

1 **Reviewer 1. 1. "more background review"** - Thanks for the references. We will add discussion of these
2 papers. **2. "return communication grows with W.."** - In k -sparsification methods, the weight update is the
3 sum of W k -sparsified gradients and can have Wk non-zero entries. In practice, the weight update is dense
4 for large W , as shown in Figure 4. There is no way to losslessly compress this back-communication to
5 under $O(Wk)$ bytes, and we are not aware of any lossy compression that achieves good performance in
6 practice. **3. "In practice how should the size of the Count Sketch table be set?"** - Theory suggests the sketch
7 size should be $O(k \log d)$, with $c = O(k)$ columns and $r = O(\log_c d)$ rows. The exact values must be found
8 experimentally. **"I expect if the model parameters themselves are sparse then the gradients are more likely
9 to be sparse"** - In general, sparsity of model parameters doesn't imply sparsity of gradients. For example,
10 consider squared loss $(y - w^\top x)^2$ with a sparse w , the gradient is $-2(y - w^\top x)x$, which might not be sparse.
11 We make no assumptions about the sparsity of the gradients, and for the models and loss functions we use
12 (and that are commonly used), the gradients are dense.

13 **Reviewer 2.(A) "does it reduce training time.."** - We have not yet benchmarked overall running time on
14 realistic network topologies, though we are actively working in this direction. Sketching does impose some
15 computational overhead, and with our current implementation we expect a speedup in overall run time
16 when communication cost is about as high (or higher) than the computation cost. This is realistic for a
17 large network (200M parameters) even over a fast (10Gbps) network. In the regime of many workers (e.g.
18 federated learning), where our algorithm performs best theoretically, network can be as slow as a few Mbps.
19 **(B) using asynchronous updates** - We have not experimented with asynchronous updates, but we expect that
20 the benefits asynchrony brings to regular SGD would also apply to Sketched SGD. We believe our theoretical
21 results could be extended to the asynchronous case, but we leave this to future work.

22 **Reviewer 3. 1. "Bounded second moment and variance assumption"** - A bounded second moment implies
23 bounded variance (via variational definition of variance), so $\sigma \leq B$. However note that when you average
24 W stochastic gradients, the variance reduces by factor W , but the second moment bound is still the same. We
25 use different parameters for variance & second moment bounds so that it is easy to argue how the quantities
26 in the upper bound change when we average stochastic gradients. See remark 3 in the paper.

27 **2. " $k \log(dT/\delta) > d$ for sufficiently large T ..."** - $T > \Omega(e^d)$ is not realistic for even our smaller model ($d \approx 10^7$).

28 **3. "Can't one send the compressed gradients back..?"** The $O(W)$ per-worker communication cost is not a
29 technicality: if each worker sends k non-zero gradient elements, the weight update (which is the sum of
30 these sparse gradients) will have up to kW non-zero elements. In practice, the weight update quickly does
31 become dense, as shown in Figure 4. In principle, one could lossily compress the back communication, but
32 we are not aware of any methods that do so while maintaining high performance.

33 **4. "... you only send back the k elements that were modified right?"** - Correct, \tilde{g}_t only has k non-zero elements, so
34 we only need to transmit $O(k)$ bytes. We will add a remark to make this clear. **5. "TODO"** - Thanks!

35 **6. "extends to methods like Adam?"** - This is a good point. The Adam update uses the ℓ_2 norm of the sum (or
36 linear combination) of gradients in the update, and count sketch can be used to approximate the ℓ_2 norm.
37 Therefore we can, in principle, mimic the Adam update. We do not know if this approximation is good
38 enough to preserve Adam's theoretical properties, or if it would work well in practice.

39 **7. Regarding momentum?** - Sorry for the poor terminology - all experiments use momentum.

40 **8. "...why momentum factor masking 'interferes with momentum'"** - in line 11 of Algorithm 3, we zero out the
41 coordinates of each u_i that were updated (this is momentum factor masking). In the limit of $k \rightarrow d$, each
42 coordinate is updated every iteration, so there is effectively no momentum.

43 **9. "communication reduction at the expense of model accuracy"** - To be clear, we get no drop in performance with
44 40x compression for the translation task. For CIFAR-10, we can make up the accuracy drop by training
45 longer. For example, training with 17x compression for the usual number of iterations gives 92.5% validation
46 accuracy. Training with 50% more iterations (reducing to 11x overall compression) restores accuracy to 94%.
47 We expect similar results to hold for the translation task at compression rates higher than 40x.

48 **10. "compare empirically against others"** - We outperform many popular methods. For example, the "Local
49 Top- k " method in Figure 4 is very similar to deep gradient compression (Lin+2017), except we do not clip
50 gradients, and we warm up the learning rate instead of the sparsity. Local Top- k does about as well as 9x-
51 compression Sketched SGD, but it achieves only $\sim 2x$ compression for large enough W . We also implemented
52 and compared to signSGD (arXiv:1810.05291), which gets 92.5% validation and 32x compression. When we
53 quantize all communication in Sketched SGD to 16 bits, we achieve 92.5% accuracy with 34x compression
54 ($r = 3, c = 100,000, P = 8$). For the translation task, we already achieve 40x compression with no loss in
55 accuracy, which is better than signSGD's 32x, and which we expect to double to 80x without loss of accuracy
56 by adding 16-bit quantization. TernGrad (arxiv:1705.07878) is another popular method, but achieves hardly
57 any compression of the back-communication from parameter server to workers for the large models we
58 consider with $d \approx 2^{26}$ - see TernGrad pg. 3.