

1 **We thank all four reviewers for their valuable feedback. All suggestions will be included in the final draft.**

2 **Reviewer1 On TPR hardware:** The reviewer’s understanding of even/odd phase is correct—but dedicated MAC units  
3 are not needed to handle the additional odd phase. FP4 numbers in the odd phase are simply  $0.5\times$  of those in the even  
4 phase. Thus, the odd-phase final MAC result is also  $0.5\times$  and can therefore be efficiently computed by subtracting 1  
5 from the exponent of the final FP16 result using the same MAC engine. **On Sec.4.2:** Thank you for finding the typo. To  
6 clarify, it is the gap between TPR’s  $\text{Var}(dL/dx)$  and FP32’s  $\text{Var}(dL/dx)$  that increases from layer50 to layer1. Our results  
7 on FP4 and "FP4+FP8" indicate that quantization causes the loss of gradient variance in Fig.6a. Also in [35], vanishing  
8 gradients are reflected by the loss of gradient variance—so we’ll look to study this correlation further. **On Depthwise**  
9 **Conv (DW):** The variance analysis and the hybrid approaches are not limited to Conv1x1 and can be applied to compact  
10 networks as well. In fact, while analyzing MobileNetV2 with DW, we found it optimal to use 4b for Conv1x1 and  
11 follow the FP8 work recommendations to use FP16 for DW layers[15]. FP16 DW does not impact latency much  
12 since it is  $<3\%$  of total FLOPs and is hard to accelerate due to limited data reuse. For MobilenetV2, we achieve 1.2%  
13 degradation on CIFAR10 and 3.0% degradation on ImageNet using FP4. These new accuracy results are especially  
14 promising—given that just INT4 inference in MobilenetV2 remains challenging [23,24]. **On SOTA models:** To our  
15 knowledge, this is the first time a SOTA-size model has ever been trained in 4b without major convergence issues. To  
16 show the robustness of our approach over a larger model set, we trained a large NLP task with transformer\_base model  
17 on WMT De-En and obtained a good BLEU score with 2 points off the baseline (25.4 vs 27.5). Overall, we believe this  
18 work lays a solid foundation upon which future FP4 innovations can be added on to address the remaining accuracy  
19 losses. **On hardware area/power:** The area projection is based on insights derived from multiple previous AI Hardware  
20 and Low-precision FPU designs (published in Symp. on VLSI circuit). In our estimates, a MAC unit that performs  
21 4-way INT4 $\times$ FP4 inner products consumes 55% of the area of the FP16 FPU while providing  $4\times$  throughput, yielding  
22 a total compute density improvement of  $7.3\times$ . Compared to FP16 FPUs, the 4b unit has simpler shift-based multipliers  
23 thanks to the power-of-2 FP4 numbers. It also benefits from the absence of addend aligners, narrower adders, and a  
24 simpler normalizer. Additionally, the simpler datapath requires fewer executing stages and saves flip-flop area.

25 **Reviewer2 On 4b overhead:** We’ve estimated that the overheads for 4b training to be  $<5\%$  of the total GEMM ops.  
26 In back-propagation, all collection approaches cost  $\mathcal{O}(1)$  FLOPs/gradient element. Compared with HFP8[15] and  
27 S2FP8, our FP4 conversion is actually simpler due to the absence of mantissa rounding—and can be done with two  
28 back-to-back instructions(1 MUL+1 custom bit OP). The overhead of statistics collection for layer-wise scaling is  
29 also comparable to that in S2FP8. Specifically, we expect  $\sim 10$  FP16 FLOPs/gradient for PACT BWD(2), Radix  
30 Conversion(3), Two-phase Rounding(3), and Layer-wise Scaling(2) overheads. These overheads are much smaller  
31 than  $\mathcal{O}(k_i \times k_j \times \text{channel})/\text{gradient}$  in convolution GEMMs (e.g. In ResNet50, the effective GEMM FLOPs is 642  
32 per gradient element). Therefore, **with the majority of FLOPs spent on GEMM, 4b training retains significant**  
33 **advantage over HFP8 and S2FP8 training due to the throughput and power & area boost in going from 8b to**  
34 **4b GEMM.** With additional optimization from our compiler [published in IEEE Micro], 4b ResNet50 training can  
35 yield at least 60-80% higher throughput vs. HFP8 training along with a 42% area and power saving. **On conversion**  
36 **hardware:** The conversion between radix-2 and radix-4 is remarkably simple for FP4. Due to the absence of mantissa  
37 bits, this can be done through simple bit operations on exponent - which rounds down the FP16 exponent to even or  
38 odd value (and handles overflow and underflow) through common clipping and rounding operations.

39 **Reviewer3 On Conv1x1:** Our mixed precision approach is not limited to only Conv1x1, which is chosen in the context  
40 of ResNet50. For compact models like ShuffleNet and MobileNet, Depthwise convolution layers will be in higher  
41 precision (see line 8). Although the Conv3x3 layers have  $2\times$  FLOPs than the bottom-conv1x1, they yield less accuracy  
42 improvements when cast in FP8 (74.7 vs. 75.0). While the performance is determined by FLOPs instead of gradient  
43 size, bottom-conv1x1 has  $4\times$  output gradients than conv3x3 and top-conv1x1, explaining its efficiency in accuracy  
44 gain when cast in FP8. **Tradeoff of TPR:** Please see line 2. **On optimal radix:** Radix-2 does not have enough range for  
45 gradients. Our experiments with radix-8 on BWD and/or UPDATE GEMM show that training cannot converge due to  
46 poor representation. In addition, for hardware implementation, radix-8 cannot adopt TPR to enhance its representation  
47 due to the non-integral exponent of dividing  $2^3$  into two phases. **On contribution of format:** In the context of radix-4, we  
48 propose efficient rounding schemes that minimize the MSE of gradient quantization. In addition, we propose the hybrid  
49 use of INT4 for FWD and FP4 for BWD. **On channel-level quantization:** We limit the quantization granularity to per  
50 layer for hardware efficiency. If each output-channel has a different scale, reduction over the channel dimension will  
51 be hardware inefficient. **On cancellation:** Fig.4a shows that the same gradient will most likely be quantized towards  
52 opposite directions in each phase. This is important to get a lower expected quantization error  $\mathbb{E}(Q_{err})$ , which is the  
53 sum of mean errors over the entire ensemble from each phase **divided by 2**. The cancellation contributes to a lower  
54 average expected value but individual values could certainly be skewed higher or lower.

55 **Reviewer4 On empirical validation:** We agree with the reviewer and would like to point that we have also evaluated  
56 other models (please refer to lines 12&16). In future, we plan to utilize AutoML techniques to make FP4 broadly  
57 applicable. **On bias:** The numerical bias caused by quantization is random and not expected to result in fairness issues —  
58 especially since it acts at lower levels than the data/target/model choices and is therefore not directly susceptible to  
59 human intervention. We do plan to validate this further as part of future work.