1 We thank all reviewers for their insightful comments, and have addressed them below.
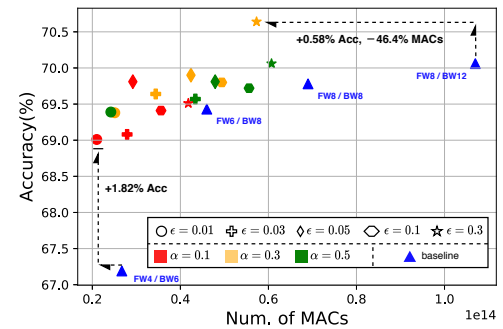
2 **[R1] FracTrain on larger and deeper models:** Thanks for the ad-
3 vice which helps to strengthen our evaluation. Given the limited
4 time, we apply FracTrain to ResNet-110/ResNet-164 on CIFAR-
5 10/CIFAR-100 and find that again FracTrain consistently outperforms

| Method | ResNet-110 | | ResNet-164 | |
| --- | --- | --- | --- | --- |
| | CIFAR-10 | CIFAR-100 | CIFAR-10 | CIFAR-100 |
| FW8/BW8 | 93.38 | 72.11 | 93.72 | 74.55 |
| FracTrain | 93.51 | 72.19 | 93.77 | 74.8 |
| Comp. Saving | 67.3% | 45.17% | 38.69% | 43.6% |

6 the FW8/BW8 baseline with 38.69%-67.3% computational savings under a lightly higher accuracy (+0.05%- +0.25%).

7 **[R1] MP on BitFusion and Integer on integer hardware:** We will clarify the evaluation metrics in the final version.
8 We used Bit-Fusion for both the MP and integer-only systems to keep the hardware parameters (e.g., dataflows) the
9 same for a fair comparison. We have conducted experiments to address your suggestion by comparing FracTrain on
10 Bit-Fusion and FW8/BW8 on an integer-only hardware Eyeriss [Y. Chen, ISCA'16] based on the simulator in Tetris [M.
11 Gao, ASPLOS'17]: for ResNet-38/74 on CIFAR-100 (accuracy in Table 3), FracTrain on Bit-Fusion still outperforms
12 FW8/BW8 on Eyeriss with +65.8%/+69.8% energy savings and +72.6%/+68.2% latency savings, when both adopting
13 the same unit energy and memory size as in Bit-Fusion for a fair comparison.

14 **[R1] Bit choices of DFQ:** It is in fact DFQ's advantage to allow adaptive allocation of higher precision to important
15 layers/inputs and lower precision to unimportant ones, and thus enable a larger range of precision choices over static
16 quantization, given the same computational cost. Furthermore, we follow your suggestion, and limit BW in DFQ no
17 more than 8 bits: compared with FW8/BW8 on (1) ResNet-74@CIFAR-10 (93.04%) and (2) ResNet-74@CIFAR-100
18 (71.01%), DFQ still achieves slightly higher accuracy (+93.11%/71.11%) with +37.3%/+43.7% computational savings.

19 **[R3] Sensitivity to hyper-params in PFQ:** Figure 2 in the ap-
20 pendix shows PFQ's insensitivity to its precision schedule hyper-
21 params under three different precision schedule strategies. Further-
22 more, we perform your suggested ablation study to evaluate PFQ-
23 FW(3,4,6,8)/BW(6,6,8,8) on ResNet-38@CIFAR-100 under various
24 $\epsilon$ (different shapes) and $\alpha$ (different colors): We can see that a good
25 accuracy-efficiency trade-off can be found in a large range of settings
26 compared with static baselines, showing PFQ's insensitivity to hyper-
27 params. It is intuitive that (1) $\epsilon$ and $\alpha$ control the accuracy-efficiency
28 trade-off, and (2) a larger $\epsilon$ and $\alpha$ (i.e., faster precision increase) lead
29 to higher training cost and higher accuracy.



30 **[R3] Modifications on Bit-Fusion:** We did not modify the BitFusion RTL. As the backpropagation can be viewed
31 as two convolution processes (for computing error and gradient, respectively), we estimate energy by executing
32 the three convolution processes of training sequentially in BitFusion. The reuse patterns optimized by BitFusion is
33 output-stationary for both gradients/activations.

34 **[R3, R4] How MACs are calculated:** Inspired by the computation complexity determined by precision in Sec-2.1 of
35 [30], we calculate the effective MACs of the dot product between a and b using (# of $MACs) * Bit_a/32 * Bit_b/32$,
36 following [J. Shen, AAAI'20], which is in proportional to bit operations. We will clarify this in the final version.

37 **[R3, R4] ML accelerators to support FracTrain:** Both dedicated ASIC (e.g., [H. Yoo, ISSCC'19] and [H. Yoo,
38 JICS'20]) or FPGA accelerators (e.g., EDD [Y. Li, DAC'20]) can help exploit FracTrain's best potential by making use
39 of its lower average precision to save both data movement and computation costs during training. After the submission,
40 we have proceeded to implement FracTrain on FPGA to evaluate its real-hardware benefits, following the design in
41 EDD, which adopts a recursive architecture for mix precision networks (i.e., the same computation unit is reused by
42 different precisions) and a dynamic logic to perform dynamic schedule. Evaluated ResNet-38/ResNet-74@CIFAR-100
43 on Xilinx ZC706 (accuracy in Table 3), FracTrain leads to 34.9%/36.6% savings in latency and 30.3%/24.9% savings in
44 energy compared with FW8/BW8. We will clarify this experiment in the final version.

45 **[R3] Support by bit-parallel accelerators:** Since the precision granularity in DFQ is block/layer-wise (e.g., a block
46 with several layers in ResNet will use the same precision), bit-parallel is feasible within each block/layer (see our
47 answer and FPGA implementation right above this one).

48 **[R4] Dataset split:** We use standard train/test datasets (50000 vs. 10000) for CIFAR-10/100 WITHOUT any special
49 split. The training trajectories in Figure 5 visualize the evolution of test accuracy during the whole training process,
50 where the x-axis captures the total computational cost to reach the current epoch instead of the number of epoch. We
51 will clarify this point in the final version.

52 **[R4] Granularity of the RNN controller:** We use per mini-batch for hardware-friendly run-time quantization in both
53 training and inference. We will clarify this in the final version.

54 **[R1, R3, R4] Typos and missing references:** Thanks a lot for pointing out! We will ensure that they are addressed and
55 proofread more carefully before camera-ready.