=0.01, SWA applied to last 25% of epochs

- Good speed-up over ARM
- We don't exploit batch-level parallelism
- Much slower than a Tesla M40 GPU (600x) …

- Profiling VGG16 convolution layers (Pynq-Z1 and ZCU111)
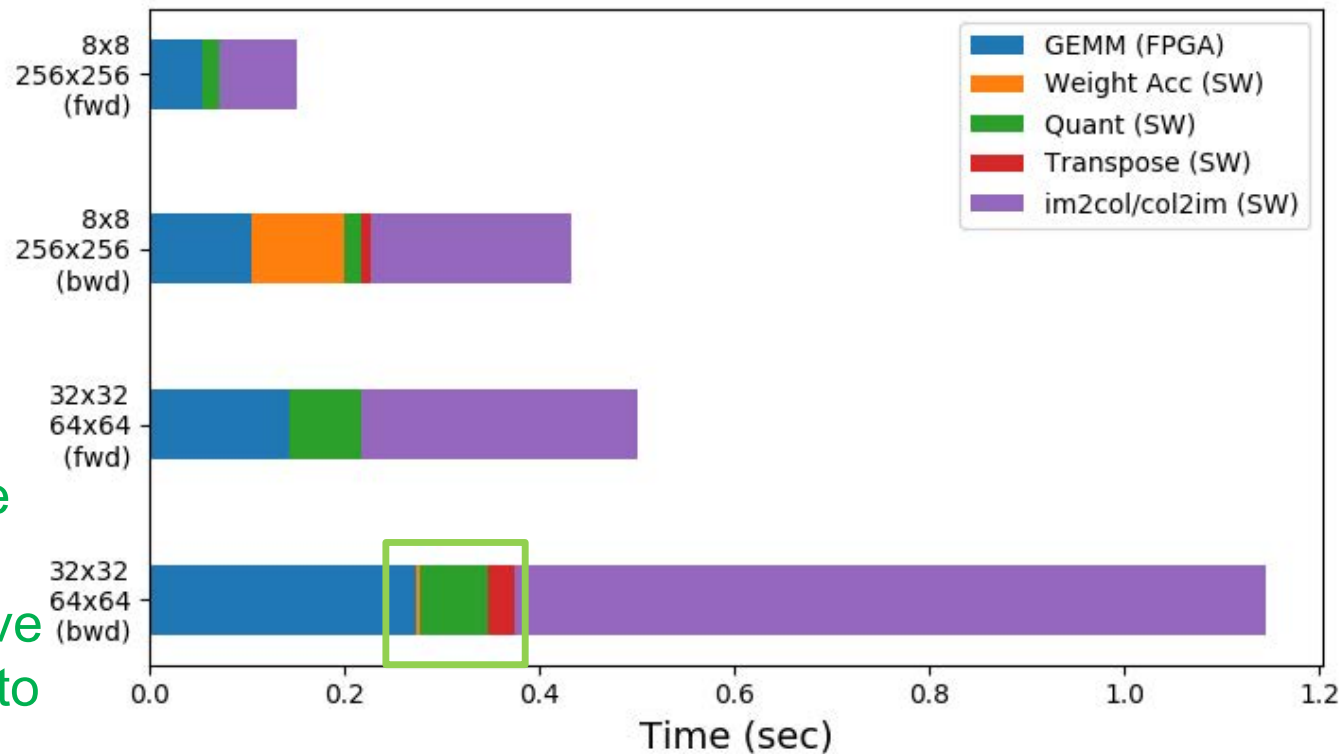


Need to move
im2col/col2im to FPGA

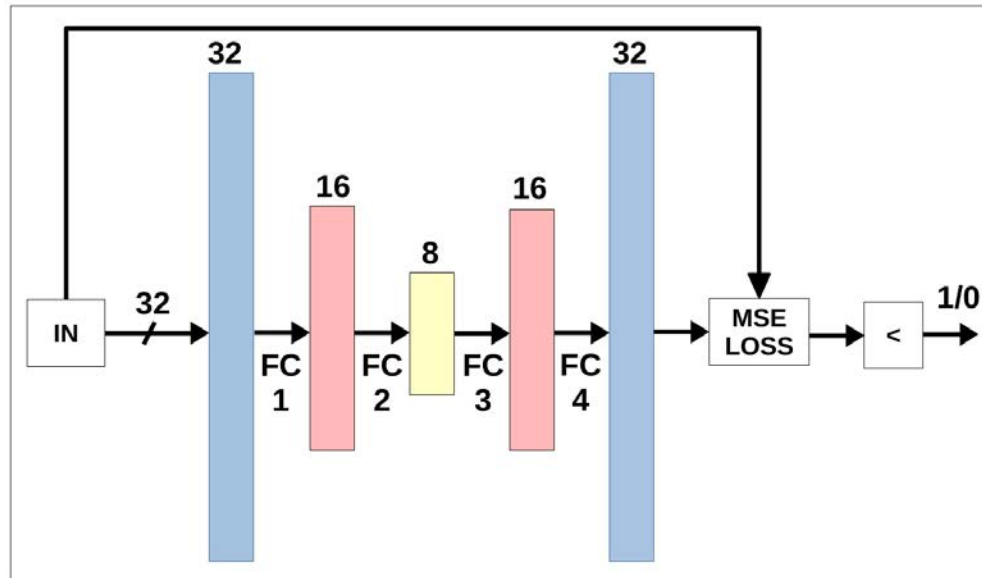- Profiling VGG16 convolution layers



Room for better
GEMM

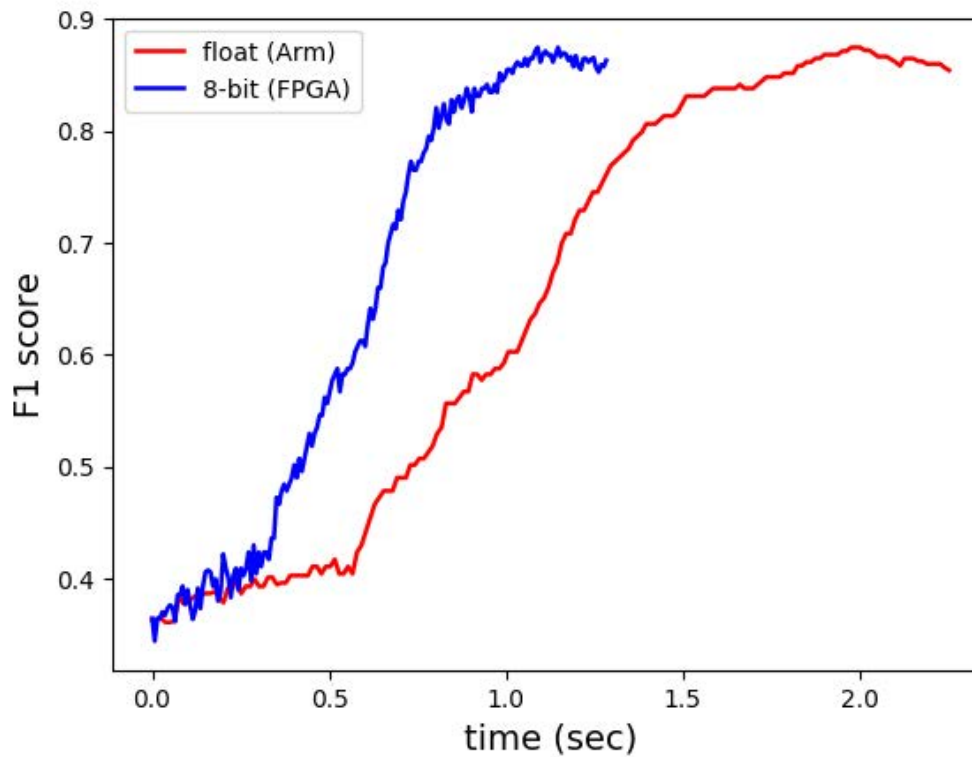- Profiling VGG16 convolution layers



Could update BFP less regularly/move quantisation to FPGA

- Trained a small Autoencoder network in low-precision

- Injected known anomalies into the validation set, and observed the time taken to train the model.

› FPGAs can be used for training at the edge!

› We demonstrate:

› Low-precision training techniques for high accuracy

› FPGA prototype with software integration

› We have highlighted room for optimisation:

- Move more computation to the FPGA

- Optimise matrix multiplication

- Lower precision

- Versal?

https://phwl.github.io/talks

THE UNIVERSITY OF
SYDNEY

[1] Stephen Tridgell, Martin Kumm, Martin Hardieck, David Boland , Duncan Moss, Peter Zipf, and Philip H. W. Leong. Unrolling ternary neural networks. *ACM Transactions on Reconfigurable Technology and Systems*, page to appear (accepted 30 Aug 2019), 2019.

[2] Sean Fox, Julian Faraone, David Boland, Kees Vissers, and Philip H.W. Leong. Training deep neural networks in low-precision with high accuracy using FPGAs. In Proc. International Conference on Field Programmable Technology (FPT), to appear. 2019.