

1 We thank the reviewers for their thoughtful and positive feedback. We are encouraged they found our work well-argued
 2 (R1) and with substantive contributions (R2). We are pleased (R3) finds our work consistently motivated by theoretical
 3 insights and (R4) identifies the idea as novel and insightful. We address reviewers comments below and we make sure
 4 to incorporate all feedbacks.

5 @ (R1) - **Baselines not tuned properly:** We used the grid for hyper-params set by their authors. For example, for
 6 A-GEM and ER-Reservoir, we had the best hyper-params reported in the appendix of the original papers in our grid.

7 @ (R1) - **Emphasize on two distinct contributions:** We agree and we'll update the introduction and conclusion section.

8 @ (R2) - **Motivate your new forgetting measure:** We also believe this definition needs more context. Regarding the
 9 forgetting measure in experiments we agree it's not new and we will provide additional context to reduce the confusion.

10 @ (R2) - **Early stopping as an important tool for CL:** Yes. Early stopping would prevent going further from the
 11 previous minima. However, there will be a trade-off since if the parameters do not change enough, the new tasks can not
 12 be learned. This is very related to the discussion in l204 as they both refer to the distance from the previous minima.

13 @ (R3) - **Experimental setting:** Task labels are provided for CIFAR-100, only to have a similar implementation with
 14 other baselines (e.g., A-GEM, ER-Reservoir). Nevertheless, we agree that this deserves explicit discussion and we
 15 will update the paper. Regarding the number of epochs for CIFAR experiment, we did so to be compatible with all the
 16 benchmarks as explained in the paper. However, we will add the results for more training epochs.

17 @ (R3) - **Reproducibility:** We apologize if you had difficulty to reproduce the results. We uploaded the polished and
 18 cleaned version of the code since we wanted to make sure the reviewers can investigate the code. However, there were
 19 some minor typos in the uploaded code and we sincerely apologize for that. We wish we could upload a new version of
 20 the code or share an anonymous link to our experiments to make it easy to verify the results. Unfortunately, according to
 21 the NeurIPS guidelines, we can not share any external links. However, as a proof of concept, you can verify experiment
 22 1's rotation MNIST benchmark using the following command which should give an average accuracy of 92.4% with all
 23 accuracy metrics above 90%.

```
24 python -m stable_sgd.main --dataset rot-mnist --tasks 5 --epochs-per-task 5 --lr 0.15 --hidens 100 --batch-size 16 --gamma 0.25 --dropout 0.25 --seed 1234
```

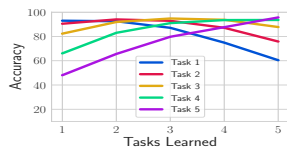
25 In addition, you can use the the following command to replicate experiment 2 on CIFAR-100:

```
26 python -m stable_sgd.main --dataset cifar100 --tasks 20 --epochs-per-task 1 --lr 0.2 --gamma 0.85 --batch-size 10 --dropout 0.1 --seed 2345
```

27 To ensure all results are reproducible we'll add external links to each experiment in the revision.

28 @ (R4) - **Extension with high-dimensional input:** We acknowledge that our results may be affected by the curse
 29 of dimensionality. However, we discussed this issue in lines 106-116 and especially in 129-133. Moreover, recent
 30 work by Fort & Ganguli ("Emergent properties of the local geometry of neural loss landscapes") suggests that: (1)
 31 The Hessian eigenspectrum is composed of a bulk plus C outlier eigenvalues where C is the number of classes, (2)
 32 Gradient aligns with this tiny Hessian subspace which implies that most of the descent directions lie along extremely
 33 low dimensional subspaces of high local positive curvature. Hence, the bound would be still tight enough for our sake
 34 as also demonstrated empirically in Fig2-c,d on high dimensional models.

35 @ (R4) - **Other CL metrics:** As noted by the reviewer, the focus of this work was on forgetting. However, we have
 36 included results for forward transfer in Fig. 1 & 2 on rotation-MNIST. We see the forgetting for stable net is much less
 37 by compromising a negligible amount of forward transfer. This is in line with our discussion on stability-plasticity
 38 dilemma. We will include this result in the paper. Thanks for pointing out.



39 **Figure 1:** Plastic Net

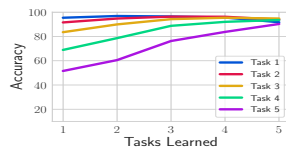


Figure 2: Stable Net

Method	Task 1	Task 2	Task 3	Task 4	Task 5
OGD	75.6	86.6	91.7	94.3	93.4
Stable-SGD	91.5	92.1	95.02	94.2	90.9

Table 1: OGD versus stable-SGD

Method	Average Acc
SGD	53.9 (± 4.2)
Stable-SGD	65.5 (± 3.7)

Table 2: CORE50 result

40 @ (R4) - **Lack comparisons with similar approaches:** We compare our method with OGD. The validation accuracy
 41 for each task at the end of CL experience on rotation MNIST with five tasks is reported in Table 1. We will add
 42 discussions on OGD [Farajtabar et al] and its extensions [Bennani et al] and other papers (EWC, SI, etc) in the revision.

43 @ (R4) - **Benchmarks - low-resolution data with supervised image classification tasks, which are a bit simple:**
 44 While we agree with the reviewer that these datasets are simple, we want to highlight that they are difficult for *sequential*
 45 *multitask* problems. For instance, the best result in our paper on rotation MNIST dataset with 20 tasks yields 78%
 46 accuracy which leaves a lot of room for improvement. Moreover, to demonstrate the feasibility of the proposed work on
 47 other applications we have included a preliminary result for five runs over NI setting of CORE50 dataset with Resnet18
 48 network in Table 2. Full results will be reported in the revision.

49 @ (R4) - **Computational metrics:** Computing the Hessian in the code is just for reproducing eigenvalues in Fig 3 and
 50 empirically showing that minima found by stable-SGD is wider. It is not a part of the algorithm. Hence, stable-SGD is
 51 as computationally cheap as normal SGD, which we believe is an advantage [over many alternatives].