
Reproducibility Report for "On Warm-Starting Neural Network Training"

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 Scope of Reproducibility

3 We reproduce the results of the paper "On Warm-Starting Neural Network Training." In many real-world applications,
4 the training data is not readily available and is accumulated over time. As training models from scratch is a time-
5 consuming task, it is preferred to use warm-starting, i.e., using the already existing models as the starting point to obtain
6 faster convergence. This paper investigates the effect of warm-starting on the final model's performance. It identifies a
7 noticeable gap between warm-started and randomly-initialized models, hereafter referenced as **the warm-starting gap**.
8 Furthermore, they provide a solution to mitigate this side-effect. In addition to reproducing the original paper's results,
9 we propose an alternative solution and assess its effectiveness.

10 Methodology

11 We reproduced almost every figure and table in the main text and some of those in the appendix. We used our
12 implementation to produce these results. In case of a mismatch of the results, we also investigated the cause and
13 proposed possible explanations. We mainly used GPUs to train our models using infrastructure offered by public clouds
14 and those that were available to us privately.

15 Results

16 Most of our results closely match the reported results in the original paper. Therefore, we confirm that the warm-starting
17 gap exists in certain settings and that the Shrink-Perturb method successfully reduces or eliminates this gap. However,
18 in some cases, we were not able to completely reproduce their results. By investigating the root of such mismatches,
19 we provide another solution to avoid this gap. In particular, we show that data augmentation also helps to reduce the
20 warm-starting gap.

21 What was easy

22 The experiments described in the paper were based on regular training of neural networks on a portion of widely-used
23 datasets, possibly from a pre-trained model. Therefore implementing each experiment was relatively easy to do.
24 Furthermore, since many of the parameters were reported in the original paper, we did not need much tuning in most
25 experiments. Finally, it is straightforward to implement and use the proposed solution.

26 What was difficult

27 Though implementing each experiment is relatively simple, the numerosity of experiments proved to be slightly
28 challenging. In particular, each of the online experiments in the original setting requires training a deep network to
29 convergence more than 30 times. In these cases, we sometimes changed the settings, sacrificing granularity to reduce
30 computation time. However, these changes did not affect the interpretability of the final results.

31 Communication with original authors

32 We briefly communicated with the authors to clarify the experiments' details, such as the convergence conditions.

33 1 Introduction

34 Training large models from scratch is usually time and energy-consuming, so it is desired to have a method to accelerate
35 retraining neural networks with new data added to the training set. The well-known solution to this problem is
36 warm-starting. Warm-Starting is the process of using the weights of a model, pre-trained on a subset of the data, as the
37 starting point of training with the complete data.

38 The paper investigates the effect of warm-starting on the final model’s accuracy and identifies a generalization gap on
39 warm-started models. The paper also provides a method to mitigate this gap by shrinking the pre-trained weights and
40 adding a random perturbation.

41 In this report, we repeat the original paper’s experiments and compare them with the reported results. Also, we extend
42 the original paper results by investigating the effect of data augmentation on this phenomenon. In particular, we establish
43 that using data augmentation might be a second solution to mitigating the generalization gap.

44 We report and discuss our results in Section 2. In section 3, we detail our experimental settings and hyperparameters.

45 2 Results & Discussion

46 2.1 Warm-Starting Generalization Gap

47 Similar to the paper, we start by demonstrating the existence of a generalization gap when using warm-starting before
48 training. We use the same set of experiments used by the authors. Unless otherwise stated, we follow the settings
49 described in the paper for our experiments.

50 In particular, in the offline setting, we first train our model on half of the training data and then further train the
51 pre-trained model on the whole dataset. We compare the resulting model with a model trained from randomly initialized
52 weights. Figure 1 depicts the test accuracy of ResNet-18 [1] during training in this setting and matches Figure 1 of the
53 paper.

54 We repeat this experiment with different datasets, models, and
55 optimizers. In particular we perform experiments on CIFAR-10
56 [2], CIFAR-100 [2], and SVHN [3]. As our model, we exper-
57 iment with ResNet-18, a three-layer perceptron, and logistic
58 regression. The same models and datasets were used in the
59 original paper. For the optimizers, we compare SGD [4] and
60 Adam [5]. In this particular experiment, we compare SGD
61 with and without momentum. In the rest of this work, unless
62 explicitly stated, SGD is used without momentum. The final
63 accuracies are reported in Table 1 similar to Table 1 of the
64 original paper. Our results are similar to the paper’s results for
65 CIFAR-10 and CIFAR-100 datasets. In particular, we observe a
66 generalization gap when using a warm-started model instead of
67 training from scratch. However, we did not observe the same gap
68 on the SVHN dataset. Furthermore, we were unable to
69 obtain a reasonable accuracy with the MLP model using SGD
70 without momentum with the reported setting on SVHN. We
71 instead report the result of using 0.005 learning rate. We did
72 not perform hyperparameter tuning for the other experiments.

73 In the online setting, we follow the original paper and train our
74 model in several steps, increasing the amount of data available
75 at each step. This setting is a more accurate simulation of the real-world problems where the training data grows over
76 time.

77 We split the training data into batches of 1000 samples and start adding them, one by one, to the pool of available data.
78 We follow two different scenarios. In one scenario, we reinitialize our model randomly after each batch is added and
79 train it from scratch. In the other scenario, we continue training the model with the parameters learned in the previous
80 step.

81 After each batch is added, we continue training our model until convergence before adding the next batch. We assume
82 convergence when the model reaches 99% training accuracy. By communicating with the original paper’s authors, we
83 confirmed that this is the same condition used in the original paper.

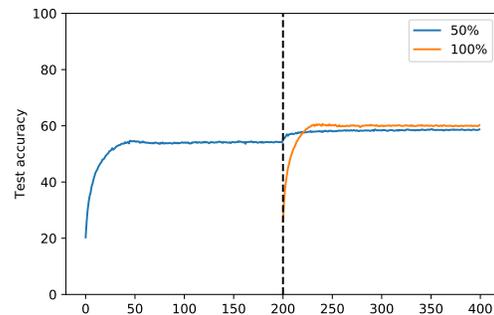


Figure 1: Test accuracy during training a ResNet-18 with SGD with and without warm-starting. The results for the randomly-initialized model has been shifted 200 epochs to overlap the part of training where the warm-started model is trained on the whole dataset.

84 As in the paper, we optimize the model using Adam optimizer with a learning rate of 0.001 on CIFAR-10. Given the
 85 discrepancy of our results on the SVHN datasets in the offline settings, we additionally perform the same experiment on
 86 this dataset. For the CIFAR-10 dataset, the generalization gap between random-initialization training and warm-start
 87 training is clearly observed (Figure 2a).

88 However, like the offline experiment, we did not observe the gap for the SVHN dataset (Figure 2b). Still, we were
 89 able to reproduce the gap by increasing the convergence accuracy threshold to 99.9% (Figure 2c). Note that 99% train
 90 accuracy is more challenging to achieve on the CIFAR-10 dataset than on SVHN and therefore requires more training,
 91 possibly leading to more over-fitting. Increasing the convergence threshold compensates for this difference.

92 This result and the fact that the authors show that the proposed
 93 Shrink-Perturb method is similar to an aggressive regularization
 94 brings up the question of whether this gap might be a side-effect
 95 of over-fitting when training on partial data. There are various
 96 known techniques to prevent overfitting. The original paper
 97 investigates the effect of some of these techniques, namely
 98 regularization, and early-stopping. We reproduced these exper-
 99 iments and explained the results below.

100 **Early-Stopping:** To investigate the effect of early-stopping,
 101 following the original paper, we trained a ResNet-18 model on
 102 half of the CIFAR-10 data and checkpointed its parameters ev-
 103 ery 20 epochs. The result is plotted in Figure 3, which matches
 104 Figure 4 of the original paper and shows the warm-starting gap
 105 can be observed even after 20 epochs of training. To decrease
 106 computational costs, we used lower granularity than the origi-
 107 nal paper to perform this experiment, saving parameters every
 108 20 epochs rather than 5 epochs. Also, we only perform each
 109 experiment once.

110 **Regularization:** Regularization is commonly used to improve
 111 generalization. The original paper explores the effect of various
 112 types of regularization. Due to time and resource limitations,
 113 we only look into weight decay, which is widely used and
 114 is a de-facto standard for training the state-of-the-art models.
 115 We repeat the offline setting experiment on CIFAR-10 with a
 116 weight decay of 0.1 on both the pre-training and main training.
 117 However, contrary to the original paper’s results, we observe
 118 that the warm-starting gap decreases when applying weight
 119 decay. We also test with weight decay values of 0.01 and 0.001.
 120 We find out that higher values of weight decay result in lower
 121 warm-starting gap. The results are reported in Table 2, which
 122 corresponds to Appendix Table 13 of the original paper.

123 **Data Augmentation:** Data augmentation is widely used to
 124 obtain state-of-the-art performance and is known to help gener-
 125 alization [6], but it is not used in the other experiments of this
 126 paper. It is specifically important to check the effect of data
 127 augmentation since it is widely used in practice. Therefore we
 128 extend the original paper’s experiments by investigating the im-
 129 pact of data augmentation. We report our results in Section 2.4.

130 2.2 Effect of Hyperparameters

131 Our results show that in some cases, a generalization gap exists
 132 when pre-training our model on a portion of the final dataset.
 133 However, when training in the online setting on SVHN, we
 134 could only observe this gap with a high enough convergence
 135 training hyperparameters, namely learning rate and batch size.

136 To investigate the effect of learning rate and batch size, we train a ResNet-18 with different values for these hyperpa-
 137 rameters. We choose the learning rate from $\{0.1, 0.01, 0.001\}$ and the batch size from $\{128, 64, 32, 16\}$. We iterate

	CIFAR-10	SGD	ADAM	MSGD
Random Init	60.3 (0.1)	80.3 (0.1)	65.5 (0.2)	
Warm Start	57.8 (0.4)	79.0 (0.1)	63.4 (0.2)	
SVHN				
Random Init	84.7 (0.1)	92.4 (0.2)	87.7 (0.1)	
Warm Start	86.3 (0.4)	93.2 (0.2)	87.2 (0.1)	
CIFAR-100				
Random Init	30.5 (0.3)	48.8 (0.2)	34.4 (0.1)	
Warm Start	27.6 (0.5)	46.1 (0.2)	30.6 (0.3)	

(a) Test accuracies for ResNet-18

	CIFAR-10	SGD	ADAM	MSGD
Random Init	38.2 (0.2)	46.5 (0.2)	45.9 (0.1)	
Warm Start	38.5 (0.3)	46.0 (0.3)	43.5 (0.4)	
SVHN				
Random Init	72.7(1.0)*	72.2 (0.8)	68.8 (1.0)	
Warm Start	67.3(1.7)*	72.5 (0.5)	70.0 (0.6)	
CIFAR-100				
Random Init	5.1 (0.2)	19.5 (0.3)	16.4 (0.1)	
Warm Start	5.1 (0.3)	18.6 (0.1)	16.6 (0.2)	

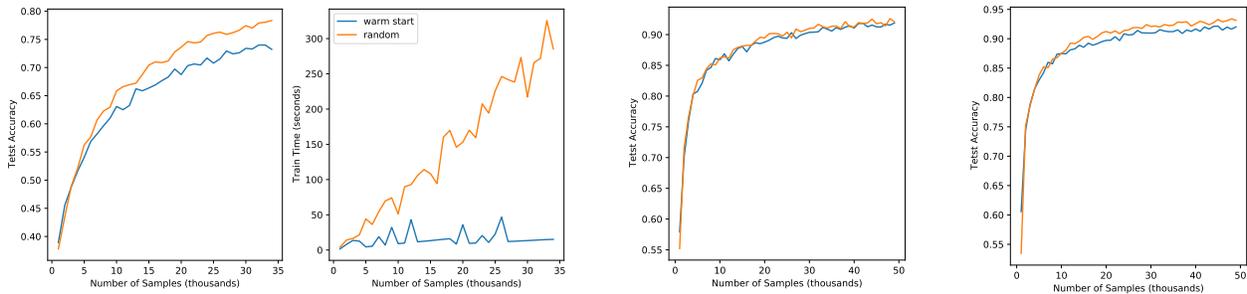
(b) Test accuracies for MLP

	CIFAR-10	SGD	ADAM	MSGD
Random Init	39.8 (0.1)	35.6 (0.3)	38.4 (0.3)	
Warm Start	39.7 (0.2)	35.2 (0.3)	38.5 (0.2)	
SVHN				
Random Init	19.8 (0.6)	24.2 (0.9)	22.4 (0.4)	
Warm Start	19.8 (0.4)	24.4 (0.8)	22.8 (0.4)	
CIFAR-100				
Random Init	16.6 (0.2)	12.6 (0.2)	17.2 (0.2)	
Warm Start	16.5 (0.1)	12.2 (0.1)	17.0 (0.1)	

(c) Test accuracies for Logistic Regression

Table 1: Test accuracies for various datasets and models and optimizer for warm-start training and training from random initialization. We use an MLP with Tanh activation with 3 hidden layers of 100 neurons. A different learning rate was used for cells marked with a star (*).

The paper investigates the effect of other



(a) For the CIFAR-10 dataset, the generalization gap between randomly initialized and warm started models is observed for the 99.0% convergence threshold. However, the training time required at each step increases almost linearly in the case of randomly initialized model, but is constant when using warm-starting.

(b) For the SVHN dataset, there is no generalization gap between randomly initialized and warm started models with the 99.0% convergence threshold.

(c) For SVHN, the generalization gap appears by increasing the convergence threshold to 99.9%.

Figure 2: Online learning experiment on CIFAR-10. the horizontal axis shows the number of samples available to train the model.

	0.1	0.01	0.001
Random Initialization	81.44	62.73	60.42
Warm Starting	81.16	61.22	58.63

Table 2: Accuracy of training a ResNet-18 with and without warm-starting for different values of weight decay.

138 over all pairs for these values. For each pair, we train over the full CIFAR-10. We also train a different model over
 139 50% of the CIFAR-10 dataset and use it to warm-start a model and train it on the whole dataset. We use a different
 140 learning rate and batch size, randomly chosen from the sets of values, in the second part of the training, i.e., for training
 141 the warm-started model. We repeat each experiment 9 times. Each model is trained to 99% training accuracy. The
 142 test accuracy is plotted against training time in Figure 4, which corresponds with Figure 3 of the paper. Note that the
 143 training time for the warm-started model corresponds to the time of the second part of the training. In other words, the
 144 time of training on half of the dataset is not included. This is justified because the goal is to assess if warm-starting
 145 leads to comparable accuracy while saving training time when a new batch of data arrives. In our results, choosing the
 146 right hyperparameters can lead to achieving comparable or even better accuracy, when using warm-starting, faster than
 147 training a randomly initialized model. This does not match the results of the paper, where the warm-started models with
 148 comparable accuracy take the same amount of training time as the randomly initialized models. While we perform less
 149 experiments in the warm-started setting, we perform the same number of experiments with random initialization as
 150 described in the original paper’s text. However, the number of points in Figure 3 of the original paper corresponding to
 151 randomly initialized models, is more than what has been described in the text. We note that performing more randomly
 152 initialized experiments might be the reason for the mismatch in our results with the original paper.

153 2.3 Shrink & Perturb Solution

154 In addition to establishing the warm-starting gap’s existence and investigating its roots, the paper also provides a method
 155 to mitigate this issue. In this method, the training starts from a shrunk and perturbed version of the pre-trained weights,
 156 so we reference it as the Shrink-Perturb method. More specifically, for a given λ and σ , the new weight is computed as

$$w_{new} = \lambda w_{pretrained} + \sigma w_{random} \quad (1)$$

157 where w_{old} is the pre-trained weight and w_{random} is the corresponding weight from a randomly initialized model.
 158 Whenever we apply the Shrink-Perturb transformation, we create a new randomly initialized model and use its weights
 159 as w_{random} .

160 We tested the effectiveness of this method in both offline and online settings. In the offline setting, we applied the
 161 Shrink-Perturb transform after pre-training on 50% of CIFAR-10. We used $\sigma = 10^{-4}$ and repeated this experiment
 162 with different values of λ . We plotted the test accuracy during training on all of the data in Figure 5. It can be seen that
 163 the method is effective and leads to even better performance than the randomly initialized model.

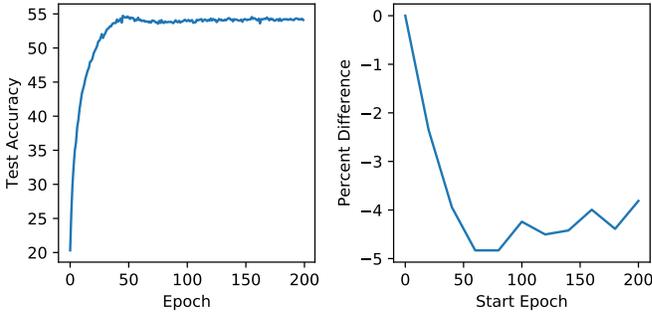


Figure 3: **Left:** Test accuracy while training on half of CIFAR-10. **Right:** Plot of test accuracy damage, as percentage difference from random initialization, against number of warm-starting epochs.

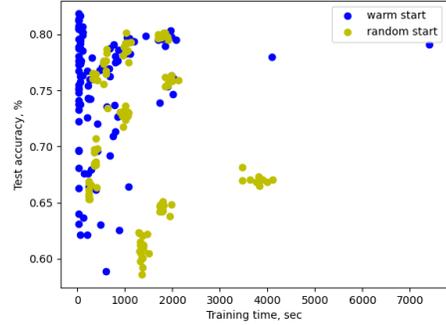


Figure 4: Training time vs Test accuracy for randomly initialized and warm-started models

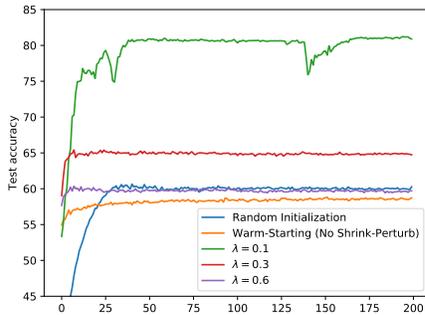


Figure 5: Test accuracy during training a ResNet-18 with SGD with warm starting and Shrink-Perturb and without it

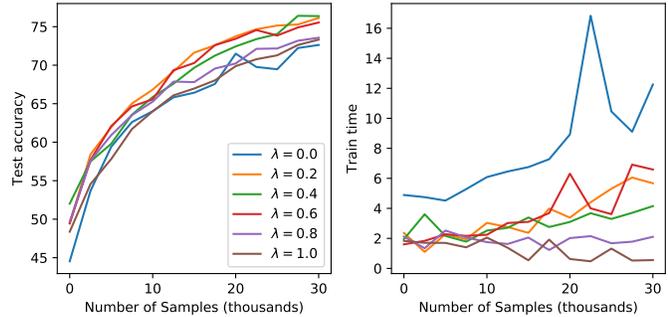


Figure 6: Test accuracy during training a ResNet-18 in the on-line setting while applying Shrink-Perturb method with different values of λ .

164 In the online setting, we applied the Shrink-Perturb transform every time a new batch of data is added. To reduce
 165 computation cost, we add data in batches of 2500 samples. The result for $\sigma = 10^{-4}$ and different values of λ is plotted
 166 in Figure 6. Figure 6 matches Figure 7 of the original paper. The result of applying the Shrink-Perturb method in the
 167 offline setting is not reported in the original paper.

168 To assess the impact of shrinking weights on the model’s performance,
 169 we fit different models to CIFAR-10. Then, we shrink the weight
 170 with different values of λ and evaluate the accuracy. Similar to the
 171 paper, we train ResNet18 and an MLP with ReLU activation with
 172 and without bias. In addition, we also train an MLP with Tanh
 173 activation with and without bias. The result is shown in Figure 7,
 174 which corresponds with Figure 6 of the paper. The only difference
 175 in our findings with the original paper’s is that we observe classifier
 176 performance damage for MLP with ReLU for $\lambda > 0.6$. Though
 177 for $\lambda > 0.8$ the damage is negligible. Also, note that shrinking the
 178 weights of an MLP without bias and ReLU activation only scales the
 179 final output, which does not affect the output labels. Therefore its
 180 immunity to shrinkage is expected. The more interesting result is that
 181 even for ResNet-18 or MLP with Tanh activation, the test accuracy
 182 is not significantly damaged for λ values greater than 0.2.

183 In order to explain why the Shrink-Perturb method is effective, the
 184 original paper compares the average gradients over the first and sec-

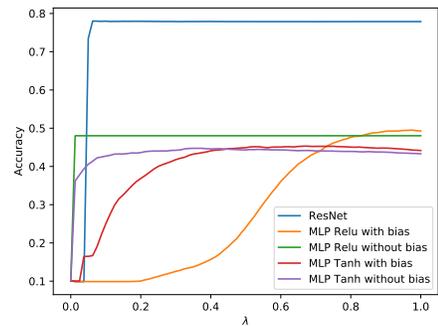


Figure 7: Test accuracy of trained models for different shrinkage coefficients λ .

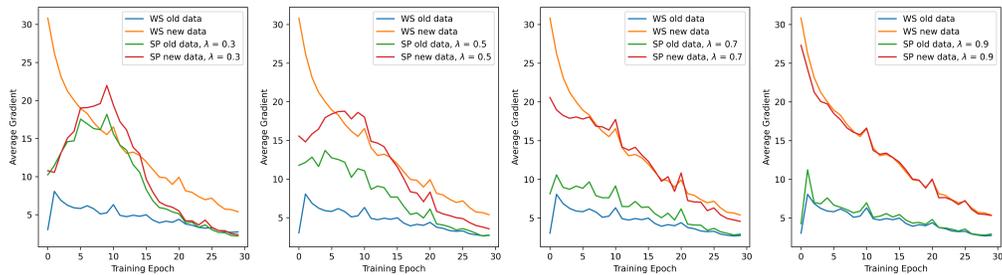


Figure 8: Average over old data or new data of L2 norm of gradient with respect to model parameters. WS refers to warm starting without and SP to warm starting with Shrink-Perturb.

185 ond half of the dataset during the training of the warm-started model in the offline setting. In particular, a ResNet-18 is
 186 first trained on half of CIFAR-10. Warm-starting from the pre-trained model, the model is trained on the full dataset
 187 while measuring the average gradient over the first and second half of the dataset simultaneously. It is observed that the
 188 gradient for the first half, which the model was pre-trained on, is substantially lower than for the second part. However,
 189 applying the Shrink-Perturb transformation eliminates this difference. We reproduced this experiment with some slight
 190 modifications. In particular, instead of averaging the gradient over part of the dataset after each batch, we did it at the
 191 beginning of each epoch. We plotted these values in Figure 8. Our result matches Figure 5 of the original paper. In
 192 particular, we confirm that the Shrink-Perturb method successfully eliminates the gap between the gradients.

193 2.4 Effect of Data Augmentation

194 Data augmentation is widely used in practice. However, it is not used for the experiments in the original paper.
 195 Therefore, we decided to assess its impact on the generalization gap for warm starting training.

196 We perform our experiments on ResNet-18 and CIFAR-10. To augment the data, we first pad the image with 4 pixels on
 197 each side and then randomly crop it back to 32x32. We then perform a random horizontal flip with probability 0.5. We
 198 also apply color jitter with brightness, contrast, and saturation factor equal to 0.25. Finally, we also apply small random
 199 rotations. All experiments were done with SGD and a learning rate equal to 0.001 in order to make the setup consistent
 200 with previous warm start experiments. The results are reported on Figure 9 .

201 It can be seen that applying the augmentation mitigates the warm-starting gap. We allow the models to train for 350
 202 epochs. However, because the learning rate is low, the models are not fully converged even after 350 epochs. We did not
 203 continue the training because of resource limitations. However, it is visible that warm-starting with data augmentation
 204 can achieve good performance faster than training from scratch.

205 To explain the effectiveness of the Shrink-Perturb solution, the original paper’s authors looked at the differences of the
 206 gradient norm for the first and the second part of the dataset, which is heavily reduced after applying Shrink-Perturb (as
 207 shown in Figure 8). Following the same direction, we checked if applying data augmentation can affect the difference
 208 as well. It can be seen in Figure 10 that, similar to the shrink perturb method, the gradient norm difference is also
 209 mitigated when using data augmentation.

210 It is clear that applying data augmentation prevented overfitting. The original warm start setup has a large divergence
 211 between train and test accuracy from the beginning of the training. On the contrary, the model trained with data
 212 augmentation has close performance on train and test datasets. We leave the investigation of other overfitting prevention
 213 techniques’ effects as future work.

214 Additionally, we note that data augmentation usually slows down the convergence and it cannot be applied to every task
 215 since for some types of data transform set cannot be defined. Due to the limits of this report, we also leave the careful
 216 comparison between data augmentations and Shrink & Perturb as future research in this area.

217 2.5 Warm-Starting Gap in Transfer Learning

218 Deep learning models require large training sets to perform well. This presents a problem in many practical cases
 219 where only limited data is available, and acquiring additional data is expensive. This has encouraged the use of transfer
 220 learning [7; 8; 9]; the practice of warm-starting from a model trained on a different dataset.

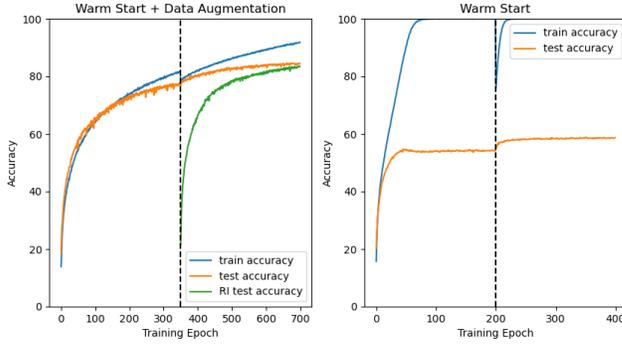


Figure 9: Test and train accuracy of models in warm start setting with and without data augmentation. RI refers to Random Initialization.

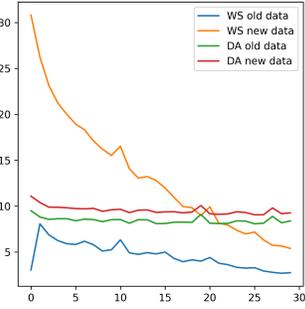


Figure 10: Gradients with respect to the first and the second half of the dataset, DA refers to Data Augmentation

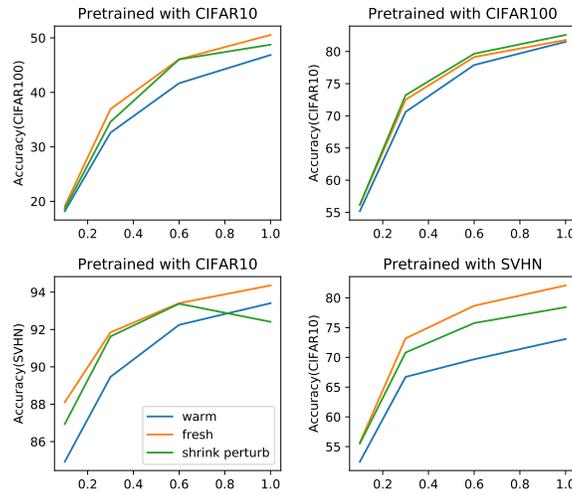


Figure 11: Test accuracy against the available portion of data p for training a model on one dataset with warm-starting from another dataset.

221 To investigate whether a similar gap is observed in transfer learning, we trained a ResNet-18 model on one dataset and
 222 used the pre-trained weights to warm-start training on a different dataset. We performed this experiment for all pairs of
 223 CIFAR-100, CIFAR-10, and SVHN datasets. To also investigate the effect of the amount of the data available, we also
 224 considered subsets of these datasets where only a fraction p of data is available. More accurately, for every two datasets
 225 and $p \in \{0.1, 0.3, 0.6, 1.0\}$, we chose a random subset from each dataset containing $p \times n$ data points, where n is the
 226 total number of data points in that dataset. We performed the described transfer learning experiment for these subsets
 227 and recorded the final test accuracy. To assess the effect of warm-starting and the Shrink-Perturb method, we plotted
 228 the final test accuracy of each of the described three settings (random initialization, warm starting, and warm-starting
 229 with Shrink-Perturb) with respect to p in Figure 11. In this experiment, we used Adam as our optimizer. This figure
 230 corresponds with Figure 9 of the original paper.

231 At each part of the training, we train our models for 200 epochs. When training CIFAR-100 starting from weights of a
 232 model trained on SVHN or CIFAR-10, the last layer is initialized randomly because of the mismatch in the number of
 233 classes. It can be seen that, as mentioned in the paper, the warm-starting gap exists in the transfer learning settings as
 234 well, and that it is worsened when the amount of data available is increased. Furthermore, the Shrink-Perturb method
 235 proves useful in this setting, as well.

236 3 Methodology

237 In this section, we define the setting we used for our experiments. There was no available code for the original paper,
 238 and we implemented everything from scratch. We use the PyTorch framework for the implementations.

239 3.1 Model descriptions

240 Most of the experiments are performed using ResNet-18 [1]. Some experiments are also performed on a Multi-Layered
241 Perceptron (MLP) and Logistic Regression. We detailed the structure of each of these models below.

- 242 • **ResNet-18:** We used an implementation of ResNet-18 tuned for CIFAR-10 dataset. We used the code from
243 https://github.com/huyvnphan/PyTorch_CIFAR-10. In all experiments, batch normalization [10] was
244 enabled.
- 245 • **MLP:** The MLP has three hidden layers, each of which has 100 neurons. Either ReLU or Tanh was used as
246 the activation function. Unless explicitly stated, the bias term is added.
- 247 • **Logistic Regression:** We implement Logistic Regression as a Multi-Layered Perceptron with no hidden
248 layers.

249 We used either Adam or SGD optimizers for training the models. More accurately, we use Adam in Table 1, Figure 2,
250 Figure 6, and Figure 11. In all other experiments, we use SGD. We use 0.001 learning rate and batches of size 128.
251 Unless otherwise stated, we used SGD without momentum and without weight decay. In cases where momentum was
252 used (such as in Table 1), the value of momentum was set to 0.9. The Adam optimizer was used with default parameters
253 from PyTorch’s implementation, namely $\beta_1 = 0.9$, and $\beta_2 = 0.999$.

254 3.2 Datasets

255 Same as the original paper, we perform experiments on CIFAR-10, CIFAR-100, and SVHN datasets. We normalize
256 each of the RGB channels by the mean and standard deviation of that channel in the CIFAR-10 dataset. Except for the
257 data augmentation experiments, we do not apply any data augmentation.

258 3.3 Hyperparameters

259 We used hyperparameters stated in the original paper in most of our experiments. In cases where we deviated from the
260 reported values, mostly due to computational resource and time limitation, we have reported them in the text where we
261 described the experiment. In case a hyperparameter is not reported in the original paper, we either communicated with
262 the authors to ask the hyperparameters, pick a value making reasonable assumptions, or try out different values and
263 report the result for all of them. In all these cases, we clarified the parameter we used in the text.

264 3.4 Experimental setup

265 We ran our experiments on both public cloud infrastructure, such as Google Colab and private GPUs that were
266 available to us. Therefore the infrastructure varies between different experiments. Our implementations for all the
267 experiments in this work is available in the Supplementary Material and also in [https://github.com/CS-433/
268 cs-433-project-2-fesenjoon](https://github.com/CS-433/cs-433-project-2-fesenjoon).

269 4 Communication with Authors

270 In the original paper [11], it was not clear what convergence condition was used to stop the training. Therefore, We
271 communicated with the authors via email and asked them to explain the convergence condition used in the experiments
272 more clearly. They stated that convergence happens when the training accuracy reaches 99%. However, for reproducing
273 the Table 1 which also uses simpler models like Logistic Regression and MLP, it is not possible to reach 99% accuracy.
274 They clarified that in this scenario, the convergence condition is met when the training accuracy stops improving.

275 5 Conclusion

276 We have verified the existence of the generalization gap in certain training settings. Additionally, we have confirmed
277 that the introduced Shrink-Perturb method can be effective in removing this gap. We did this by repeating experiments
278 of the original paper and performing some experiments of our own. However, we also encountered cases where we
279 were not able to reproduce the warm-starting gap or where the Shrink-Perturb method was not very successful. In
280 addition, we reproduced several experiments to investigate the effect of hyper-parameters, such as learning rate, on
281 this phenomenon. Finally, we have shown that applying data augmentation can also help to remove this gap. To allow
282 others to reproduce our results, we have detailed our experiments and have released our code.

283 **References**

- 284 [1] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” in *2016 IEEE Conference*
285 *on Computer Vision and Pattern Recognition (CVPR)*. Las Vegas, NV, USA: IEEE, Jun. 2016, pp. 770–778.
286 [Online]. Available: <http://ieeexplore.ieee.org/document/7780459/>
- 287 [2] G. E. Hinton, “Learning multiple layers of representation,” *Trends in Cognitive Sciences*, vol. 11, no. 10, pp.
288 428–434, Oct. 2007. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S1364661307002173>
- 289 [3] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, “Reading digits in natural images with
290 unsupervised feature learning,” in *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*,
291 2011. [Online]. Available: http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf
- 292 [4] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning representations by back-propagating errors,” *Nature*,
293 vol. 323, no. 6088, pp. 533–536, Oct. 1986. [Online]. Available: <http://www.nature.com/articles/323533a0>
- 294 [5] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd International Conference on*
295 *Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*,
296 Y. Bengio and Y. LeCun, Eds., 2015. [Online]. Available: <http://arxiv.org/abs/1412.6980>
- 297 [6] C. Shorten and T. M. Khoshgoftaar, “A survey on image data augmentation for deep learning,” *Journal of Big*
298 *Data*, vol. 6, no. 1, p. 60, 2019.
- 299 [7] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, “Learning and Transferring Mid-level Image
300 Representations Using Convolutional Neural Networks,” in *2014 IEEE Conference on Computer Vision*
301 *and Pattern Recognition*. Columbus, OH, USA: IEEE, Jun. 2014, pp. 1717–1724. [Online]. Available:
302 <https://ieeexplore.ieee.org/document/6909618>
- 303 [8] J. Huang, J. Li, D. Yu, L. Deng, and Y. Gong, “Cross-language knowledge transfer using multilingual deep neural
304 network with shared hidden layers,” in *2013 IEEE International Conference on Acoustics, Speech and Signal*
305 *Processing*, May 2013, pp. 7304–7308, iSSN: 2379-190X.
- 306 [9] M. Long, H. Zhu, J. Wang, and M. I. Jordan, “Unsupervised domain adaptation with residual transfer networks,”
307 in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, ser. NIPS’16.
308 Barcelona, Spain: Curran Associates Inc., Dec. 2016, pp. 136–144.
- 309 [10] S. Ioffe and C. Szegedy, “Batch Normalization: Accelerating Deep Network Training by Reducing Internal
310 Covariate Shift,” in *International Conference on Machine Learning*. PMLR, Jun. 2015, pp. 448–456. [Online].
311 Available: <http://proceedings.mlr.press/v37/ioffe15.html>
- 312 [11] J. T. Ash and R. P. Adams, “On warm-starting neural network training,” 2020.