

COMPARING FINE-TUNING AND REWINDING IN NEURAL NETWORK PRUNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Neural network pruning is a popular technique for reducing inference costs by removing connections, neurons, or other structure from the network. In the literature, pruning typically follows a standard procedure: train the network, remove unwanted structure (*pruning*), and train the resulting network further to recover accuracy (*fine-tuning*).

In this paper, we explore an alternative to fine-tuning: *rewinding*. Rather than continuing to train the resultant pruned network (fine-tuning), *rewind* the remaining weights to their values from earlier in training, and re-train the resultant network for the remainder of the original training process. We find that this procedure, which repurposes the strategy for finding lottery tickets presented by Frankle et al. (2019), makes it possible to prune networks further than is possible with fine-tuning for a given target accuracy, provided that the weights are rewound to a suitable point in training. We also find that there is a wide range of suitable rewind points that achieve higher accuracy than fine-tuning across all tested networks. Based on these results, we argue that practitioners should explore rewinding as an alternative to fine-tuning for neural network pruning.

1 INTRODUCTION

Pruning (Le Cun et al., 1990; Hassibi et al., 1993; Han et al., 2015; Li et al., 2017; Liu et al., 2019) is a popular set of techniques that reduce the memory or computational demands of a neural network by removing weights, filters, neurons, or other structures while maintaining the network’s accuracy. In practice, pruning can reduce the parameter counts of large-scale networks for vision and NLP tasks by an order of magnitude while maintaining the accuracy of the original network (Han et al., 2016). Pruning allows for accelerating neural network inference: smaller neural networks can fit on smaller devices, and in some situations can use fewer FLOPS, reducing energy/power consumption and inference latency (Han, 2017).

Many state-of-the-art pruning methods follow a standard *fine-tuning* framework: (1) train the network, (2) prune a given percentage of network structures according to a heuristic, and (3) *fine-tune* the resulting network with additional training for several epochs to recover its accuracy. Instantiated with a simple heuristic that prunes weights with the lowest magnitude, this framework can reduce the size of common vision networks by an order of magnitude with no loss in accuracy (He et al., 2018; Liu et al., 2019; Gale et al., 2019).

Other techniques use variations of this framework, such as *iterative pruning* (Han et al., 2015), which repeats steps 2 and 3 iteratively, alternating between pruning and fine-tuning the network until reaching a desired sparsity or accuracy level. This is in contrast with *one-shot pruning*, which prunes to its target sparsity in a single step.

In practice, finding smaller or more accurate pruned networks requires more epochs of fine-tuning, thereby increasing pruning’s computational cost. However, this investment may be justified if the network will be used extensively for inference, or needs to reach a particular size to run on certain platforms (e.g. mobile devices). Therefore, the process of pruning a neural network necessitates a careful tradeoff between three criteria:

ACCURACY. The accuracy of the resultant pruned network.
INFERENCE COST. The cost of performing inference with the pruned network.
SEARCH COST. The cost of finding the pruned network (i.e. the amount of training required).

Many pruning techniques focus exclusively on finding the sparsest network that matches the accuracy of the original network (Han et al., 2015; Molchanov et al., 2017; He et al., 2018). However, other important regions exist in the tradeoff space – for instance, trading off ACCURACY for even more sparsity (decreasing INFERENCE COST), or constraining the cost of finding the pruned network (SEARCH COST) to be below some threshold.

In this paper, we demonstrate a technique, *rewinding*, that can find smaller networks (lower INFERENCE COST) than fine-tuning for any given target ACCURACY. Further, given the same target sparsity (INFERENCE COST) and total epochs of training (SEARCH COST), rewinding is able to produce networks that are at least as accurate as those produced by fine-tuning, provided a suitable point to rewind to.

Rewinding. The lottery ticket hypothesis (Frankle et al., 2019) demonstrates that typical neural networks contain far smaller subnetworks that emerge early in training (i.e., in the first 10% of the way through training), and can train in isolation to the same accuracy as the full network. To find these subnetworks, Frankle et al. propose a variation on the fine-tuning framework: (1) train the network to convergence at final epoch T , (2) prune structures according to a heuristic, (3) *rewind* the remaining weights to their values from an intermediate epoch t in training, and (4) re-train the pruned and rewind subnetwork for the final $T - t$ epochs of training, using the same training parameters and learning rate schedule as was used during the original training. Rewinding (which Frankle et al. use to find lottery tickets) offers an alternative framework to fine-tuning for neural network pruning.

Our contribution. In this paper, we demonstrate the value of rewinding as a general pruning framework and compare rewinding and fine-tuning on CIFAR-10 and ImageNet networks. We study rewinding to epochs all throughout training, and show that even rewinding to later in training (which Frankle et al. do not study) can result in more accurate networks than fine-tuning.

Rewinding expands the frontier of the tradeoff space between the pruning criteria: the standard network-agnostic rewinding framework we use is able to prune a ResNet-50 by 79% with no loss in accuracy, matching more complex and computationally expensive state-of-the-art results (He et al., 2018; Gale et al., 2019). Further, rewinding is able to find sparser neural networks for any given target accuracy than fine-tuning can find.

Rewinding involves picking a suitable rewind point. We show that across all tested networks, rewinding achieves higher accuracy than fine-tuning across a wide range of rewind points. For a fixed re-training budget, desired sparsity, and given a rewind point between 30% to 75% of the way through training, rewinding finds more accurate pruned networks than fine-tuning.

Based on these results, we argue that practitioners should explore rewinding as an alternative to fine-tuning for neural network pruning.

2 METHODOLOGY

We select a variety of standard vision networks, and train and prune them with standard hyperparameters. We collect 3 independent trials of all experiments.

2.1 EXPERIMENTS

Rewinding experiments. Rewinding introduces one hyperparameter: the epoch t that the network is rewind to. For example, consider a network that is normally trained for T epochs. We train the network for the original T epochs, prune, rewind the weights to their values at epoch t , and re-train the resultant network for the remaining $T - t$ epochs of training. When rewinding weights to a specific epoch, we also rewind the learning rate schedule to the same epoch. The rewinding algorithm is presented in Appendix A, Algorithm 1.

Fine-tuning experiments. There are several hyperparameters that must be chosen when pruning with fine-tuning, namely the number of epochs t to fine-tune for, and the learning rate schedule during that re-training. In our experiments, we explore many values of t , which controls the cost of obtaining the final, pruned network (SEARCH COST). We follow the standard practice of fine-tuning using the final learning rate used during the original training phase (Liu et al., 2019). The fine-tuning algorithm is presented in Appendix A, Algorithm 2.

Reinitialization experiments. We include a baseline of reinitializing and re-training the pruned network from scratch, for the same total number of epochs as the rewinding experiments and the fine-tuning experiments. This reinitialization baseline is from Liu et al. (2019). In the iterative case, it is not within our computation budget to reinitialize and re-train every distinct discovered network structure at every sparsity level: reinitialization re-trains a masked subnetwork, and iterative pruning generates a different mask for each different rewind point and sparsity level (as opposed to the singular mask per sparsity level generated by one-shot pruning). Instead, we use two approaches, and report the best performing of the two. First, we use the sparsity discovered by the best-performing network from the previous iterative pruning iteration, under the assumption that reinitializing from the best-performing network structure will result in the best accuracy. Second, we simply use a one-shot pruning and reinitialization from the original network. The reinitialization algorithm (following the choice of using one-shot pruning) is presented in Appendix A, Algorithm 3.

2.2 PRUNING FRAMEWORK

Pruning heuristic. We study a sparse pruning heuristic in which we remove individual weights with the lowest magnitudes from the network, following Han et al. (2016); Zhu & Gupta (2017); Gale et al. (2019); Liu et al. (2019); and Frankle & Carbin (2019). We prune all weights in the network *globally* (Lee et al., 2019; Frankle & Carbin, 2019), removing weights with the lowest magnitudes without consideration for the layer in which the weight resides.

Iterative pruning. Pruning frameworks can be applied iteratively, repeatedly pruning and fine-tuning or rewinding to gradually compress the network. In our iterative pruning experiments, we follow the strategy from Han et al. (2016) and Frankle & Carbin (2019) to remove 20% of weights per iteration, repeating this process as necessary to reach the desired level of sparsity.

2.3 NETWORKS AND HYPERPARAMETERS

Training hyperparameters. We study both CIFAR-10 (Krizhevsky & Hinton, 2009) and ImageNet (Russakovsky et al., 2015) networks. All networks are trained with SGD with Nesterov momentum with $\beta = 0.9$. We use the standard ResNet training schedule for CIFAR-10, which trains for 182 epochs with a batch size of 128, using an initial learning rate of 0.1, decaying by a factor of 10 at epochs 91 and 136 (He et al., 2016). On ImageNet, we use the schedule from Goyal et al. (2017) with a batch size of 1024, training for 90 epochs with an initial learning rate of 0 linearly scaling up to 0.4 over the first 5 epochs, then decaying by a factor of 10 at epochs 30, 60, and 80.

Networks. We use a standard implementation of a ResNet-20 (He et al., 2016) provided alongside Tensorflow (Abadi et al., 2016), available online¹. With the CIFAR-10 training schedule, this implementation trains the ResNet-20 to $91.58\% \pm 0.20\%$ accuracy.

We use a slight modification of the VGG-16 network as described by Simonyan & Zisserman (2014). Our VGG-16 includes Batch Normalization (Ioffe & Szegedy, 2015) after each convolution and before each activation. Also, rather than multiple fully connected layers at the end, we use a single fully connected layer. This implementation trains to $93.24\% \pm 0.19\%$ accuracy with our CIFAR-10 training schedule.

We use a standard implementation of a ResNet-50 provided alongside Tensorflow, available online². With the ImageNet training schedule, this implementation trains the ResNet-50 to $76.13\% \pm 0.13\%$ accuracy, in line with the baseline accuracy reported in the state-of-the-art pruning approaches (He et al., 2018; Gale et al., 2019).

¹<https://github.com/tensorflow/models/tree/v1.13.0/official/resnet>

²<https://github.com/tensorflow/tpu/tree/98497e0b/models/official/resnet>

Network	Pruning Method	Finetune Sparsity	Rewinding Sparsity	Ratio
ResNet-20	One-shot	36.00%	67.23%	1.95×
ResNet-20	Iterative	67.23%	83.22%	1.95×
VGG-16	One-shot	94.50%	96.48%	1.56×
VGG-16	Iterative	94.50%	98.85%	4.77×
ResNet-50	One-shot	59.04%	73.79%	1.56×
ResNet-50	Iterative	48.80%	79.03%	2.44×

Table 1: Sparsest networks with no accuracy loss from the original network (max across all trials). The ratio column is the ratio between the *density* (inverse of sparsity) of the fine-tuning network and the density of the rewinding network.

2.4 LIMITATIONS OF EXPERIMENTS

This paper studies only vision networks with a single pruning technique (unstructured magnitude pruning), comparing against a single variant of the fine-tuning framework. We do not consider any variations on this fine-tuning framework, such as techniques that gradually prune throughout training. We also do not consider structured pruning methods that produce structurally dense networks more easily amenable to acceleration.

Our rewinding framework also precludes some potentially useful techniques, such as disentangling the rewinding epoch from the number of re-training epochs, e.g. training for shorter or longer than the original training schedule or using a different learning rate. We believe that other tradeoffs between the pruning criteria can be made by exploring this space and tuning the re-training schedule, but do not explore these modified schedules in this paper in order to provide a fair comparison against fine-tuning.

3 SPARSE NETWORKS

3.1 MATCHING THE ORIGINAL ACCURACY

Rewinding can find sparser networks (INFERENCE COST) that meet that the accuracy of the original network (ACCURACY) than fine-tuning can. This metric follows standard approaches to pruning which allow for large budgets to produce the most accurate and efficient network (Han et al., 2015): pruning is often motivated as producing networks which are just as accurate as the original, unpruned network, allowing for an unbounded SEARCH COST.

Methodology. To find the sparsest network with no loss in accuracy achievable by a given technique, we ran a parameter sweep across 10 evenly distributed re-training epochs between 0 and the original training time T of the network (90 for ImageNet networks and 182 for CIFAR-10 networks) for rewinding, and between 0 and $2T$ for fine-tuning and reinitialization (allotting fine-tuning up to $2\times$ the re-training time as rewinding if it needs longer to train). We report the sparsity level at which at least one of the three trials for a given technique and number of re-training epochs matched the accuracy of the original unpruned network.

Results. Table 1 shows the sparsest network achievable with no loss in accuracy (i.e. a resultant accuracy equal to that of the original network). On all networks considered, rewinding finds a sparser network without loss in accuracy than fine-tuning finds. Beyond just outperforming fine-tuning, these results are competitive with state-of-the-art pruning approaches tailored specifically for ResNet-50: using a bespoke, highly tuned framework which prunes by specific layerwise rate at specific points throughout the regular training schedule, Gale et al. (2019) achieve 80% sparsity with a 0.17% accuracy drop on a ResNet-50. Using reinforcement learning to find layerwise pruning rates (rather than the simple global magnitude heuristic presented in this paper), He et al. (2018) achieve 80% sparsity with a 0.02% drop in accuracy on a ResNet-50. Our approach, using standard pruning hyperparameters (Han et al., 2015) and with just the addition of rewinding instead of fine-tuning, achieves 79.03% sparsity (7 iterations of pruning by 20%) with a 0.06% increase in accuracy.

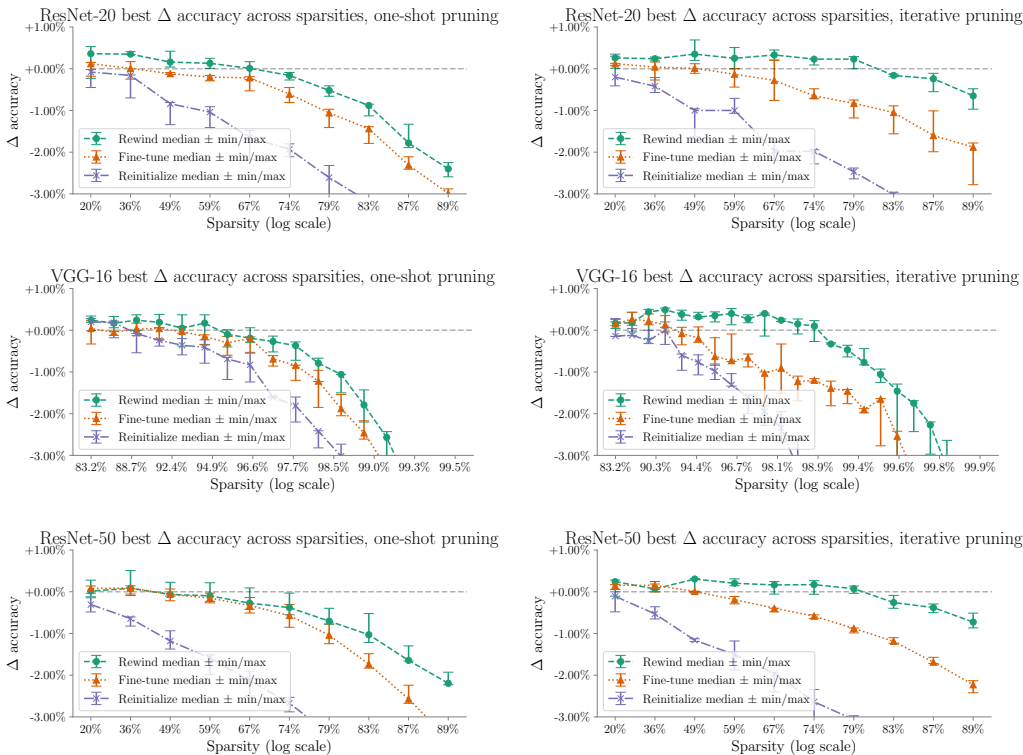


Figure 1: The most accurate networks achievable by different pruning and re-training methods across different target sparsity levels.

3.2 PRUNING WITH AN ACCURACY TARGET

Beyond just producing sparser networks that match the original accuracy, rewinding can find sparser networks that match any accuracy threshold.

Methodology. Using the same experimental methodology as Section 3.1, we plot the most accurate achievable networks across a wide range of target sparsities in Figure 1. Accuracies reported in this section are the median accuracy of the best performing set of trials across all re-training epochs in the parameter sweep (e.g. the median of the accuracies for all 3 trials of re-training for 80 epochs, rather than the best 3 trials across all choices of re-training epochs).

As a baseline, we also consider re-initializing the discovered pruned network and re-training it by extending the original training schedule to the same total number of training epochs as fine-tuning trains for. Liu et al. (2019) found that for many pruning techniques, pruning and fine-tuning results in the same or worse performance as simply training the pruned network from scratch for an equivalent number of epochs. To address these concerns, we include comparisons against random reinitializations of networks with the discovered pruned structure, trained for the original T training epochs plus the extra t epochs that networks were re-trained for. The methodology for these experiments is detailed in Section 2.1.

Results. Rewinding always achieves either equivalent or higher accuracy than fine-tuning, with a particularly distinct advantage at higher sparsities. At lower sparsities, the techniques perform equivalently well, recovering the original accuracy of the unpruned network. At higher sparsities, the techniques result in different accuracies: rewinding consistently achieves higher accuracy than fine-tuning, which in turn outperforms reinitialization. These results are consistent across all tested networks and pruning methods.

4 PICKING THE REWIND EPOCH

To use rewinding when pruning a neural network, one must first select a point in training to rewind the network’s weights to. While rewinding achieves higher accuracy than fine-tuning for many choices of rewind points, it does not do so for all.

Rewinding has regions of low performance. Rewinding to the very beginning of training results in significantly lower accuracy than rewinding to a small amount into training (Frankle et al., 2019). Similarly, rewinding to a point too late through training means that the full accuracy benefits of rewinding will not be realized: our rewinding framework does not allow adding on additional training time after the fact.

Further, the choice of rewind point is permanent: selecting a different rewind point to achieve better accuracy requires re-training from the new rewind point. On the other hand, one can continue fine-tuning to achieve its best available accuracy: it is easy to continue fine-tuning indefinitely until either a desired accuracy is met or learning plateaus. A given rewind point imposes a floor on re-training time in order to get the accuracy benefits, along with a ceiling on the accuracy achievable from that rewind point.

Methodology. To quantitatively determine where to rewind to, we consider two different metrics:

- *Safety*: a measure of the accuracy gap between rewinding to a given point and fine-tuning with an equivalent re-training budget. More formally, the safety of a rewind point t and a sparsity s is defined as the difference between the accuracy achievable by pruning to s sparsity, rewinding to epoch t in training, and re-training for exactly $T - t$ epochs (where T is the original training time of the network), and the best accuracy achievable by fine-tuning for at most $T - t$ epochs.
- *Dominance*: a measure of the accuracy gap between rewinding to a given point and fine-tuning given an unlimited budget. Formally, the dominance of a rewind point t and a sparsity s is defined as the accuracy achievable by pruning to s sparsity, rewinding to epoch t and re-training for $T - t$ epochs, and the maximum accuracy achievable by pruning to s sparsity and fine-tuning for any amount of time (up to $2 \times$ the original training time).

Results. Figure 2 and Figure 3 show the safety and dominance respectively of rewinding across networks, pruning methods, and sparsity levels. In these plots, the y axis shows the rewind point as a percentage of the original training time of the network – i.e., 0% is rewinding to the beginning of training, and 100% is rewinding to the end. Red points mean that rewinding achieves higher accuracy for that sparsity and re-training budget; blue means that fine-tuning achieves higher accuracy.

An alternative presentation of these data, instead showing accuracy curves in a form similar to Figure 1, can be found in Appendix C.

Analysis. At medium to high sparsities, rewinding to anywhere between 30% to 75% of the way through training consistently outperforms fine-tuning for an equivalent number of re-training epochs (Figure 2). The cross-network dominant region extends from 30% to 60% of the way through training (Figure 3): essentially, rewinding to anywhere in that region and re-training results in higher accuracy than any amount of fine-tuning. Given a rewind point, higher accuracy can often (though not always) be attained by searching for better-performing rewind points closer to the beginning of training (i.e. moving down on Figure 3); re-training time can be reduced, at the cost of accuracy, by rewinding closer to the end of training and re-training for fewer epochs (i.e. moving up on Figure 3).

These regions vary across networks and datasets: for instance, ResNet-50 has a much broader safe and dominant region than VGG-16. The regions also vary across sparsity levels: rewinding performs better relative to fine-tuning on sparser networks, a result also seen in Section 3.2, Figure 1.

Edge cases. The top-most part of the plots of safety in Figure 2 have no strong tendency towards either fine-tuning or rewinding. Fine-tuning and rewinding have approximately the same behavior in the limit of few to no re-training epochs, since (especially in the last leg of the learning rate schedule, when weights are near their final values) they are trained from approximately the same starting point, using approximately the same learning rate.

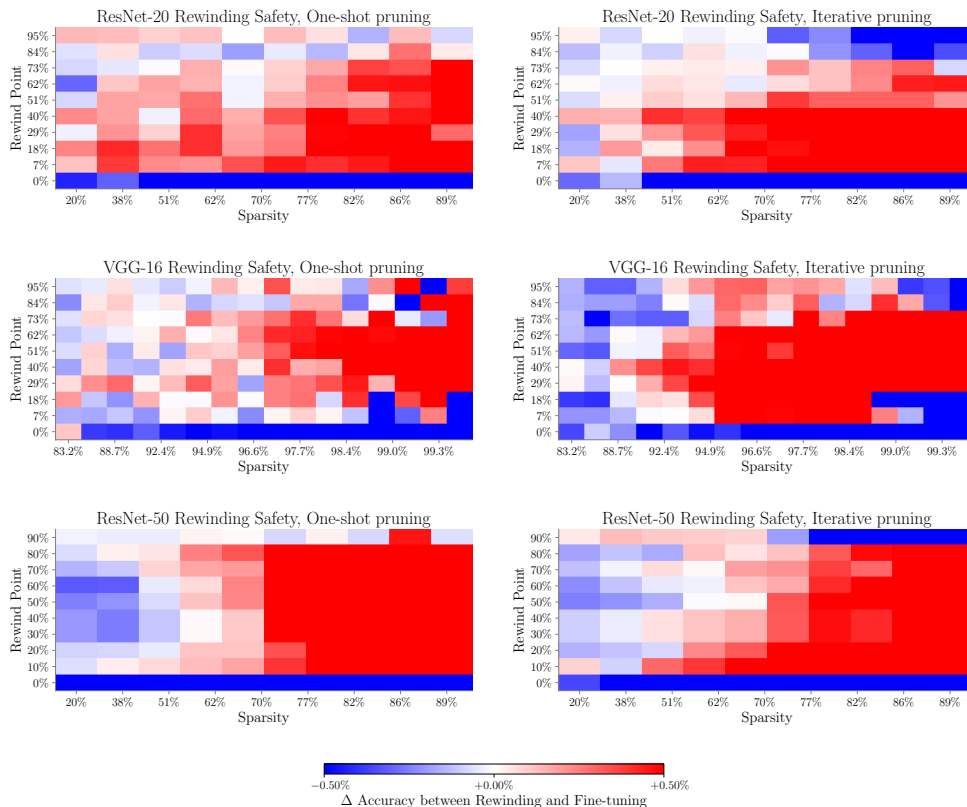


Figure 2: Safety across networks and pruning methods. Red points denote rewinding achieving higher accuracy than fine-tuning for any amount up through the same re-training budget; blue points denote the opposite.

The top region of Figure 3 is sharply blue, indicating that rewinding to the end of training and not re-training achieves significantly lower accuracy than fine-tuning reaches. Magnitude pruning of neural networks damages the function that the neural network learned. Without enough time spent re-training, the network is not capable recovering its original accuracy.

The bottom-most parts of all plots are also sharply blue, which means that rewinding loses its advantage as the weights are rewound to the very beginning of training. This behavior is predicted by Frankle et al. (2019), which demonstrates that networks rewound to the very beginning of training do not typically reach the same accuracy as the unpruned network, whereas networks rewound to a small amount of the way through training do reach that accuracy.

Dense networks. Fine-tuning typically outperforms rewinding in denser networks, as shown by the left side of the plots in Figure 2 and Figures 3 tending towards blue. This is likely due to these networks being trivially prunable, as discussed in Section 3.2. In these situations, fine-tuning is not just a method of recovering the original accuracy of the network: it also behaves like an additional phase of training. In essence, for dense networks fine-tuning can act as a crutch for a non-optimal training schedule, and thereby achieve higher accuracy than rewinding. This advantage falls away at higher sparsities: for deeply pruned networks, there is always a wide range of rewind points such that rewinding achieves higher accuracy than fine-tuning.

Necessary re-training time. These results show a floor on the amount of re-training time for rewinding to be a worthwhile technique: rewinding requires a re-training budget of at least 25% of the original training time to consistently achieve higher accuracy than an equivalent amount of fine-tuning. In situations where the re-training budget is strictly limited to less than 25% of the original training time, fine-tuning outperforms rewinding. However, standard approaches currently

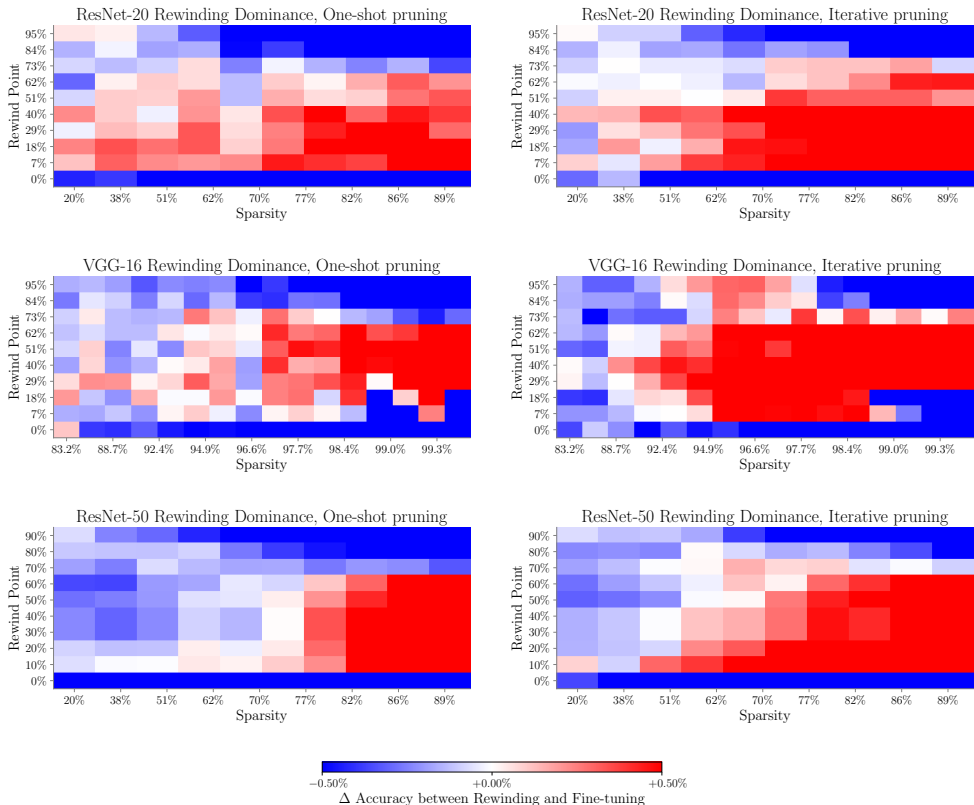


Figure 3: Dominance across networks and pruning methods. Red points denote rewinding achieving higher accuracy than fine-tuning for any amount; blue points denote the opposite.

fine-tune for at least a quarter of the original training time (Li et al., 2017; Liu et al., 2019), and often much longer (Han et al., 2015). Further, 25% of re-training time is often not enough for fine-tuning either: as can be seen by the extra accuracy accumulated by fine-tuning as more re-training time is added (i.e. by the difference between the safety and dominance plots), fine-tuning’s accuracy does not always plateau with 25% re-training time. More visualizations can be found in Appendix B.

5 CONCLUSION

Rewinding expands the frontier in the tradeoff space between the three pruning criteria, ACCURACY, INFERENCE COST, and SEARCH COST. For sparse magnitude pruning on a variety of vision networks, rewinding finds sparser networks than fine-tuning can for any given accuracy threshold. Using this technique, a ResNet-50 can be pruned by 79% with no loss in test accuracy, comparable with state-of-the-art techniques but without using any ResNet-50-specific hyperparameters. Using rewinding requires selecting a point to rewind to ahead-of-time, but we find that there is a wide range of suitable rewinding points that work across all tested networks, at moderate to high sparsities. Rewinding to anywhere between 30% to 75% of the original training time is a *safe* and efficient use of re-training time, achieving higher accuracy than an equivalent amount of fine-tuning. Rewinding to anywhere between 30% to 60% is *dominant*, outperforming any amount of fine-tuning.

Based on these results, we argue that practitioners should explore rewinding as an alternative to fine-tuning for neural network pruning, in order to find sparse and accurate neural networks.

REFERENCES

- Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pp. 265–283, 2016. URL <https://www.usenix.org/system/files/conference/osdi16/osdi16-abadi.pdf>.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rJl-b3RcF7>.
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M. Roy, and Michael Carbin. The lottery ticket hypothesis at scale. *CoRR*, abs/1903.01611, 2019. URL <http://arxiv.org/abs/1903.01611>.
- Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. *CoRR*, abs/1902.09574, 2019. URL <http://arxiv.org/abs/1902.09574>.
- Priya Goyal, Piotr Dollár, Ross B. Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: training imagenet in 1 hour. *CoRR*, abs/1706.02677, 2017. URL <http://arxiv.org/abs/1706.02677>.
- Song Han. *Efficient Methods and Hardware for Deep Learning*. PhD thesis, Stanford University, 2017.
- Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 28*, pp. 1135–1143. Curran Associates, Inc., 2015. URL <http://papers.nips.cc/paper/5784-learning-both-weights-and-connections-for-efficient-neural-network.pdf>.
- Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL <http://arxiv.org/abs/1510.00149>.
- Babak Hassibi, David G. Stork, Gregory Wolff, and Takahiro Watanabe. Optimal brain surgeon: Extensions and performance comparisons. In *Proceedings of the 6th International Conference on Neural Information Processing Systems, NIPS’93*, pp. 263–270, San Francisco, CA, USA, 1993. Morgan Kaufmann Publishers Inc. URL <http://dl.acm.org/citation.cfm?id=2987189.2987223>.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, June 2016. doi: 10.1109/CVPR.2016.90.
- Yihui He, Ji Lin, Zhijian Liu, Hanrui Wang, Li-Jia Li, and Song Han. Amc: Automl for model compression and acceleration on mobile devices. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37, ICML’15*, pp. 448–456. JMLR.org, 2015. URL <http://dl.acm.org/citation.cfm?id=3045118.3045167>.
- A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Master’s thesis, Department of Computer Science, University of Toronto*, 2009.

- Yann Le Cun, John S. Denker, and Sara A. Solla. Optimal brain damage. In *Advances in Neural Information Processing Systems*, pp. 598–605. Morgan Kaufmann, 1990.
- Namhoon Lee, Thalaiyasingam Ajanthan, and Philip Torr. SNIP: SINGLE-SHOT NETWORK PRUNING BASED ON CONNECTION SENSITIVITY. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=B1VZqjAcYX>.
- Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017. URL <https://openreview.net/forum?id=rJqFGTslg>.
- Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network pruning. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rJlnB3C5Ym>.
- Dmitry Molchanov, Arsenii Ashukha, and Dmitry Vetrov. Variational dropout sparsifies deep neural networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17*, pp. 2498–2507. JMLR.org, 2017. URL <http://dl.acm.org/citation.cfm?id=3305890.3305939>.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- Michael Zhu and Suyog Gupta. To prune, or not to prune: exploring the efficacy of pruning for model compression. *arXiv preprint arXiv:1710.01878*, 2017.

SUPPLEMENTARY MATERIAL: COMPARING FINE-TUNING AND REWINDING IN NEURAL NETWORK PRUNING

Anonymous authors

Paper under double-blind review

A PRUNING ALGORITHMS

The following 3 algorithms provide more formal descriptions of the three experiments (rewinding, fine-tuning, and reinitialization) described in Section 2.

Pruning Frameworks

A neural network is a function $f(x; W, mask)$ with weights $W \in \mathcal{R}^d$ and pruning mask $mask \in \{0, 1\}^d$ such that the network parameters are the element-wise product $W \odot mask$.

RANDOM-INITIALIZATION() creates a random initialization of the weights in the network.

TRAIN($W, mask, e, h$) returns the weights resulting from training a masked network $W \odot mask$ for e epochs using the hyperparameters in h (e.g. learning rate, batch size, etc.).

PRUNE($W, mask, h$) removes weights from masked network $W \odot mask$ and returns an updated mask, using the hyperparameters in h (e.g. the pruning heuristic and percentage of weights to prune).

Algorithm 1 Prune with rewinding.

```

1: function REWIND(train_epochs, train_hparams, prune_iters,
                  prune_hparams, retrain_epochs)
2:    $W_0 \leftarrow$  RANDOM-INITIALIZATION()
3:    $mask \leftarrow 1^d$ 
4:    $rewind\_epoch \leftarrow train\_epochs - retrain\_epochs$ 
5:    $W^* \leftarrow$  TRAIN( $W_0, mask, rewind\_epoch, train\_hparams$ )
6:    $W \leftarrow$  TRAIN( $W^*, mask, retrain\_epochs, train\_hparams$ )
7:   for iter in prune_iters do
8:      $mask \leftarrow$  PRUNE( $W, mask, prune\_hparams$ )
9:      $W \leftarrow$  TRAIN( $W^*, mask, retrain\_epochs, train\_hparams$ )
10:  return  $W, mask$ 

```

Algorithm 2 Prune with fine-tuning.

```

1: function FINE-TUNE(train_epochs, train_hparams, prune_iters, prune_hparams,
                    retrain_epochs, finetune_hparams)
2:    $W_0 \leftarrow$  RANDOM-INITIALIZATION()
3:    $mask \leftarrow 1^d$ 
4:    $W \leftarrow$  TRAIN( $mask, train\_epochs, train\_hparams$ )
5:   for iter in prune_iters do
6:      $mask \leftarrow$  PRUNE( $W, mask, prune\_hparams$ )
7:      $W \leftarrow$  TRAIN( $W, mask, finetune\_epochs, finetune\_hparams$ )
8:  return  $W, mask$ 

```

Algorithm 3 Reinitialize.

```

1: function REINITIALIZE(train_epochs, train_hparams, prune_iters,
                        prune_hparams, retrain_epochs)
2:    $W_0 \leftarrow \text{RANDOM-INITIALIZATION}()$ 
3:    $mask \leftarrow 1^d$ 
4:    $W \leftarrow \text{TRAIN}(W_0, mask, train\_epochs, train\_hparams)$ 
5:    $mask \leftarrow \text{PRUNE}(W, mask, prune\_hparams)$ 
6:    $W'_0 \leftarrow \text{RANDOMINITIALIZATION}()$ 
7:    $W \leftarrow \text{TRAIN}(W'_0, mask, train\_epochs + retrain\_epochs \times prune\_iters, train\_hparams)$ 
8:   return  $W, mask$ 

```

B FINE-TUNING TIME

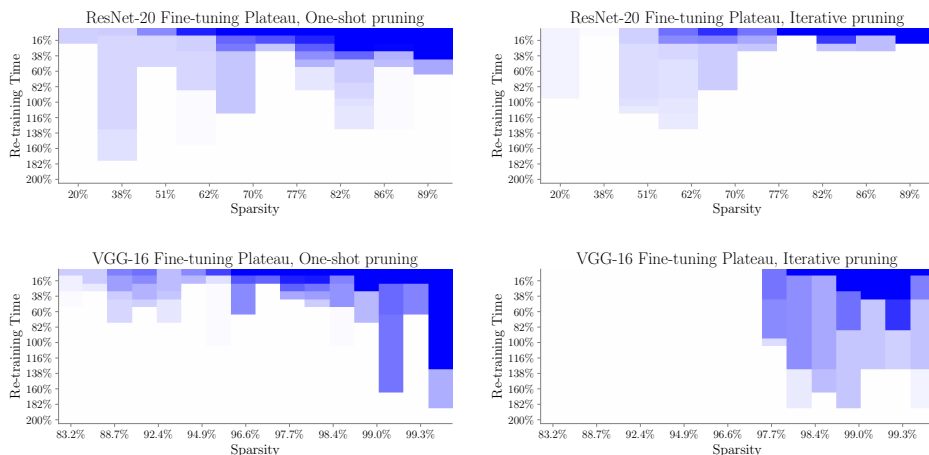
Background. Many standard approaches to pruning only fine-tune for up to 25% of the training time. For instance, Liu et al. (2019) fine-tune CIFAR-10 networks for 40 epochs against an original training schedule of 160 epochs (25%), and fine-tune Imagenet networks for 20 epochs against an original training schedule of 90 epochs (22%). Li et al. (2017) also fine-tune for 40 epochs against the original 160 on CIFAR-10 networks. This is not always the case though: Han et al. (2015) fine-tune for a total time up to $3\times$ the original training time of the network, although this is in the context of iterative pruning.

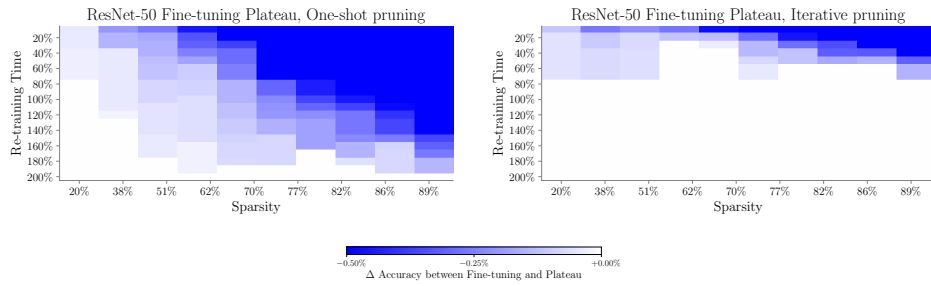
Data. Consider the plots below, which replicate the plots from Section 4, but instead considering fine-tuning’s room for improvement. The y axis is re-training time: points along a given row represent the difference between fine-tuning for at most $y\%$ of the original training time of the work and fine-tuning for up to $2\times$ the original training time for the network. White points denote points where fine-tuning has plateaued; blue points denote points where more fine-tuning can help.

Analysis We find that in many cases, a quarter of the original re-training time is not enough for fine-tuning’s accuracy to plateau. In ResNet-50 one-shot pruning, training does not seem to plateau

We also find that the required re-training time for accuracy to plateau is sparsity-dependent. This is unsurprising, as more pruning brings the network further away from the function it initially learned. However, standard approaches in the literature do not acknowledge this fact, and instead use a fixed fine-tuning schedule, rather than a sparsity-dependent one or a schedule that stops training upon accuracy plateau.

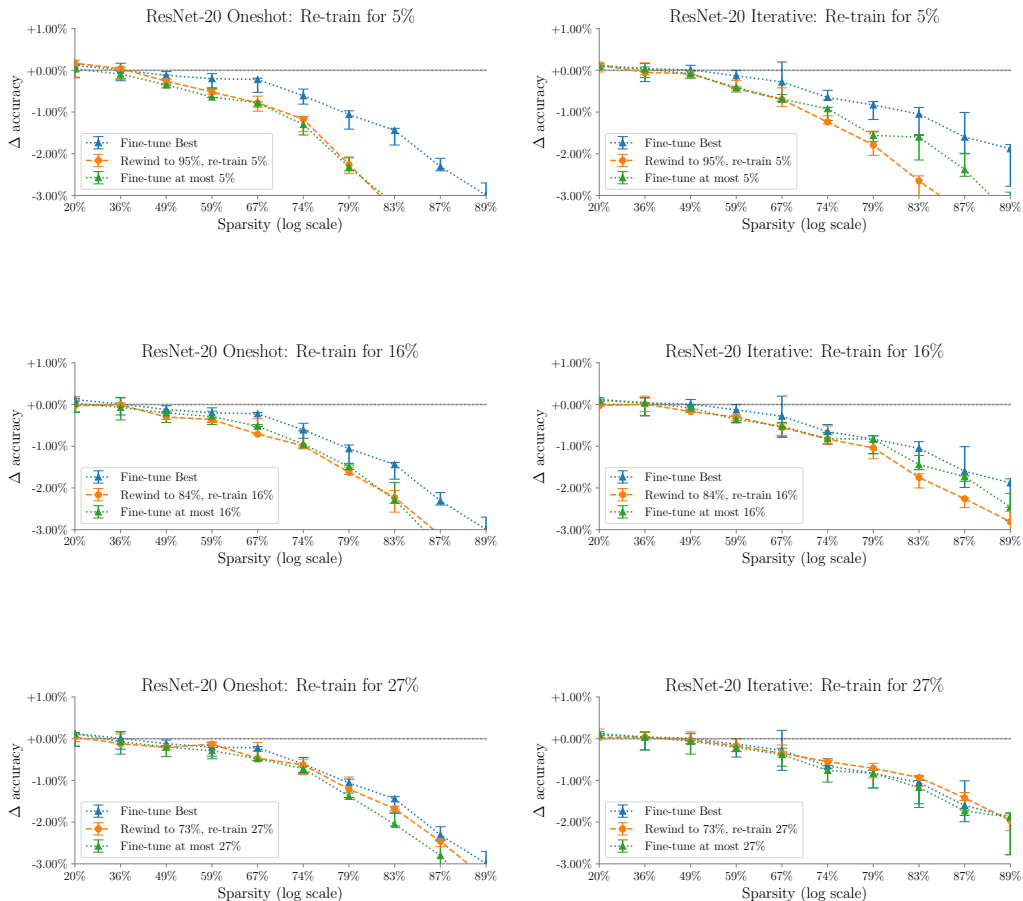
More views and analysis of the same underlying data can be found in Appendix C and Appendix D.

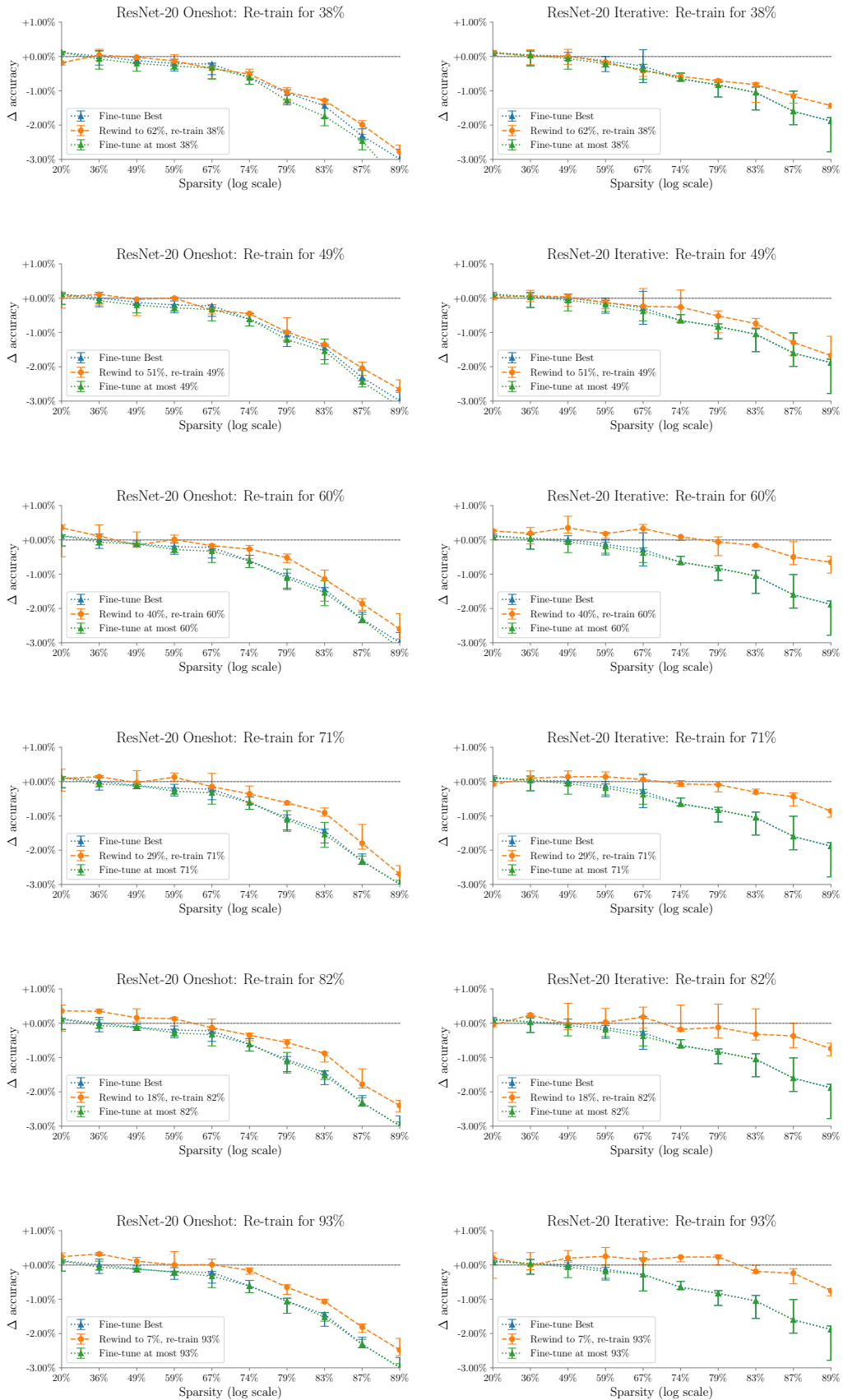


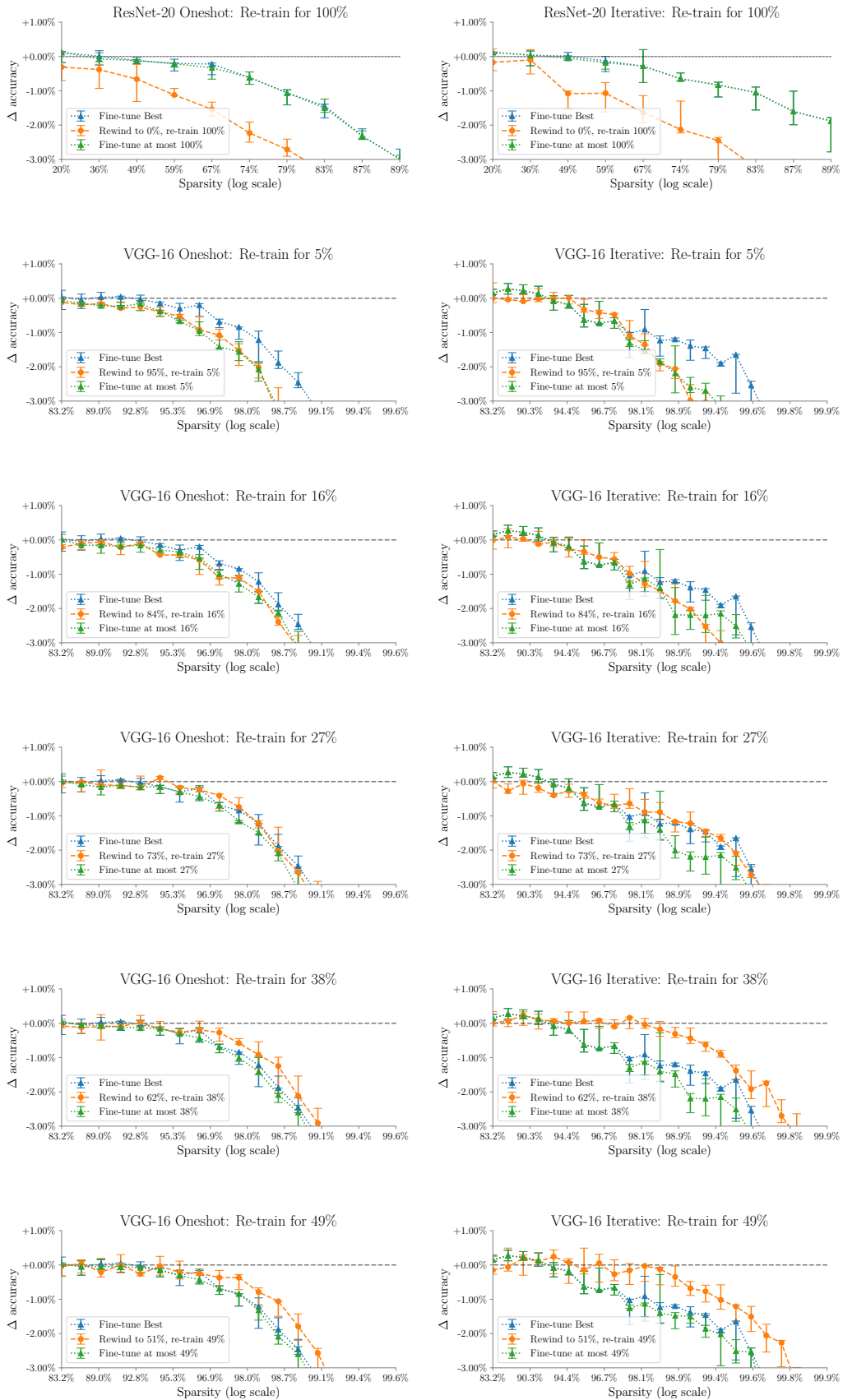


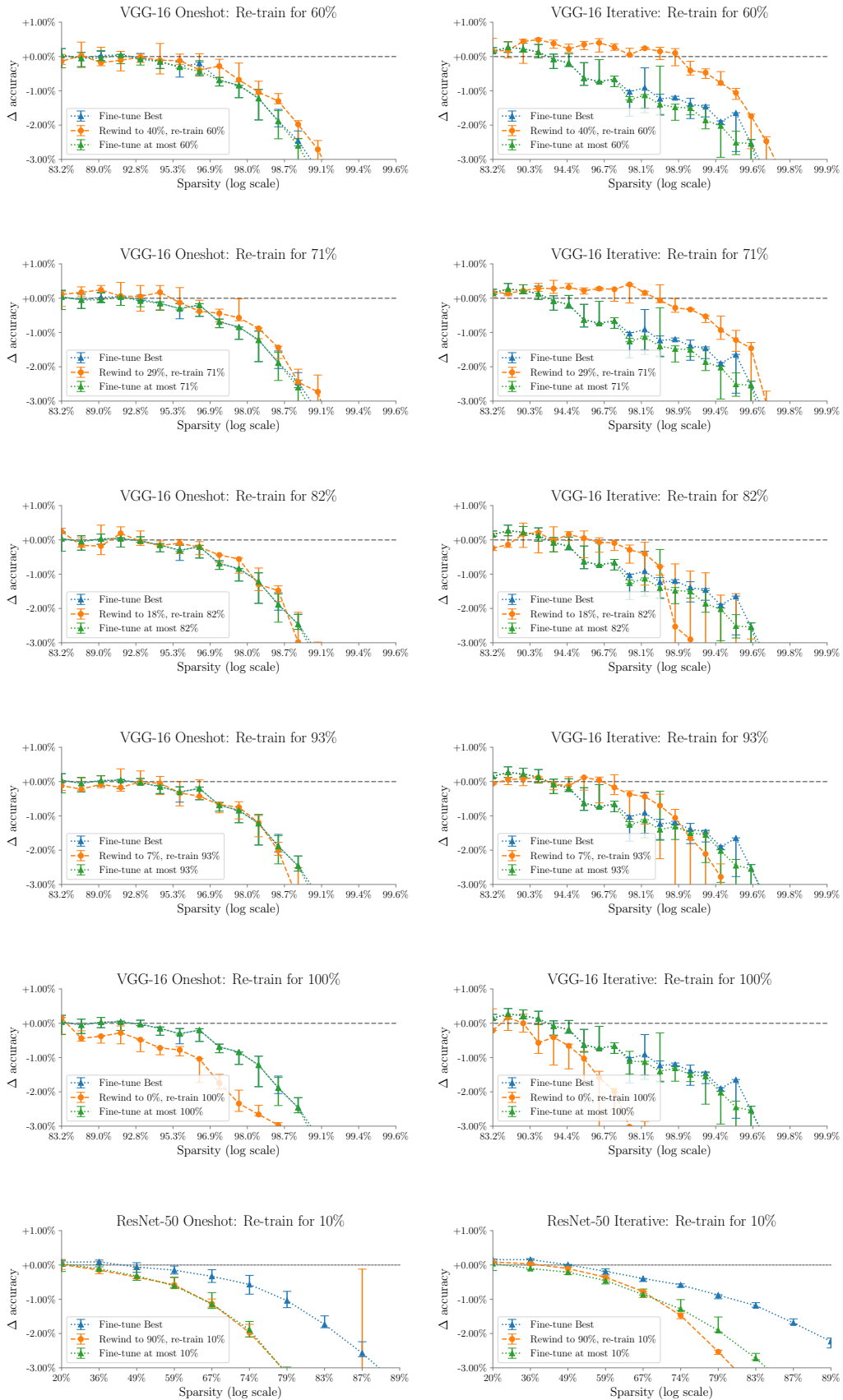
C BUDGET-LIMITED ACCURACIES

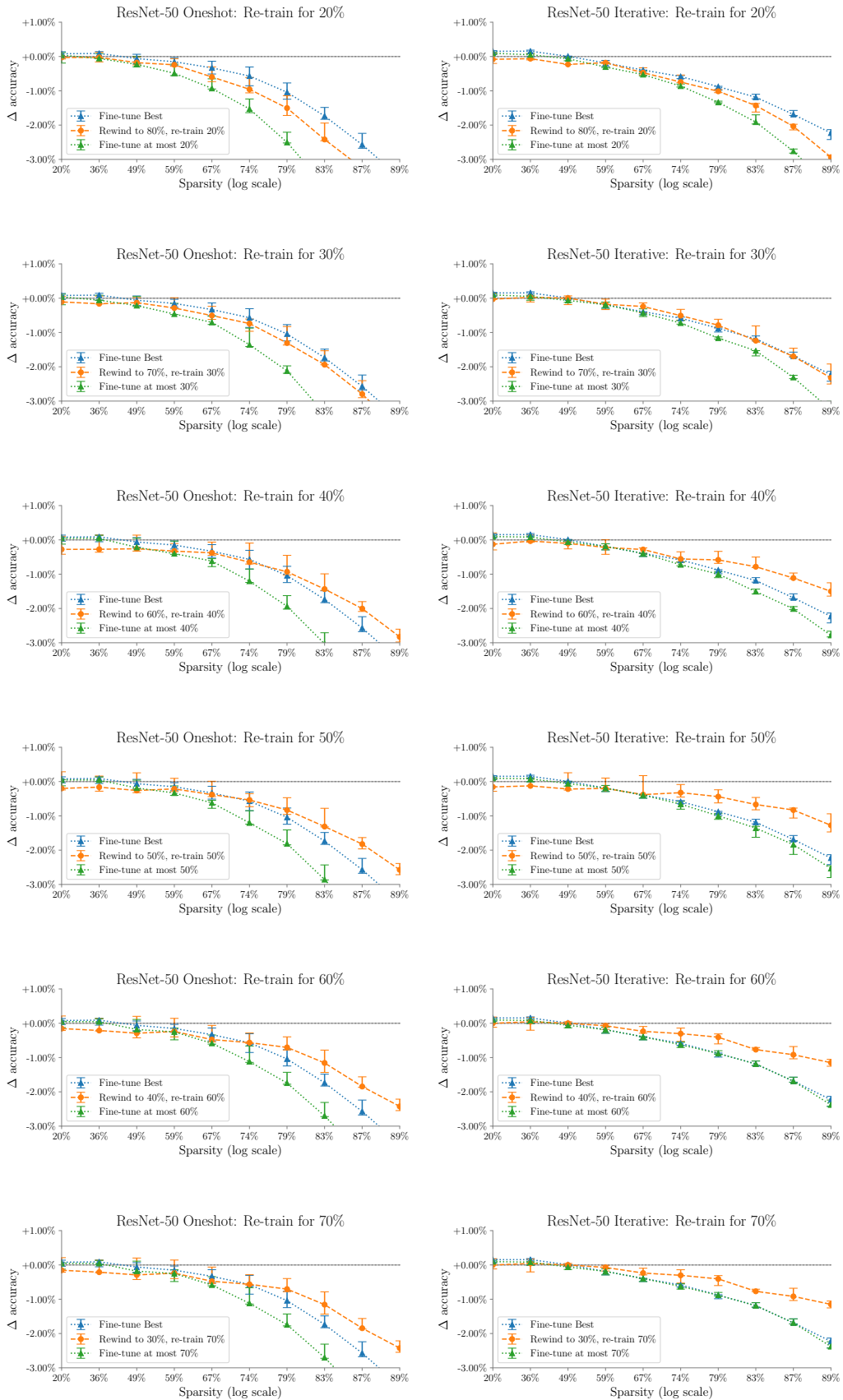
In this appendix, we consider the accuracy attained by rewinding to specific point in training, fine-tuning for any amount up to and including the number of epochs for which that rewind network is re-trained, and fine-tuning for any amount at all. These plots can be considered versions of the plots in Figure 1, but with rewinding to specific points and with budget-constrained fine-tuning. The intention of these plots is to show the absolute accuracy attainable by picking specific rewind points, and to provide another view of the safe and dominant regions as described in Section 4. These plots can be seen in the plots below.

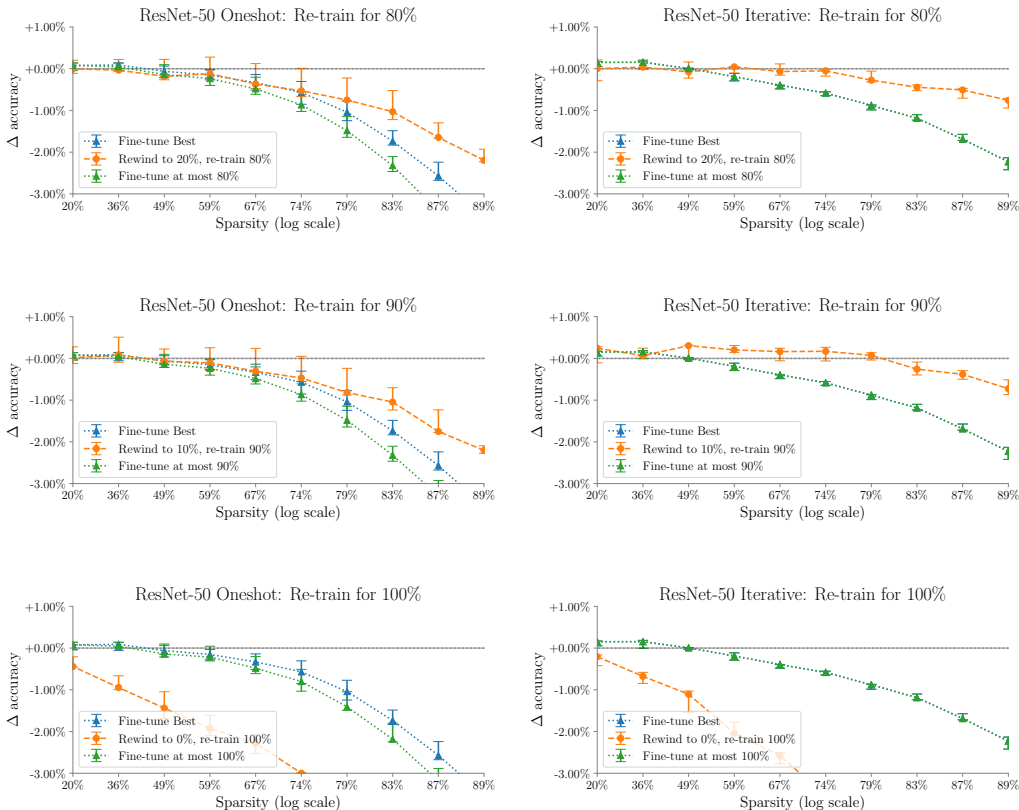












D EQUAL BUDGET RE-TRAINING

In this appendix, we compare rewinding and fine-tuning using the same number of re-training epochs for each, showing that epoch-for-epoch, rewinding outperforms fine-tuning in most situations, with the exception of iterative re-training for a small number of re-training epochs. This is similar to the concept of safety presented in Section 4, but with an exactly equal budget rather than comparing rewinding with the best-performing set of fine-tuning trials. We include these data for two reasons. First, it is informative to know how many epochs of fine-tuning are necessary for different networks and sparsity levels. Also, we include these data to show a more fine-grained comparison between rewinding and fine-tuning: although fine tuning can stop early, this is both not entirely true for iterative fine-tuning, and also is not representative of the real accuracy of fine-tuning since a max across trials is not an unbiased estimator (e.g. consider the case when learning from fine-tuning plateaus early: since each point is an independent set of 3 trials, the max of all medians picks the set of trials that happened to perform the best, rather than the true median performance).

The plots below compare fine-tuning and rewinding across networks and sparsity levels. In these plots, the y -axis is the change in accuracy between the pruned network and the original network (higher is better), and the x -axis is the number of epochs that the pruned network was re-trained for according to Algorithms 1 and 2 (in iterative pruning, the x axis also reflects the iterative nature of iterative pruning – e.g., after 5 iterations of pruning with fine-tuning for 80 epochs, the total number of re-training epochs is 80×5). For fine-tuning, the x -axis denotes the number of fine-tuning epochs. For rewinding, the x -axis denotes the epoch $T - x$ that the network was rewound to (where T is the original training time of the full network) and the subsequent x epochs of re-training.

Iterative fine-tuning can overfit One interesting feature of the below graphs is that fine-tuning seems to overfit when ran for too long. Rewinding and replaying seems to be immune to overfitting from gradient descent: rather than accumulating gradients from hundreds or thousands of training epochs, networks pruned with the rewinding/replaying framework only receive gradient updates from a single full training run. While for one-shot pruning this is an easy problem to fix (simply

stopping early, or using a high accuracy checkpoint), this is a thornier problem for iterative pruning, since it is not possible to undo extra training iterations from previous iterations in the iterative pruning process by simply restoring from a checkpoint.

