

1 We thank all reviewers for their thoughtful comments, and respond to specific concerns below!

2 **REVIEWER 1:** Note that the submission already discusses papers that carry out quantization both in the forward and
3 backward pass ([6,15], in Lines 87-96 of our related work section). We will add the Banner et al. work to this discussion,
4 which presents a newer 8-bit quantization scheme (at the cost of some drop in performance). We note that while our
5 work is related, our contribution is not the quantization: indeed, we just use a simple fixed-point quantization scheme.

6 Instead, our contribution is in our backprop **algorithm**. As we discuss below, we believe this algorithm is both **novel**
7 and **significant** because it limits the accumulation of error from quantization, while still delivering savings in memory.

8 **The algorithm is novel** because backprop has thus far always been implemented the same way—by using the same
9 version of activations for forward and backward passes. Attempts at quantization have adopted this implementation and
10 operated “locally”, by only modifying per-layer operations. Breaking from this, our algorithm takes a look at the entire
11 computation process in backprop, and shows it is beneficial to perform the forward pass exactly during training, and
12 that it is possible to do so while still saving memory. As our analysis shows, this minimizes the effect of approximation
13 error by preventing a cascading effect of errors building up from layer to layer.

14 **The algorithm is significant** because it allows the use of much higher rates of approximation and quantization, and
15 thus greater memory savings. Note that in our experiments, we show that even 4-bit fixed point quantization allows
16 successful training, even though we’re using perhaps the simplest quantization function. This is precisely because we
17 are able to do the forward-pass in full-precision. As our experiments show, doing the same quantization naively in both
18 forward and backward passes simply fails.

19 In summary, the new algorithm prevents per-activation approximation errors from propagating across layers in deep
20 networks. It allows greater levels of memory savings with even simple quantization schemes, and we believe will allow
21 greater flexibility in exploring new kinds of approximation and quantization schemes for training.

22 **REVIEWER 2: - Effect of Activation Functions:** Our algorithm is designed specifically for RELU-like activations
23 whose gradient depends on the sign of the activation, but not the value. This is because otherwise, we would incur
24 additional errors while computing the gradient to input layers (eq 6), which would cause errors to build up during the
25 backward pass. Currently, it can’t be applied to layers with sigmoid / tanh activations.

26 **- Applied to RNNs:** Our method can be applied to networks that have occasional sigmoid-like activations by just
27 leaving those layers un-approximated (e.g., we don’t approximate the last softmax layer in our current experiments).
28 But RNNs and transformer networks have sigmoid/tanh activations in nearly every layer, and so the current version of
29 our method would not work on these.

30 We realize that this is a bit disappointing. But as R3 also points out, ReLU-based networks cover a very large class of
31 architectures that are widely used in many application domains. Our method will thus have real practical impact for
32 many researchers and practitioners who train such kinds of models. Also, we believe our method can be a starting point
33 for future work that targets RNNs, etc. Thus, we think this paper will be of interest to the NeurIPS audience.

34 **- Densenet:** If a network is fully-dense (every layer connects to every preceding layer), then our method would offer no
35 savings. But note that DenseNets typically have sequences of dense blocks (with dense connections within blocks), and
36 so W would be the size of the block, not the size of the entire network. Other networks use skip connections, but only
37 from a sparse subset of layers, and would thus also have $W \ll L$ and allow for significant memory savings.

38 **- Other optimizers:** We chose momentum because this was the optimizer used by the baselines. Our method works
39 equally well with Adam and RMSProp. As additional analysis in the revision will show (see response to R3 below), this
40 is because the errors in gradients due to our approximation are much lower than from the randomness of SGD itself.

41 **REVIEWER 3:- More Analysis:** We’ll add visualizations that give a deeper explanation of why our method works: in
42 addition to just showing training accuracy, we have computed results for the errors (when using our method vs exact
43 training) in the actual gradients of individual layer weights. We will plot these, and compare them to the error due to
44 SGD itself—i.e., the variance in the same gradients when computed on different mini-batches for the same model. This
45 will show that our approximation errors are an order of magnitude lower than SGD variance, and help demonstrate why
46 our approximation enables accurate training.

47 **- Other memory-saving methods:** Note that the main prior method for memory saving during training is checkpointing.
48 This is equivalent to exact training in terms of accuracy: the disadvantage being that it’s slower (the memory-speed
49 trade-off can vary based on how frequently layers are recomputed). The other most common approach is to simply
50 quantize / approximate activations as and when they’re computed. We compare to this strategy as our ‘naive quantization’
51 baseline (for equivalent quantizations, we show this simply doesn’t work).