

MOTIVATION

Deep neural nets (DNNs) achieve state of the art performance tasks in machine learning, but challenges exist in deploying DNNs on embedded platforms.

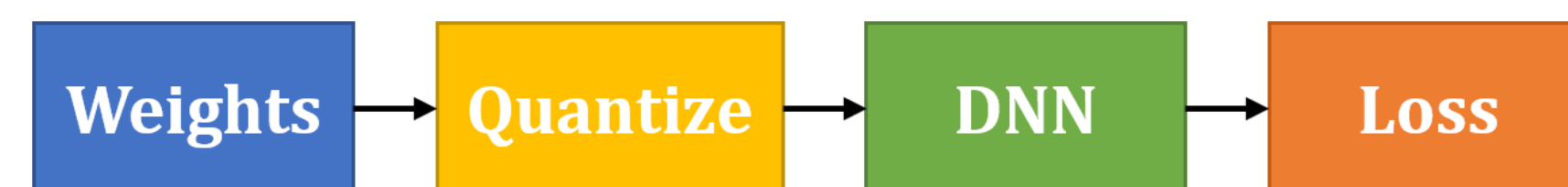
- ❑ Large model size of millions of parameters requires large memory
- ❑ Many operations can lead to large power draw and latency
- ❑ Full precision weights not feasible in mobile systems.

PROJECT AIM

- ❑ Train DNNs to retain high accuracy under quantization
- ❑ Use binary and ternary quantization to limit memory, power consumption
- ❑ Employ regularization to ensure weights are robust against quantization

QUANTIZATION-AWARE TRAINING^{[1][2]}

Forward Pass:



Backpropagation:



- ❑ Calculate forward pass with quantized weights
- ❑ Replace quantizer with identity mapping during backpropagation

METHOD

Training Procedure

1. Train network for 30 epochs using quantization-aware update:

- $w_{t+1} = w_t - \alpha \frac{\partial L(w)}{\partial Q(w)}$
- $L(w)$ is the loss, $Q(w)$ are quantized weights

2. Add regularization term to the loss

- $L_R(w) = L(w) + \frac{\partial^2 L(w)}{\partial w^2} (w - Q(w))^2$
- Regularized update:

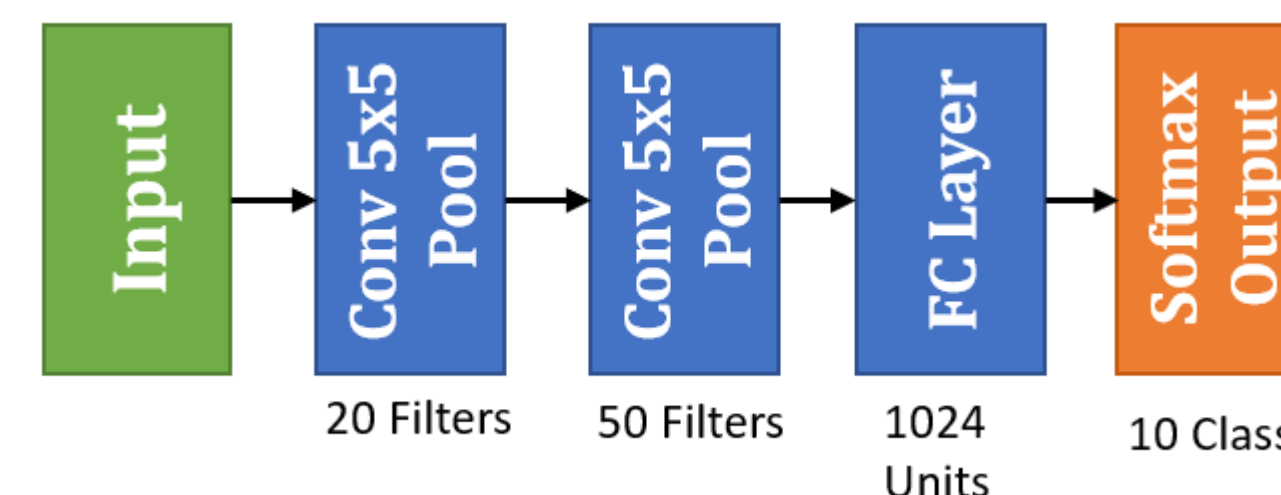
$$w_{t+1} = w_t - \alpha \frac{\partial L(w)}{\partial Q(w)} - \frac{\partial^2 L(w)}{\partial w^2} (w - Q(w))$$

- Regularizer forces broad minimum solution (low curvature)

3. Train for 20 epochs using $L_R(w)$

DATASET AND DNN ARCHITECTURE

- ❑ FashionMNIST 10-class dataset of greyscale images of clothing (T-shirt, sneaker, dress)
- ❑ Lenet-5 Network Architecture^[3] (2.5 million parameters):



PRELIMINARY RESULTS

- ❑ With our approach, DNN retains high performance under quantization of weights and biases

Quantizer	Quantizer Levels	Test Accuracy
None (baseline)	n/a	92.3% (baseline)
Ternary	[-0.1, 0, 0.1]	92.17 ±.12%
Binary	[-0.1, 0.1]	91.8 ±.17%

Accuracy of Lenet-5 under ternary and binary quantization. Quantizer levels chosen using K-means on weights of normally trained network.

ONGOING & PLANNED WORK

- ❑ Extend regularization scheme to other network architectures such as VGG, ResNet, etc, and other datasets like CIFAR10, ImageNet.
- ❑ Quantize activations in addition to weights and bias values.
- ❑ Deploy and evaluate model on embedded platform.

REFERENCES

- [1] M. Courbariaux, Y. Bengio, J.P. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," Advances in Neural Information Processing Systems, 2015.
- [2] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnornet: Imagenet classification using binary convolutional neural networks," European Conference on Computer Vision, 2016.
- [3] Y. LeCun, "LeNet-5, convolutional neural networks," November 2013.

ACKNOWLEDGEMENTS

This work is supported in part by the NSF I/UCRC NCSS SenSIP site and Raytheon Missile Systems.