### Running 2-bit quantized CNN models on ARM CPUs



Davis Sawyer Co-founder & CPO

May 4, 2023

2

 $\bigcap$ 

Deeplite



#### **Optimize AI for the Edge**



Founded 2019



# Edge Computing Challenges

#### 

#### **High Computational Complexity**

Millions of expensive floating-point operations for each input classification are needed.

#### 0

#### **Memory Footprint**

Huge amounts of weights and activations with limited onchip memory and bandwidth.



#### **Power consumption**

Deep learning requires significant power and can easily consume battery life





#### In-memory / near memory Computation

- We must bring computation at the very edge (on-chip, in-memory) Otherwise we lose a lot only on memory bottleneck
- Size is the biggest challenge in this environment



## Industrial EdgeAI Challenges

- Big models/datasets 🙁
  - Expensive Training Step
- Very sensitive on accuracy drop depends on the use case
- Expecting real time execution on their edge HW 🙁
- Generic and automated approaches which work on all use cases





End-to-end Computer Vision Pipeline

Flexible platform with custom and opensource options Sophisticated and production-tested model optimization

## Neutrino Model Optimization



**10x improvement** (more FPS, inferences/sec)

#### at lower cost

(less W, less MB, less \$)

**Pre-trained** 

Model

#### Our Low-bit Quantization Contribution



- QAT with 1 and 2-bit precision (weights and activations) for both detection and classification models.
- Supports any precisions (1,2,3, etc.)
- Mixed precision approach to minimize the accuracy drop of quantized models.

- Custom ultra-low precision convolution operators to accelerate speed and memory throughput of quantized layers
- End-to-end framework to deploy and execute mixed precision ultra-low bit quantized models on Armv7 and Armv8 Cortex-A processors.

### Issues with low bit quantization

- High accuracy loss for 2-bit quantization
- Mixed precision quantization
  - Different layers to have different bit precision to maximize accuracy
    - Sensitivity analysis
    - Minimize switching between different precision
    - trade-off between accuracy and speed
  - Automatic way with 2 bit, 8-bit and FP32 layers

## Deeplite Neutrino Quantizer

Sample Applications	Architecture	Deeplite Neutrino Quantization			precision	Accuracy	Dataset
		Original Size	Optimized Size	Improvement	precision	Change	Dutuset
Image classification	ResNet18	42MB	2.9MB	x15	2w/2a	<b>0.00%</b> (Top1)	CIFAR100
	VGG19	76MB	5MB	x15.3	2w/2a	<b>&lt;1.50%</b> (Top1)	CIFAR100
	ResNet18	42MB	2.9MB	x15	2w/1a	<b>&lt;1.00%</b> (Top1)	VWW
	ResNet50	97MB	17MB	x5.6	2w/2a, 8w/8a¹	<b>~1.50%</b> (Top1)	ImageNet
Object Detection	VGG16_SSD	90MB	5.6MB	x16	2w/2a	<b>&lt;0.01</b> (mAP)	widerface
	VGG16_SSD	100MB	6.2MB	x16	2w/2a	<b>&lt;0.02</b> (mAP)	VOC 2012
	Yolo5_6n	7MB	0.95MB	x7	2w/2a, FP32 <sup>1</sup>	<b>&lt;0.02</b> (mAP)	Custom (Person Detection)

<sup>1</sup>Mixed-Precision

# DeepliteRT

- Custom ultra-low precision convolution operators to accelerate speed and memory throughput of quantized layers
- End-to-end framework to deploy and execute mixed precision ultralow bit quantized models on Armv7 and Armv8 Cortex-A processors.



### Low Bit Conv2D Implementation

FP32	2 bit	Packed 32 bit	Multiplication output	Unpacked

## DeepliteRT

- Uses intrinsic from the Neon vectorized instruction set for both Armv7 and Armv8 architectures to target 32-bit and 64-bit Arm CPU devices.
- Efficient tiling and parallelization schemes used to improve upon the performance of the vectorized kernels.
- On the ResNet18 model running on the low-power Arm Cortex-A53 CPU in the Raspberry Pi 3B+, our overall implementation realizes speedups of up to 2.9x on 2-bit and 4.4x on 1-bit over an optimized floating-point baseline

## ResNet18 on VWW



Accuracy/performance benchmark of DeepliteRT on ResNet18 model on VWW dataset (2A/2W- weights and activations quantized to 2 bits, 1A/2W- activations quantized to 1 bit and weights quantized to 2 bits)

### Industrial Use case #1



Object detection - Yolo5 based model

#### **Detection** Performance









YOLOv5s COCO – Raspberry Pi 4B (4x Cortex-A72)



## Combining All Techniques

Deplice | EIW2022: Running 2 bit quantized CNN models on ARM CPUs

## Conclusion

- 2-bit model running on commodity hardware is presented
- Mixed precision approach is used to minimize the accuracy loss
- Novel method is proposed to run the 2-bit model on Arm Cortex A devices
- Benchmark on classification and Object detection models are presented
- Future work
  - Extend to other hardware
  - Performance Improvement



# Thank you

davis@deeplite.ai