# Fast Implementation of 4-bit Convolutional Neural Networks for Mobile Devices

Anton Trusov<sup>1, 3</sup>, Elena Limonova<sup>1, 2, 3</sup>, Dmitry Nikolaev<sup>2, 3, 4</sup>, Vladimir V. Arlazarov<sup>2, 3</sup>

<sup>1</sup>Moscow Institute of Physics and Technology

<sup>2</sup>Institute for Systems Analysis, FRC CSC RAS

<sup>3</sup>Smart Engines Service LLC

<sup>4</sup>Institute for Information Transmission Problems RAS

Se smart engines

Dmitry Slugin<sup>2, 3</sup>.



# Introduction

Low-bit quantized neural networks (QNNs) allow us to:

- accelerate inference;
- decrease model size;
- perform real-time computation on low-powered devices;
- follow paradigm of edge intelligence.

Low-bit QNNs do not suit end devices of general architecture well:

- one can not directly use BLAS libraries for matrix multiplication;
- CPUs allow only 8-bit (or multiple) access and computations;
- no known efficient CPU implementations for lower than 8-bit quantization.

- We provide a novel algorithm for fast inference of 4-bit quantized neural network on CPU, based on a fast multiplication (used in convolution and fully-connected layers).
- We experimentally prove its efficiency for ARM architecture.

4-bit QNN

Linear quantization method:

$$\begin{split} \hat{w}_i &= \left\lfloor \frac{w_i}{s} \right\rfloor - z \\ s &= \frac{\max(\max_i w_i, 0) - \min(\min_i w_i, 0)}{2^p - 1} \\ z &= \min(\min_i w_i, 0), \end{split}$$

where  $\hat{w}_i$  denotes quantized values,  $w_i$  are floating-point values, s is scale factor, z is a zero-point (offset), p is a number of bits used in quantized values

### **Quantized multiplication**

Let's consider the quantized approximation of matrix multiplication R = WX:

$$\begin{aligned} r_{ij} &= \sum_{k=1}^{D} w_{ik} x_{kj} \approx \sum_{k=1}^{D} s_w (\hat{w}_{ik} - z_w) s_x (\hat{x}_{kj} - z_x) \\ &= s_w s_x \Big( \sum_{k=1}^{D} \hat{w}_{ik} \hat{x}_{kj} - z_w \sum_{k=1}^{D} \hat{x}_{kj} - z_x \sum_{k=1}^{D} \hat{w}_{ik} + D z_x z_w \Big) \end{aligned}$$

where  $r_{ij}$  denotes values of R matrix,  $w_{ik}$  and  $x_{kj}$  are values of W and X matrices,  $\hat{w}_{ik}$  and  $\hat{w}_{ik}$  are their quantized approximations,  $s_w$  and  $s_x$  are scale factors,  $z_w$  and  $z_x$  are zero-points and D is a depth of multiplication.

# Quantized multiplication micro-kernel

1	3	5	7	
2	4	6	8	
9	11	13	15	
10	12	14	16	



9	17	25	
10	18	26	
11	19	27	
12	20	28	
13	21	29	
14	22	30	
15	23	31	
16	24	32	
	10 11 12 13 14 15 16	J     II       10     18       11     19       12     20       13     21       14     22       15     23       16     24	17     23       10     18     26       11     19     27       12     20     28       13     21     29       14     22     30       15     23     31       16     24     32

(b) LHS

Figure 1: The order or values of right (1a) and left (1b) matrices in temporal buffers



(a) The smaller kernel

res0

ş

rhs

4

5

(b) The bigger kernel

rhs

Figure 2: The kernel layout.

## Quantized convolutional layer

- Perform Im2col transformation to turn convolution into matrix multiplication.
- Compute matrix multiplication (4-bit to 16-bit):

$$\hat{r}_{ij} = \sum_{k=1}^{D} \hat{w}_{ik} \hat{x}_{kj} - z_w \sum_{k=1}^{D} \hat{x}_{kj} - z_x \sum_{k=1}^{D} \hat{w}_{ik} + Dz_x z_w,$$

• Save floating-point scale factor:

$$s_r = s_w s_x$$

# Quantized neural network

We can combine quantized (Q) and not quantized (F) layers in neural network:

- $\textbf{F} \rightarrow \textbf{F}.$  Input: floating-point activation. Output: floating-point activation.
- F → Q. Input: floating-point activation. It is quantized to 4-bit integer with scale factor s<sub>x</sub>. Output: 16-bit integer activation, scale factor s = s<sub>x</sub>s<sub>w</sub>.
- Q → Q. Input: 16-bit integer activation, scale factor s<sup>\*</sup>. Activation is quantized to 4-bit integer with scale factor s<sub>x</sub>.
  Output: 16-bit integer activation, scale factor s = s<sub>x</sub>s<sub>w</sub>s<sup>\*</sup>.
- Q → F. Input: 16-bit integer activation, scale factor s\*. Activation is converted back to floating-point by multiplication by s\*. Output: floating-point activation.

# Experiments

# Matrix multiplication test

Height	Width	Depth	Floating point 32-bit time	Int 32-bit time	Unsigned int 8-bit time	Unsigned int 4-bit time	Unit
		10	8.3	8.3	4.6	3.9	$\mu s$
	100	40	21	21	12	8.4	$\mu s$
		100	53	52	28	19	$\mu s$
		10	32	33	17	14	$\mu s$
8	400	40	90	90	50	33	$\mu s$
		100	0.21	0.21	0.15	0.11	ms
		10	0.13	0.14	0.070	0.058	ms
	1600	40	0.35	0.36	0.25	0.19	ms
		100	2.2	2.4	0.78	0.55	ms
		10	15	15	9.1	6.2	$\mu s$
	100	40	44	43	24	15	$\mu s$
		100	0.11	0.10	0.054	0.032	ms
	400	10	62	61	35	24	$\mu s$
24		40	0.18	0.17	0.096	0.058	ms
		100	0.44	0.42	0.24	0.15	ms
	1600	10	0.24	0.24	0.14	0.097	ms
		40	0.72	0.71	0.43	0.28	ms
		100	3.2	3.1	1.2	0.76	ms

Measured on ODROID-XU4 single-board computer with Samsung Exynos5422 ARM processor

# QNN test

- 36 MRZ character recognition from MIDV-500 dataset
- ODROID-XU4 single-board computer with Samsung Exynos5422 ARM processor

### IDD<<T010089212<<<<<<<<<<<

#	Layer	Activation	Parameters	
	type	function		
1	Conv	ReLU	8 filters 5 $ imes$ 5, stride 1 $ imes$ 1	
2	Conv	ReLU	8 filters 3 $ imes$ 3, stride 1 $ imes$ 1	
3	Conv	ReLU	8 filters 3 $ imes$ 3, stride 2 $ imes$ 2	
4	Conv	ReLU	16 filters 3 $ imes$ 3, stride 1 $ imes$ 1	
5	Conv	ReLU	16 filters 3 $ imes$ 3, stride 2 $ imes$ 2	
6	Conv	ReLU	24 filters 3 $ imes$ 3, stride 1 $ imes$ 1	
7	FC	SoftMax	36 neurons	

Model	Accuracy	Accuracy	Convolution	Total
	synthetic, %	MIDV, %	time, ms	time, ms
CNN	99.8	95.6	0.99	1.22
QNN-8	99.7	95.4	0.55	0.74
QNN-4	99.2	95.0	0.45	0.63
QNN-32	-	-	1.16	1.47

# Results

- Our 4-bit quantized matrix multiplication works about 3 times faster than floating-point multiplication from Eigen library and 1.5 times faster than 8-bit quantized multiplication similar to gemmlowp library
- Our 4-bit QNN works about 2 times faster than traditional CNN and 1.2 times faster than 8-bit QNN of the same architecture.
- The real-world problem of OCR recognition on the MIDV-500 dataset demonstrates 95.0% accuracy, while the floating-point network gives 95.6% accuracy.

# Thank you for your attention