

Efficient Deep Convolutional Neural Networks Accelerator without Multiplication and Retraining

Weihong Xu^{1,2}, Xiaohu You², Chuan Zhang^{1,2}

¹Lab of Efficient Architectures for Digital-communication and Signal-processing (LEADS)

²National Mobile Communications Research Laboratory, Southeast University, Nanjing, China

wh.xu@seu.edu.cn

April 17, 2018

Outline

1. Motivation
2. Related Work and Problem Formulation
3. Proposed Quantization and Hardware Co-design
4. Results and Analysis
5. Conclusion

Outline

1. Motivation
2. Related Work and Problem Formulation
3. Proposed Quantization and Hardware Co-design
4. Results and Analysis
5. Conclusion

Motivation

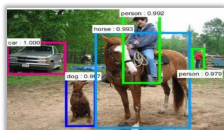
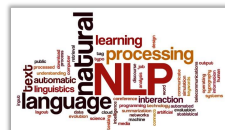


Image Segmentation



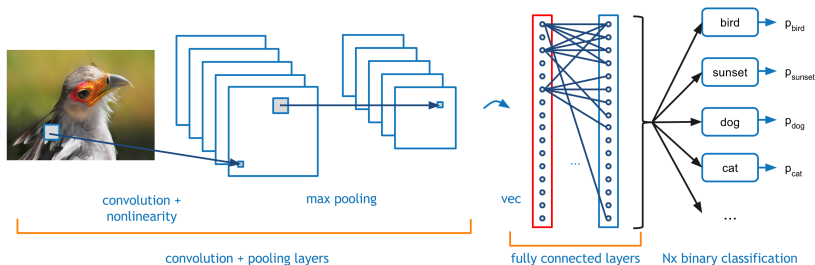
Game



NLP



Convolutional Neural Networks



- **Convolution:** feature extraction by convolving various filters over input image
- **Fully-connected:** linear transform over input features
- **Pooling and Non-linear:** perform down sampling and non-linear function

Major Challenges

- **Computation-intensive:** convolution takes up over 95% of overall complexity
 - $\mathcal{O}(N^2 K^2)$ complexity per image \rightarrow Prohibitive complexity
 - Floating point MAC is expensive \rightarrow Low energy efficiency
- **Memory-intensive:** FC layers contribute 90% parameters
 - Densely connected networks \rightarrow Millions of weights
 - Massive data movement \rightarrow Bandwidth limitation

Outline

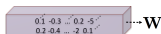
1. Motivation
2. Related Work and Problem Formulation
3. Proposed Quantization and Hardware Co-design
4. Results and Analysis
5. Conclusion

Low-precision Neural Networks

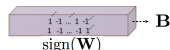
• Binarized Neural Networks

- Binary weights $\{-1, +1\}$ with scaling factor α
- Activation: 32-bit float
- α is determined by L_1 -norm of weights
- Accuracy degradation: **19%** on AlexNet

(1) Binarizing Weight

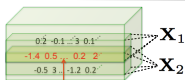


$$\frac{1}{n} \|W\|_{\ell_1} = \alpha$$



(2) Binarizing Input

Inefficient

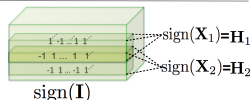


Redundant computations in overlapping areas

$$\frac{1}{n} \|X_1\|_{\ell_1} = \beta_1$$

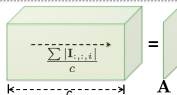
$$\frac{1}{n} \|X_2\|_{\ell_1} = \beta_2$$

K



(3) Binarizing Input

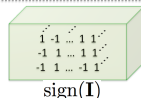
Efficient



$$A * K = \beta_1$$

$$A * K = \beta_2$$

K



¹[Rastegari, Ordonez, Redmon, et al., ECCV 2016]

Low-precision Neural Networks

- **Ternary Weight Nets**

- Ternary weights $\{-1, 0, +1\}$ with scaling factor α
- Activation: 32-bit float
- Adding zero value increases expressive abilities of weights
- Accuracy degradation: 3.7% on AlexNet

- **Objective of BNNs and TWNs**

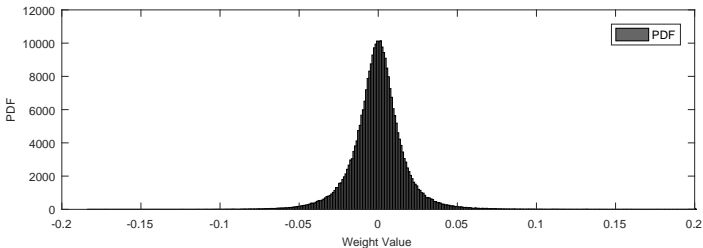
- Minimize distance between full precision weights \mathbf{W} and the ternary weights \mathbf{W}^t using scaling factor α :

$$\alpha^*, \mathbf{W}^{t*} = \arg \min_{\alpha, \mathbf{W}^t} \|\mathbf{W} - \alpha \mathbf{W}^t\|^2$$

¹[Li, Zhang, and Liu, *arXiv* 2016]

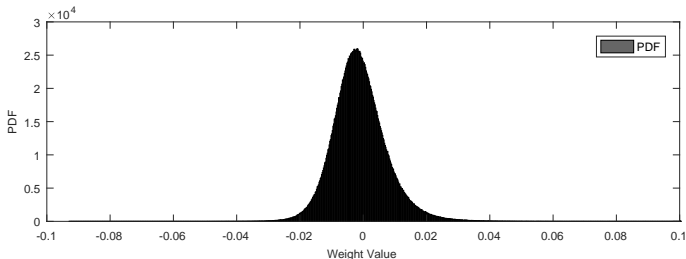
Non-Linear Quantization

- Distribution of weights in 5th layer of VGGNet



Non-Linear Quantization

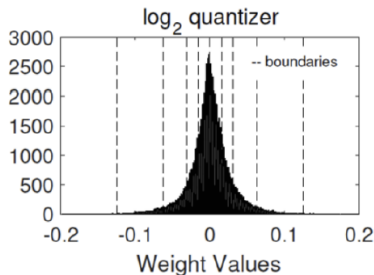
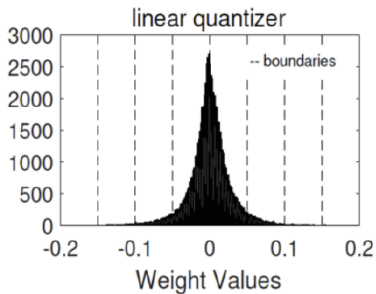
- Distribution of weights in 15th layer of VGGNet



- Near normal distribution
- Deeper layers tend to have smaller weights

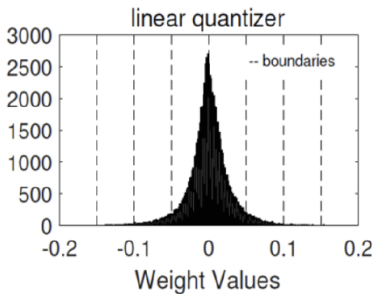
Non-Linear Quantization

- An intuitive perspective

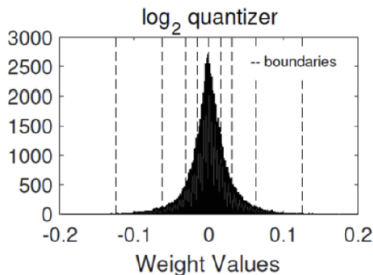


Non-Linear Quantization

- An intuitive perspective



Inefficient

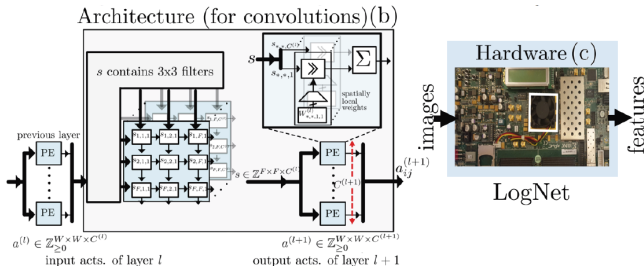


More efficient

Non-Linear Quantization

• LogNet

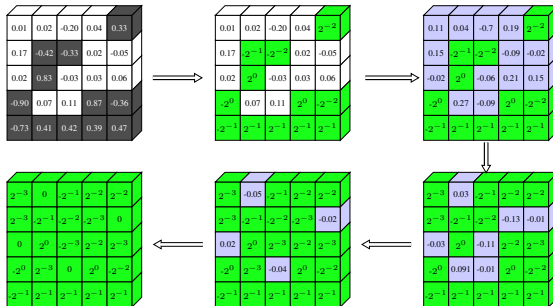
- Weights: 4-bit, Activation: 32-bit
- No scaling factor $\alpha \rightarrow$ Hardware friendly
- Substitute MAC with Shift and Add
- Accuracy degradation: 4.9% on AlexNet without Retraining
- Accuracy degradation: 4.6% on VGG16 with Retraining



¹[Lee, Miyashita, Chai, et al., ICASSP 2017]

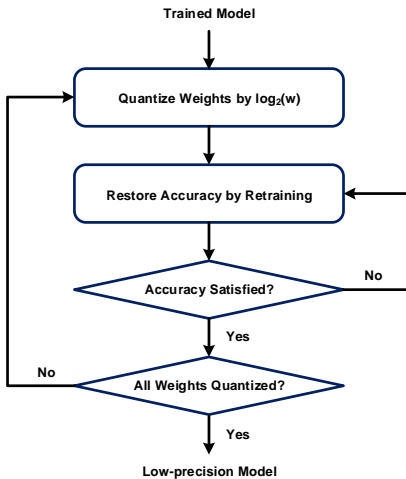
Non-Linear Quantization

- Incremental Network Quantization
 - Incremental retraining on Log domain
 - Weights: 5-bit, Activation: 4-bit
 - Accuracy degradation: **1.16%** on VGG16

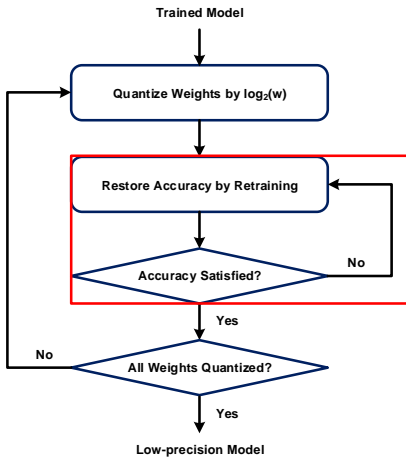


¹[Zhou, Yao, Guo, et al., ICLR 2017]

Problem Formulation

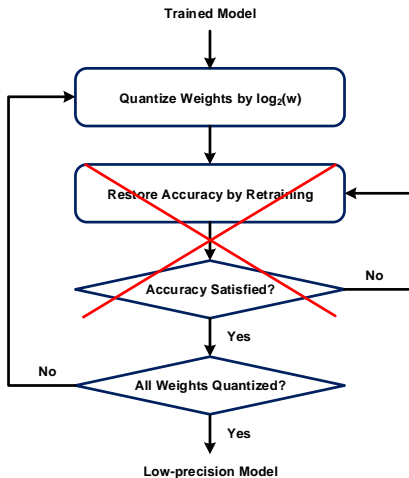


Problem Formulation



- Retraining is expensive!

Problem Formulation



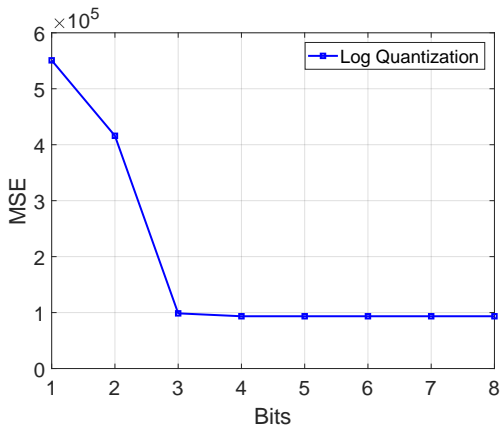
- How to skip retraining?

Outline

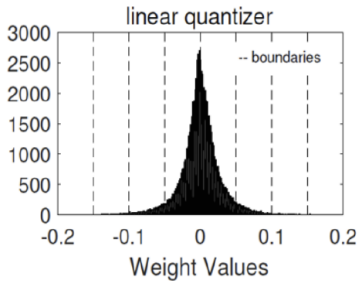
1. Motivation
2. Related Work and Problem Formulation
3. Proposed Quantization and Hardware Co-design
4. Results and Analysis
5. Conclusion

Non-uniform Quantization

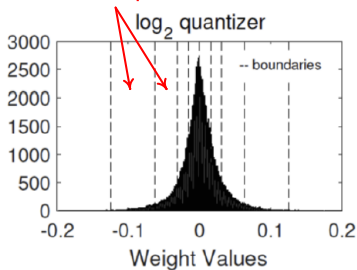
- More Log Bits \neq Less Quantization Error



Non-uniform Quantization



Unable to quantize



Proposed Non-uniform Quantization

- **Non-linear Quantization with Codebook**

$$\hat{w}_i = \sum_{n=1}^N \phi_n [\text{idx}_{i,n}]$$

- $\text{idx}_{i,n}$: i th segment of \hat{w}_i
- N codebooks

- **Codebook Structure**

$$\phi_n = \left[0, 2^{-1}, 2^{-2}, \dots, 2^{-(2^{B_n}-1)} \right]$$

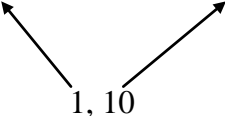
- **Quantize weights to codebook index idx**
- **Only process codebook index during inference**

Proposed Non-uniform Quantization

- **Example:** To quantize value 0.75
 - Log domain quantization: $2^{\text{round}(\log_2(0.75))} = 2^{-1} = 0.5$
 - Increasing bits don't help!

Proposed Non-uniform Quantization

- **Example:** To quantize value 0.75
 - Log domain quantization: $2^{\text{round}(\log_2(0.75))} = 2^{-1} = 0.5$
 - Increasing bits doesn't help!
- Reduce quantization error with $N = 2, B_1 = 1, B_2 = 2$
 - Codebook $\phi_1 = \{0, 2^{-1}\}, \phi_2 = \{0, 2^{-1}, 2^{-2}, 2^{-3}\}$

$$\phi_1 = \{0, 2^{-1}\}, \phi_2 = \{0, 2^{-1}, 2^{-2}, 2^{-3}\},$$


The diagram shows the expression "1, 10" at the bottom. Two arrows originate from this expression: one points diagonally up and to the left towards the codebook ϕ_1 , and the other points diagonally up and to the right towards the codebook ϕ_2 .

- Quantized value: $\hat{w}_i = 1, 10 = 2^{-1} + 2^{-2} = 0.75$

Proposed Non-uniform Quantization

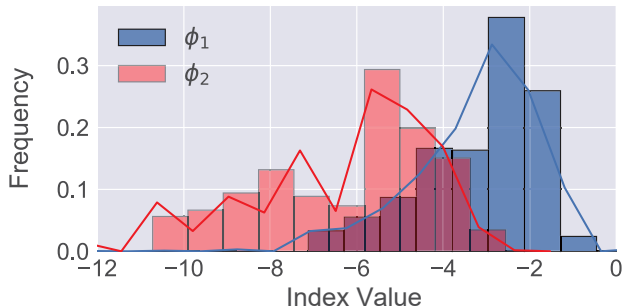


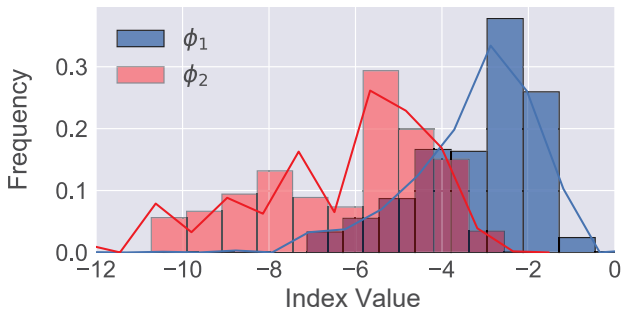
Figure: Index value distribution of FC layer in VGGNet16

- **Codebook index values tend to be centered within a range**
- **More bits are required without optimization**
 - 3 bits for ϕ_1 , 4 bits for ϕ_2 for this case

Proposed Non-uniform Quantization

- Offset β_n to cover wider range

$$\phi_n = \left[0, 2^{-1-\beta_n}, 2^{-2-\beta_n}, \dots, 2^{-(2^{B_n}-1)-\beta_n} \right],$$

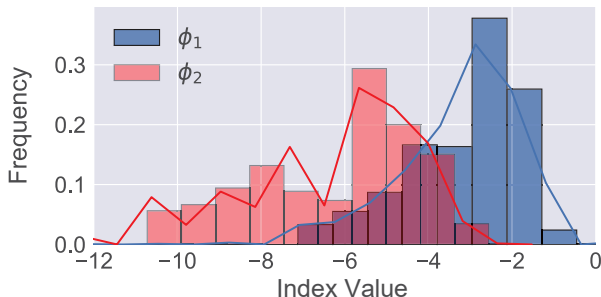


$$\beta_1 = 0 \quad \beta_2 = 2$$

Proposed Non-uniform Quantization

- Offset β_n to cover wider range

$$\phi_n = \left[0, 2^{-1-\beta_n}, 2^{-2-\beta_n}, \dots, 2^{-(2^{B_n}-1)-\beta_n} \right],$$



- Reduce to 3 bits for ϕ_1 , 3 bits for ϕ_2

Proposed Non-uniform Quantization

- MSE criterion to determine optimal offset β_n :

$$\beta_n = \arg \min_{\beta_n} \frac{1}{I} \sum_{i=0}^{I-1} \|\hat{w}_i - w_i\|^2,$$

- Weights in the same layer share the same offsets
- Only require N offset values for a layer
- Increase quantization resolution

Efficient MAC Operation

- MAC based on shift and add

$$y = \hat{w}_i * x_i + b = \sum_{n=1}^N \phi_n [\text{idx}_{i,n}] * x_i + b.$$

- Codebook elements are all power of 2 or zero
- Shift and add instead of bulky multiplier
- One multiplication = N shift and $N - 1$ addition

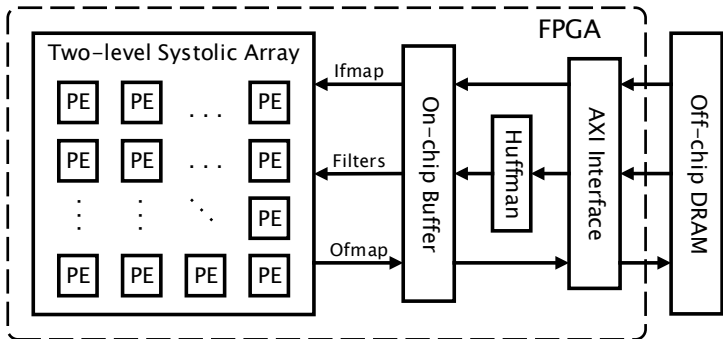
Efficient MAC Operation

- Normalized energy and area cost comparison for single MAC unit for $N = 2, B_1 = B_2 = 3 \rightarrow (3, 3)$

	Power	Area
Shift-add MAC	$1\times$	$1\times$
Fixed-point MAC	$7.3\times$	$14.5\times$

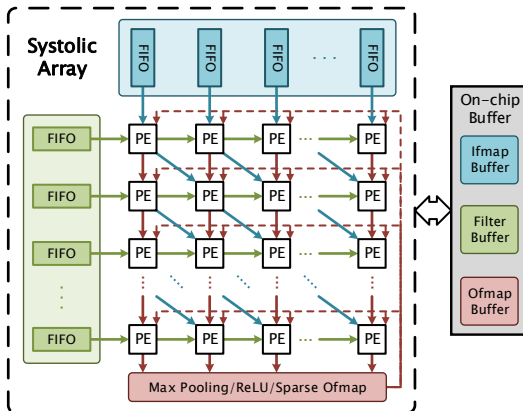
Hardware Architecture

- Huffman Coding → Lossless compression
- Two-level Systolic Array → Maximize data reuse



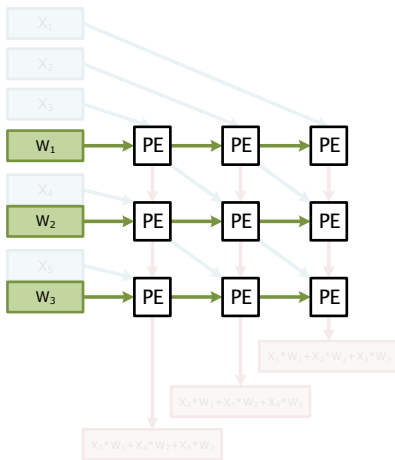
Two-level Systolic Array

- 14×14 PE array
- Row Stationary (RS) dataflow \rightarrow Minimize data movement



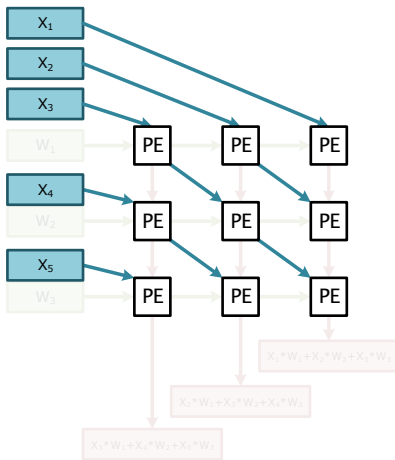
Dataflow of Systolic Array

1. Weights Broadcast



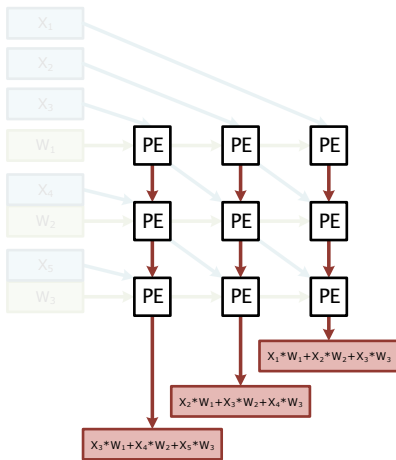
Dataflow of Systolic Array

2. Data Input (16-bit fixed)

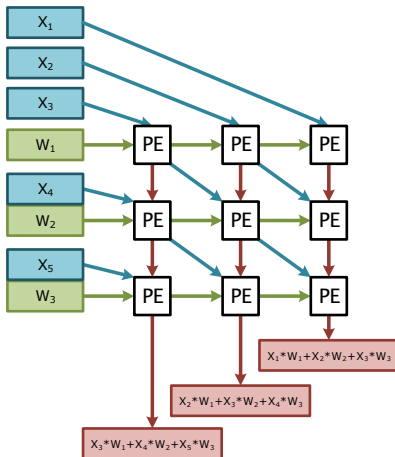


Dataflow of Systolic Array

3. Data Output (Activation: 16-bit fixed)

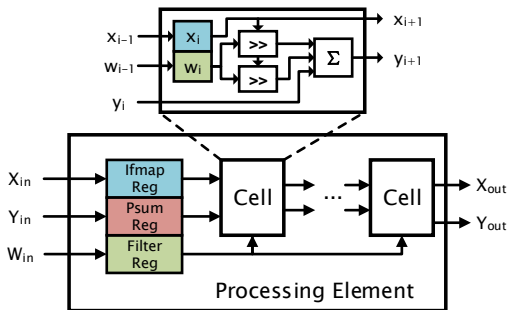


Dataflow of Systolic Array



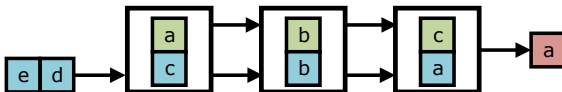
Processing Element

- Each PE contains 5 Cells
- Cell implements shift-add MAC operation
- 1-D systolic convolution \rightarrow Higher throughput



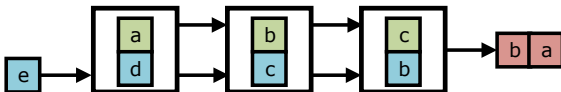
Dataflow of PE

- Weights stay
- Input data move systolically



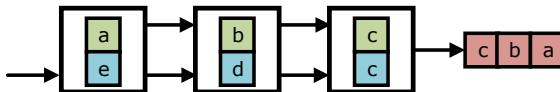
Dataflow of PE

- Weights stay
- Input data move systolically



Dataflow of PE

- Weights stay
- Input data move systolically



Outline

1. Motivation
2. Related Work and Problem Formulation
3. Proposed Quantization and Hardware Co-design
4. Results and Analysis
5. Conclusion

Test on AlexNet

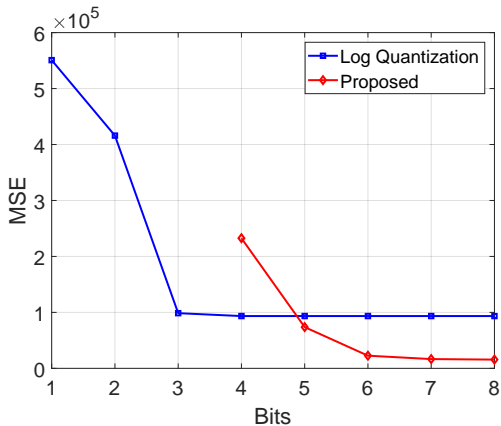
- Codebook size $N = 2$ without **Retraining**

Model	Codebook	Top-1/top-5 Accuracy	Degradation
AlexNet	Baseline	56.55%/79.09%	— / —
	(3, 2)	41.76%/66.22%	−14.79% / − 12.87%
	(4, 2)	48.36%/72.33%	−8.19% / − 6.76%
	(3, 3)	54.98%/77.89%	−1.57% / − 1.20%
	(4, 4)	55.45%/78.64%	−1.10% / − 0.45%

[†] Top-1/top-5 error are tested with single center crop.

Test on AlexNet

- Quantization MSE comparison



Validation on ImageNet

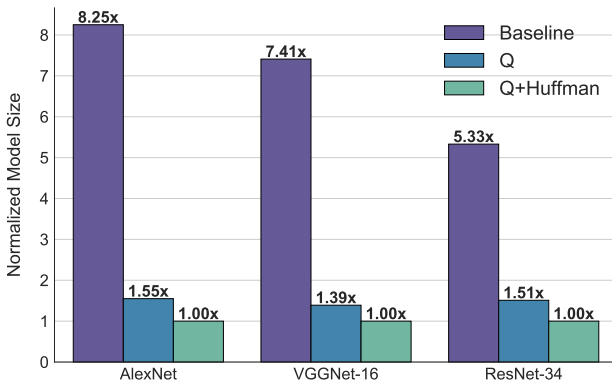
- Quantize pretrained AlexNet, VGGNet16, ResNet34 model from Pytorch
- Codebook size $N = 2$ with $B_1 = B_2 = 3$

Model	Method	Bit-width	Degradation	Retraining
AlexNet	Baseline	32	—/—	No
	Proposed	(3, 3)	−1.57%/ − 1.20%	No
	LogNet	5	—/ − 3.70%	No
VGGNet-16	Baseline	32	—/—	No
	Proposed	(3, 3)	−2.23%/ − 1.95%	No
	Fixed-point	16	−3.58%/ − 2.49%	No
ResNet-18	Baseline	32	—/—	No
	Proposed	(3, 3)	−1.97%/ − 1.17%	No
	ShiftCNN	(4, 4)	−3.21%/ − 2.05%	No
	TWNs	2	−2.56%/ − 1.80%	Yes

[†] Top-1/top-5 error are tested with single center crop.

[*] Degradation is taken from original papers.

Model Compression



Implementation Results

Design	Qiu2016	Zhang2016	This work
Platform	Zynq XC7Z045	Virtex-7 VX690t	Virtex-7 VX690t
Clock(MHz)	150	150	150
Quantization	16-bit fixed	16-bit fixed	(3, 3)
LUT	186, 251	$\approx 300,000$	107995
FF	127, 653	$\approx 300,000$	117795
DSP	2240	2833	0
BRAM	1024	1248	1279
Throughput (GOP/s)	187.8	636.0	238.2

¹[Qiu, Wang, Yao, *et al.*, ISFPGA 2016]

²[Zhang, Fang, Zhou, *et al.*, ICCAD 2016]

Outline

1. Motivation
2. Related Work and Problem Formulation
3. Proposed Quantization and Hardware Co-design
4. Results and Analysis
5. Conclusion

Conclusion

- **A framework to implement low-precision CNNs**
 - Non-uniform quantization with multiple codebooks and offset
 - Retraining-free quantization approaches
 - Multiplier-free shift-add convolution
- **Efficient hardware architecture**
 - Two-level systolic to maximize data reuse
 - Huffman compression to reduce memory bandwidth
 - 1-D systolic PEs to obtain high throughput

Reference

1. M. Rastegari, V. Ordonez, J. Redmon, *et al.*, "Xnor-net: Imagenet classification using binary convolutional neural networks," in *ECCV*, 2016
2. F. Li, B. Zhang, and B. Liu, "Ternary weight networks," *arXiv preprint arXiv:1605.04711*, 2016
3. E. H. Lee, D. Miyashita, E. Chai, *et al.*, "Lognet: Energy-efficient neural networks using logarithmic computation," in *ICASSP*, 2017
4. A. Zhou, A. Yao, Y. Guo, *et al.*, "Incremental network quantization: Towards lossless cnns with low-precision weights," in *ICLR*, 2017
5. J. Qiu, J. Wang, S. Yao, *et al.*, "Going deeper with embedded fpga platform for convolutional neural network," in *ISFPGA*, 2016
6. C. Zhang, Z. Fang, P. Zhou, *et al.*, "Caffeine: Towards uniformed representation and acceleration for deep convolutional neural networks," in *ICCAD*, 2016
7. D. A. Gudovskiy and L. Rigazio, "Shiftnet: Generalized low-precision architecture for inference of convolutional neural networks," *arXiv preprint arXiv:1706.02393*, 2017

Thanks for Your Attention!

Q & A