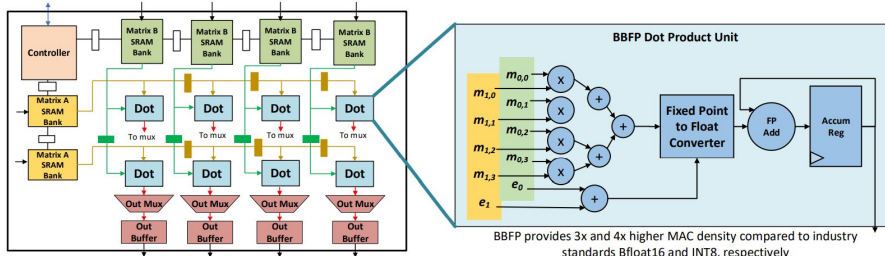


1 We thank the reviewers for their appreciation of our work and helpful feedback. Here, we address open questions.

2 **[Reviewer 7] Reproducibility of hardware:** Following figure shows a systolic tensor core architecture containing
3 multiple BBFP dot product units. Each dot product unit has a significantly lower circuit footprint compared to
4 conventional float due to the shared exponent. The math per bounding-box is mostly performed in fixed-point format
5 and the cost of dynamic scaling is amortized over n (block size). We'll amend the camera-ready with more details.



6 **[Reviewer 7 and 1] Conversion to BBFP:** the axis along which exponent is shared is always the inner dimensions
7 of a matrix multiply: for $A \times B$ we have shared exponents along the rows of A and the columns of B . If the inner
8 dimension is not dividable by the bounding box size, we pad the last bounding box with zero.

9 **[Reviewer 3 and 2] Comparison with FlexPoint, INT8, INT4, and AdaptiveFloat:** As discussed in the paper,
10 leveraging shared exponent to lower hardware cost is not a new idea. However, there is a very large search space of
11 designs across bit-width, exponent selection policy, block granularity, bounding box policy—coupled with hardware
12 implementation—that makes BBFP non-obvious. In this paper, we showed how using a balanced fine-grained approach
13 can provide a robust recipe that works across various models. Due to the low user-friction of BBFP format and its
14 high performance, this datatype has been adopted by different teams and deployed in large-scale production. FlexPoint
15 coarse-grained approach, however, results in significant accuracy drops (in presence of outliers) and incur a high-friction
16 pass to recover accuracy. We evaluated the accuracy versus MAC density trade-off for different benchmarks in Figures 1
17 and 5. This comparison includes industry-standard datatypes (bfloat16, INT8, and INT4). Different datatypes' circuitry
18 is compared in Table 1. We'll append Table 1 to include float8, float4, and posits in the revision. Overall, float8
19 has $1.57\times$ and $3.46\times$ more area overhead compared to BBFP16 and BBFP12, respectively. As for comparison with
20 AdaptiveFloat, in ResNet50-ImageNet with an average of 4 bits (sign plus mantissa), BBFP preserves 96.7% of the
21 baseline accuracy whereas AdaptiveFloat preserves 95.2%. Also, BBFP uses a uniform bit-width for all layers whereas
22 AdaptiveFloat adjusts the bit-width per layer, making it more complicated on hardware. Finally, AdaptiveFloat reports up
23 to $1.14\times$ area improvement compared to INT8 whereas BBFP has a $4\times$ lower overhead.

24 **[Reviewer 3 and 2] Fine-tuning process and calibration:** Quantization involves discretization processes such as
25 rounding and truncation that result in null gradients. The idea of straight-through estimator is to replace the gradient of
26 those operators with identity matrix. We'll elaborate more in the revision. Integer-based inference requires calibration
27 due to their fixed dynamic range—tensors must be scaled to the proper range to be represented. BBFP does not require
28 this type of calibration as it already has a scaling exponent. As for bounding-box selection, all experiments in section 4
29 including Tables 3, 4, and Figure 5 have been performed with a fixed bounding box size of 16 (no calibration involved).

30 **[Reviewer 2] Elaboration on Figure 1:** BBFP12 has 4-bit sign-magnitude mantissas and an 8-bit shared exponent.
31 BBFP12 dot-product circuitry is comparable to INT4. Unlike BBFP, INT format follows a two's complement represen-
32 tation. Two's complement MAC costs more in area/energy compared to a sign-magnitude MAC of the same bit-width.
33 We'll make the notation more explicit in the revision. MAC energy/area cost are measured based on a topographical
34 synthesis of these MAC units on TSMC 16nm FF+. Power simulations are based on an input toggle rate of 50% with
35 50% static probability @ 1GHz. Figure 1 accuracy numbers are for Resnet50-ImageNet.

36 **[Reviewer 1] Non-matmul operations and support for strided and atrous convolution:** BBFP is designed to
37 improve dot product performance. All vector operations such as sigmoid, activation functions, point-wise addition are
38 performed in-situ on HW in float16 (with conversions to BBFP supported automatically on HW). Any operation that
39 can be broken down into dot products (including strided/dilated convs) is supported by BBFP. Scatter/gather is one way
40 of reshaping the input data for strided/dilated convs (SRAM address striding and crossbar are other alternatives).

41 **[Reviewer 1] Zero representation:** Zero is represented by having all mantissa bits being 0 for a given value (shared
42 exponent can be any value). BBFP mantissas do not have an implicit leading bit and all mantissa bits are explicitly
43 represented. BBFP does not have a representation for NaN/Inf, but this does not impact DNN inference accuracy.

44 **[Reviewer 7 and 2] KL Divergence and QNSR:** The KL divergence from BBFP to float32 is computed after computing
45 the pertinent normalized histogram of values in each encoding format. KL divergence is defined as $KL(P \parallel Q) =$
46 $\sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right)$. We will replace QNSR with QSNR in the revised paper to avoid negative values.