

1 First, we would like to thank the reviewers for their time and effort.

2 **Reviewer 1:** From your brief review, we cannot tell why you believe there to be no take-home message. We describe
3 the warm-start problem, which you note is both important and understudied. We then describe a simple solution that
4 remedies the issue. The original concern *is* resolved by the shrink and perturb trick.

5 Settings described in the introduction are just batch online learning scenarios, which we faithfully simulate in our main
6 experiments (e.g. Figure 7). Please see response 3 under reviewer 4. Final-version plots will be made printer friendly.

7 **Reviewer 2: 1.** Please see Section 3.1 for a detailed analysis of both learning rate and batch size. Figure 3 shows that,
8 while there exist some hyperparameter values for which warm-started models perform as well as randomly-initialized
9 models, they take as long to train as randomly-initialized models. Tuning these parameters can close the generalization
10 gap, but only by sacrificing the computational efficiency we would like to see from warm starting.

11 Separately, it is true that it is commonplace to train models using SGD with a learning rate schedule. To address
12 this concern, we performed an online learning experiment where, starting from 100 CIFAR-10 samples, batches of
13 10000 points are supplied to the learner in sequence. Using a fresh, warm, or shrink-perturb initialization, we train the
14 ResNet using SGD for 350 epochs using the standard learning rate schedule. As expected, a generalization gap between
15 warm-started and randomly-initialized models still exists and shrink-perturb is able to remedy the problem (Fig A).
16 Similar to warm-start plots in the Appendix, the final copy will include a detailed analysis of scheduled optimization.
17 Models were all trained for the same number of epochs, but randomly-initialized models take much longer to converge.

18 **2.** We use Pearson to study the similarity between warm-start and final-solution weights. Other notions of correlation,
19 like cosine similarity, Spearman, and euclidean distance provide similar plots. We will show these in the final version.

20 **3.** It may be difficult to tell from Figure 4, but these models that have not been trained much actually perform significantly
21 worse than fully-trained models. This might be more clear in Table 3, where we apply aggressive regularization to
22 prevent overfitting and still observe a generalization gap. We will make Figure 4 more clear in light of this discussion.

23 **Reviewer 3: 1.** In the Appendix we have extensive experiments using shrink perturb with multilayer perceptrons, rather
24 than only ResNets. We also experiment with batch normalization and weight decay, and will add RNNs (see **1** below).

25 **Reviewer 4: 1.** Indeed, this article focused on experiments with image datasets using ConvNets and MLPs. In
26 response to this note, we performed an online learning experiment using a very different architecture and dataset: a
27 two-layer, bidirectional RNN on the IMDB movie reviews dataset, where the task is to predict whether a movie review
28 is positive or negative. We iteratively supply batches of 500 samples to the model, and compare randomly-initialized,
29 warm-started, and shrink-perturb models in the figure below. We show that there exists a performance gap between
30 randomly-initialized and warm-started models, which the shrink and perturb trick closes (Fig B). The analysis here
31 is truncated due to time constraints, but the trend is clear; we will include a thorough analysis in the final copy.

32 **2.** It is not true that this phenomenon does not happen in vision literature. For example, active learning work in image
33 classification requires models to be trained from scratch at each round [7, 8]. Not all articles say this explicitly, but
34 popular deep active learning github repositories also show that random reinitialization after selection is necessary [9,10].
35 Fine tuning on ImageNet works because the source dataset is large and the target datasets are, comparatively, significantly
36 smaller. Please see Section 4.2 which discusses the pre-training setup in detail. We show that when the target dataset
37 is large, it is often better to initialize models from scratch than it is to pre-train on a source dataset. Figure 9 shows
38 that shrink-perturb is effective in this setting as well: models initialized with shrink-perturb perform as well or better
39 than either randomly initializing or pre-training, regardless of whether pre-training is helpful.

40 **3.** We believe there might be some confusion here. While we do “toy” two-phased experiments for illustrative purposes
41 in Table 2.1 and Figures 3-5, our primary results are for full online learning scenarios. These experiments are like Figure
42 7, where samples are iteratively provided to the learner in batches of 1000. Each time the model receives new data, it is
43 appended to the current training set and is trained to convergence before receiving the next batch. These experiments are
44 consistent with what is usually done in active and online learning work. Also note that large tabular figures, like figure
45 8, also show this fully-online setting. Here, to concisely report results, we are showing only the final accuracies and
46 train times, i.e., after dozens of rounds of online learning. That is, the plots in Figure 7 correspond to a single column of
47 Figure 8. This is likewise true for all of the similar appendix figures. We will make this more clear in the final version.

48 **4.** The warm-start problem is not induced by overfitting. Figure 4 shows that if we do not train to convergence, we still
49 observe a generalization gap. As mentioned above, Table 3 also shows that when overfitting is avoided, by applying
50 regularization, a generalization gap is still present. Also note that the shrink-perturb trick is able to efficiently resolve
51 the generalization gap present in regularized models as well, as shown in Appendix 6.2.5 and 6.2.6.

