

1 We thank the reviewers for the positive comments, and would like to provide answers to the questions raised.

2 We would like to first address Reviewer #1’s questions regarding the inference time and a case study for IoT devices.

3 The reviewer states that not all compression techniques provide good computation gains due to a lot of constraints.

4 This is why in our evaluation, the target is a custom hardware implementation, and thus the issues above of efficiently

5 targeting a particular CPU or GPU don’t apply on custom hardware. Instead, the #Gates estimations measure hardware

6 efficiency and they are derived from a hardware accelerator design mapped to a Stratix 10 FPGA device. We additionally

7 ensure that all designs with different quantization methods have the same processing performance (line 246 in the paper).

8 A hardware design using focused quantization (FQ) requires fewer gates than its counterparts (LQ-Net, ABC-Net) to

9 achieve the same performance, hence is more power efficient. Conversely, assuming that they all use the same amount

10 of logic resources, FQ thus has a higher throughput than the compared methods, as we can perform more computations

11 in parallel. While it is possible to provide certain performance statistics (*e.g.* energy consumed per image, frames per

12 second given the same chip area, *etc.*) for a particular ASIC design, yet the intricacy involved (*e.g.* the technology node,

13 circuit design, memory bandwidth, *etc.*) is beyond the scope of this paper. The reviewer suggests us to provide an

14 example of how our technique relates to future deployment of CNN accelerators in IoT devices. We believe future IoT

15 systems bring heterogeneous compute: CNN accelerators as co-processors to assist general purpose CPUs. Note that

16 IoT CNN accelerators today typically operate on common quantizations with higher bit-widths, *e.g.* 8-bit fixed-point

17 weights and activations in Eyeriss v2 [1]. The updated paper will address this and include in Table 4 additional #Gates

18 comparison with 8-bit fixed-point weights and activations.

19 Reviewer #3 suggested that as focused quantization introduces high-precision hyperparameters μ_- , μ_+ , σ_- and σ_+ ,

20 the dot products would require more high-precision multiply-adds. We would like to thank the reviewer for catching

21 this, and will clarify this accordingly. In all of our experiments in the supplementary source code, to simplify inference

22 computation, the hyperparameters μ_- and μ_+ are quantized to the nearest powers-of-two values, and σ_- and σ_+

23 are constrained to be equal in the optimization process. Effectively in our implementation of the dot-product, all

24 high-precision parameters, *i.e.* σ_- , σ_+ and the layer-wise learnable scalar α , can thus be fused into batch normalization

25 during inference. The reviewer also pointed out that Figure 4 is missing σ values, which can be explained by the same

26 reason above. We will clarify this in both Figure 2 and Figure 4.

27 In Table 1 below, we present the top-1 and top-5 accuracies for the pruned, encoded but not quantized models, as kindly

28 suggested by Reviewer #3. These results will later be added to the paper. We are working on 3-bit results as kindly

29 requested by Reviewer #2, and will include them in later changes.

30 Finally, we would like to additionally make the following minor changes to further improve the quality of the paper:

- 31 1. The notation of the form “/N” (*e.g.* “/128”) in Figure 2 and Figure 4 denotes all numbers in the same
- 32 grid share a common denominator N , and thus the true values are the numbers scaled by $\frac{1}{N}$. We will add
- 33 descriptions to the figures to make it clear to the readers.
- 34 2. We will add the y-axis labels (“frequency”) to Figure 2 and Figure 3 as kindly suggested by Reviewer #3.
- 35 3. We will adopt the numerical citation style, replace “DNNs” in the title with “CNNs”, and fix caption spacing
- 36 and the notation as kindly suggested by Reviewer #1.

37 We thank the reviewers again for your detailed comments and help on improving the quality of our paper.

Table 1: Additional results for Table 1, newly added rows are shaded.

Model	Top-1	Δ	Top-5	Δ	Sparsity (%)	Size (MB)	CR (\times)
ResNet-18	68.94	—	88.67	—	0.00	46.76	—
Pruned	69.24	0.30	89.05	0.38	74.86	8.31	5.69
ResNet-50	75.58	—	92.83	—	0.00	93.82	—
Pruned	75.10	-0.48	92.58	-0.25	82.70	11.76	7.98
MobileNet-V1	70.77	—	89.48	—	0.00	16.84	—
Pruned	70.03	-0.74	89.13	-0.35	33.80	6.89	2.44
MobileNet-V2	71.65	—	90.44	—	0.00	13.88	—
Pruned	71.24	-0.41	90.31	-0.13	31.74	5.64	2.46

38 References

39 [1] Y. Chen, J. S. Emer, and V. Sze. Eyeriss v2: A flexible and high-performance accelerator for emerging deep neural

40 networks. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, 2018.