

1 We thank the reviewers for their time and productive comments. However, there is one misunderstanding that might
2 have affected the reviewers' scores: reviewers 1 and 3 have both misunderstood the baseline. When compared like-
3 for-like, our results outperform the baseline by a large margin (details below). While this misunderstanding is surely
4 a shortcoming of our presentation, the reviewers criticised that we did not improve on baseline performance, when
5 we in fact did. Leaving the misunderstanding aside, reviewers found the work to be "really well written and every
6 decision is well motivated" (Reviewer #3), found that the proposed implementation "greatly improves the efficiency for
7 generating random subspace at each training step" (Reviewer #2) and wrote that the contributions are "novel relative to
8 prior work [...] and can inspire future work in the area" (Reviewer #1) without mentioning any other major concern.
9 Given a correct understanding of the baseline, it seems likely that their overall scoring would have been more positive.

10 To expand further on the misunderstanding, Reviewers #1 and #3 pointed out that our results did not seem to be
11 consistent with the results published by Li et al. [1]. The misunderstanding most likely stems from the fact that Li
12 et al. [1] reported the achieved accuracy as a percent value relative to the SGD baseline while we reported the absolute
13 percent accuracy, without normalising against the SGD baseline. In particular, for a 20x reduced CIFAR-10 LetNet, [1]
14 reported 90% of the original 58% SGD accuracy which amounts to an $0.9 \cdot 57\% = 51.3\%$ accuracy in absolute terms
15 (see Figure S14b in [1] which presents the absolute accuracies). This is consistent with our reported 58.35% accuracy
16 for a 10x reduced Resnet-8-CIFAR-10 in Table 1, where the 7% improvement can be attributed to increased efficiency
17 of the ResNet architecture, as Reviewer #3 expected. Similarly, our MNIST baseline of 80% reported in Figure 1 for
18 $d=250$ is consistent with the absolute accuracy reported in Figure S6 in [1] (note that the network dimensionality is
19 $D=100K$, but the subspace dimensionality is $d=250$ only). We realize now that our reporting should have made this
20 subtle issue of the percent notation in [1] clearer. We will add the relative accuracy levels to the manuscript to ease the
21 direct comparison with prior art.

22 *Minor:* All suggested improvements are gratefully received and we will incorporate the feedback into our revision.

23 * Reviewer #2 asked whether the substantial IPU hardware speedup over CPUs was due to the accelerated PRNG or
24 could be merely explained by the forward-backward pass acceleration.

25 > While the IPU accelerates the forward-backward pass of the network, we found that the main bottleneck on CPU
26 hardware is indeed the PRNG (particularly for large subspace dimensions $d > 1000$). To rule out the possibility that the
27 measured speedup can be attributed to the forward-backward acceleration only, we benchmarked the throughput of our
28 implementation on a GPU V100 accelerator that, unlike the IPU, does not have an on-chip PRNG. We found that the
29 GPU provided no throughput improvement relative to the CPU baseline. We will include these additional results in the
30 respective Section 4.2.

31 * Reviewer #2 noted that the throughput "31 images per second" does not match with the "100 epochs / 67 minutes"
32 statement in Section 4.2.

33 > Thank you, this is a mistake. We accidentally mixed the images per second throughput for $d=10k$ with the wall-clock
34 time figure for $d=1k$. The correct throughput for $d=1k$ is 1366 and 112 images per second on IPU and CPU respectively.

35 * Reviewer #1 asked if there was "any comparison to FPD approach in terms of parallelization?"

36 > Both approaches can be parallelized in the same way since FPD can be seen as a special case of our algorithm where
37 $\varphi_t \equiv \varphi_0$. Our more efficient distributed implementation can thus be seen as a technical contribution that can also
38 benefit the investigations of intrinsic dimensionality in [1]. We will update the discussion to point this out.

39 * Reviewer #3 asked whether the low performance of the NES baseline stems from a small number of random samples.

40 > Indeed, the NES baseline in Figure 2 used the same very low-dimensional number of $d=250$ samples, while more
41 samples would certainly increase the approximation quality. The low-dimensional comparison at $d=250$ demonstrates
42 the superiority of gradient based RBD optimization over NES black-box sampling in this setting. We will adjust the
43 caption of Figure 2 to underline this point.

44 * Reviewer #3 asked about a "potential bottleneck in the distributed version of the training, as all the random numbers
45 need to be generated on the main worker".

46 > This is an good observation that motivates a trade-off between increased compute through PRNG versus reduced
47 communication between workers. Notably, however, our implementation does not require a central main worker but the
48 PRNG generation can be shared between workers in a decentralised way to load balance potential PRNG bottlenecks
49 (see Algorithm 1, right).

50 [1] Chunyuan Li et al. "Measuring the Intrinsic Dimension of Objective Landscapes". In: Sixth International Confer-
51 ence on Learning Representations. 2018. URL: <https://openreview.net/forum?id=ryup8-WCW>.