EXAMINING THE VALUE OF NEURAL FILTER PRUNING – RETROSPECT AND PROSPECT

Anonymous authors

Paper under double-blind review

Abstract

Neural network filter pruning is one of the major methods in model compression and acceleration. Despite the remarkable progress in the past several years, there is an ongoing debate concerning the value of filter pruning – Some works in 2019 argue that filter pruning is of no value since they found training the pruned network from scratch can achieve similar or even better performance than pruning a pretrained model. This argument fundamentally challenges the value of many filter pruning works. However, to date, the community has not formally responded to such acute questioning. In this paper, we present extensive empirical analyses to show the seeming contradiction is due to suboptimal learning rate schedule settings. We introduce more strict comparison setups and show filter pruning still has value within the same training epoch budgets. Apart from justifying the value of filter pruning empirically, we further examine the reason behind it and discover that the poor trainability caused by pruning is largely responsible for the sub-optimality of the learning rate schedule, thus calling for an urgent need to recover trainability after pruning. This paper does not target new SOTA performance of filter pruning. Instead, we focus on clarifying the existing mysteries in filter pruning towards a better understanding.

1 INTRODUCTION

Pruning is a time-honored methodology to reduce parameters in a neural network without seriously compromising its performance (Reed, 1993; Sze et al., 2017). The prevailing pipeline of pruning comprises three steps: 1) **pretraining**: train a dense model; 2) **pruning**: prune the dense model based on certain rules; 3) **finetuning**: retrain the pruned model to regain performance. Most existing research focuses on the second step, seeking the best criterion to select unimportant weights so as to incur as less performance degradation as possible. This 3-step pipeline has been practiced for more than 30 years (Mozer & Smolensky, 1989; LeCun et al., 1990) and is still extensively adopted in today's pruning methods (Sze et al., 2017).

These said, several works (Crowley et al., 2018; Liu et al., 2019) questioned the necessity of inheriting weights from a pretrained model because they empirically found the small model trained from scratch can match (or sometimes outperform) the counterpart pruned from the pre-trained large model. This acutely challenges the past wisdom as well as our common belief about pruning. As far as we know, there is *no* formal response to this critical conflict. A theoretical-level understanding of this problem is even more elusive.

Meanwhile, the pruning community has been observing even more open questions. Specifically, Renda et al. (2020); Le & Hua (2021) found that the learning rate (LR) in finetuning holds a critical role in the final performance. A proper learning rate schedule (*e.g.*, a larger initial LR 10^{-2} vs. 10^{-3} with step-decay schedule) can improve the top-1 accuracy of a pruned ResNet-34 model (He et al., 2016) by more than 1% on ImageNet (Deng et al., 2009). This discovery calls for more attention being paid to the finetuning step when comparing different pruning methods. Unfortunately, they did not present more theoretical insights to explain its occurrence. This also remains an open question in the community up to date.

In this paper, we will show these two open questions actually point to the same one. Specifically, we rerun the experiments of (Crowley et al., 2018; Liu et al., 2019) and find simply using a larger fine-tuning LR (10^{-2} vs. 10^{-3} and decay it) can *significantly* improve the final performance. Compared

Table 1: Top-1 accuracy comparison of different implementations of the L_1 -norm pruning (Li et al.,
2017) on ImageNet. We adopt the torchvision models as unpruned models for fair comparison.
ResNet-34-A speedup: 1.18×. ResNet-34-B speedup: 1.32×. The results of (Li et al., 2017) and
(Liu et al., 2019) are directly cited from their papers. The best cases in our training from scratch and
pruning are randomly repeated for 3 times (\pm indicates stddev) to prevent random variation.

Implementation	Unpruned (%)	Pruned model	Scratch (%)	Pruned-Finetuned (%)	Finetuning LR schedule
(I : at al 2017)	73.23	ResNet-34-A	(Not reported)	72.56	20 epochs, initial 10^{-3} , fixed
(L1 et al., 2017)		ResNet-34-B	(Not reported)	72.17	20 epochs, initial 10^{-3} , fixed
(Lin et al. 2010)	72.21	ResNet-34-A	73.03	72.56	20 epochs, initial 10^{-3} , fixed
(Liu et al., 2019)	75.51	ResNet-34-B	72.91	72.29	20 epochs, initial 10^{-3} , fixed
	79.91	ResNet-34-A		72.91	20 epochs, initial 10^{-3} , fixed
Our rerun			$73.51_{\pm 0.12}$	72.94	90 epochs, initial 10^{-3} , fixed
Our iciuii	10.01			73.88	90 epochs, initial 10^{-3} , decay
				$73.92_{\pm 0.03}$	90 epochs, initial 10^{-2} , decay
				72.50	20 epochs, initial 10^{-3} , fixed
Our rerun	72.21	PacNat 3/ B	73 16	72.58	90 epochs, initial 10^{-3} , fixed
Our lerun	10.01	Keshel-34-D	13.10 ± 0.12	73.61	90 epochs, initial 10^{-3} , decay
				$73.62_{\pm 0.04}$	90 epochs, initial 10^{-2} , decay

to the improved pruning performance, training from scratch does *not* compete or surpass pruning anymore (see Tab. 1 and Tab. 2 on ImageNet). Why does this happen? What is the theoretical reason behind this change?

This paper is meant to present more precise answers to these questions. The key perspective shift we take is that we consider the weights inherited from pruning as a kind of initialization for the following finetuning. Then a natural question comes: *Does pruning provide a "good" initialization for the subsequent finetuning?* We will show the study of this question help unveil the above two open problems.

Specifically, we tap into the notion of *dynamical isometry* (Saxe et al., 2014) (which describes a kind of nice property in neural networks that are easy to optimize) as a faithful analysis tool. We carefully design an explanatory experiment using a linear MLP (multi-layer perceptron) network to demonstrate how finetuning LR affects the final performance by affecting dynamical isometry. In brief, we observe the finetuning process can recover dynamical isometry inherently; a larger LR can help recover it faster (or better), hence the better final performance. The proposed explanation is validated by our empirical results and resonates with many empirical observations.

Contributions. (1) We empirically demonstrate the questioning about the value of inheriting weights in structured pruning in previous works is inaccurate and point out that the direct cause is improperly using a small finetuning LR. Our finding justifies the value of inheriting weights in structured pruning. (2) On top of the empirical finding, more importantly, we present a theoretical explanation through examining the dynamical isometry of networks in pruning. This explanation is empirically validated by our carefully designed control experiments. (3) Our investigation promotes a perspective shift for pruning: looking at pruning as a kind of initialization for the subsequent finetuning can help us better understand pruning.

2 RELATED WORK

Conventional pruning. Pruning aims to remove as many parameters as possible in a neural network meanwhile maintaining its performance. There are many ways to categorize pruning methods. The most popular two are grouping by pruning structure and methodology.

(1) In terms of pruning structure, pruning can be specified into unstructured pruning (Han et al., 2015; 2016) and structured pruning (Wen et al., 2016; Li et al., 2017; He et al., 2017). For the former, a single weight is the basic pruning element. Unstructured pruning can deliver a high compression ratio; whereas, without regularization, the pruned locations usually spread *randomly* in the network, which is hard to exploit for acceleration. On the opposite, structured pruning introduces certain patterns in the pruned locations, which benefit subsequent acceleration while cannot achieve as much compression. Choices between unstructured and structured pruning depend on specific application needs. For structured pruning, there are still many sub-groups (Mao et al., 2017). In the literature,

without specific mention, structured pruning means filter pruning or channel pruning. This paper focuses on **structured (filter) pruning** because the "no value of inheriting weights" argument is mainly discussed in this context (Liu et al., 2019).

(2) In terms of pruning methodology (*i.e.*, how to select unimportant weights to prune), pruning falls into two paradigms in general: importance-based and penalty-based. The former prunes weights based on some established importance criteria, such as magnitude (for unstructured pruning) (Han et al., 2015; 2016) or L_1 -norm (for filter pruning) (Li et al., 2017), saliency based on 2nd-order gradients (*e.g.*, Hessian or Fisher) (LeCun et al., 1990; Hassibi & Stork, 1993; Theis et al., 2018; Wang et al., 2019a; Singh & Alistarh, 2020). The latter adds a penalty term to the objective function, drives unimportant weights towards zero, then removes those with the smallest magnitude. Note, the two groups are *not* starkly separated. Many methods take wisdom from both sides. For example, (Ding et al., 2018; Wang et al., 2019b; 2021) select unimportant weights by magnitude (akin to the first group) while also employing the regularization to penalize weights (akin to the second group). There is no conclusion about which paradigm is better, yet empirically, the state-of-the-art pruning methods are closer to the second paradigm, *i.e.*, deciding weights via training instead of some derived formulas. Although no theories have formally discussed the reason, we can take a rough guess with the knowledge from this paper: Training can recover dynamical isometry, which is beneficial to subsequent finetuning.

For more comprehensive literature, we refer interested readers to several surveys: an outdated one (Reed, 1993), some recent surveys of pruning alone (Gale et al., 2019; Blalock et al., 2020) or pruning as a sub-topic under the general umbrella of model compression and acceleration (Sze et al., 2017; Cheng et al., 2018a;b; Deng et al., 2020).

Pruning at initialization (PaI). Recent years have seen several new pruning paradigms. The most prominent one is pruning at initialization. Different from the conventional pruning, which prunes a *pretrained* model, PaI methods prune a *randomly initialized* model. Existing PaI approaches mainly include (Lee et al., 2019; 2020; Wang et al., 2020; Frankle et al., 2021; Ramanujan et al., 2020) and the series of lottery ticket hypothesis (Frankle & Carbin, 2019; Frankle et al., 2020). Interested readers may refer to (Wang et al., 2022) for a comprehensive summary about PaI.

This topic is relevant to this work mainly because one PaI paper (Lee et al., 2020) also examines pruning using the tool of dynamical isometry. The similarity between our paper and theirs is that we both employ dynamical isometry as a tool to examine the property of network pruning. However, our paper is significantly different from theirs in many axes: (1) Basic setting. The most obvious difference is that we focus on pruning a *pretrained* model while (Lee et al., 2020) focuses on pruning at initialization (PaI). They are two different tracks in pruning (as such, PaI methods typically do not compare with the methods of pruning pretrained models) and the latter was shown to consistently underperform the former (Frankle et al., 2021; Wang et al., 2022). (2) Motivation. Despite the same tool (mean JSV in Eq. (1)), (Lee et al., 2020) uses it to select unimportant weights to prune (*i.e.*, for a new pruning criterion), while we use it to analyze why finetuning LR has a significant impact on final performance. The role of finetuning LR in pruning is not mentioned at all in their paper. (3) Proposed technical method. (Lee et al., 2020) focuses on unstructured pruning, while we focuses on *structured pruning*. This further leads to fundamental difference when designing the dynamical isometry recovery (DIR) methods – In (Lee et al., 2020), their proposed method is to use iterative optimization for approximated isometry (due to the irregular sparsity); while in our case, since the pruned filers can be completely removed from the network, one of our DIR method (OrthP) has closed-form solution and can achieve exact isometry.

2.1 EMPIRICAL STUDY: LARGER FINETUNING LR IS CRITICAL

As far as we know, mainly *two* papers question the value of inheriting weights from a pretrained model: (Crowley et al., 2018; Liu et al., 2019). Both papers draw two similar conclusions. (1) Inheriting weights from a pretrained model in pruning has *no* value, *i.e.*, training from scratch the small model can match (or outperform sometimes) the counterpart pruned from a big pretrained model. (2) Given the fact of (1), what really matters in pruning may lie in the pruned *architecture* instead of the inherited weight values. As such, both papers propose to view pruning as a form of neural architecture search (Zoph & Le, 2017; Elsken et al., 2019). In this section, we first reexamine

,			,		U	~
Network	PR	Params reduc. (%)	FLOPs reduc. (%)	Scratch (%)	Pruned-Finetuned-1 (%)	Pruned-Fintuned-2 (%)
	0	0	0	69.76^{\dagger}	/	/
	0.1	9.56	9.58	$70.15_{\pm 0.02}$	70.43 ± 0.02	$70.48_{\pm 0.10}$
	0.3	28.32	28.18	$68.90_{\pm 0.10}$	$\overline{69.29_{\pm 0.07}}$	$69.54_{\pm 0.09}$
ResNet-18	0.5	47.03	46.20	$67.03_{\pm 0.01}$	$\overline{67.36_{\pm 0.03}}$	$67.71_{\pm 0.05}$
	0.7	65.99	64.93	$64.21_{\pm 0.10}$	$\overline{63.72_{\pm 0.03}}$	$64.45_{\pm 0.04}$
	0.9	84.75	83.52	$56.70_{\pm 0.17}$	$53.49_{\pm 0.10}$	$55.89_{\pm 0.11}$
	0.95	89.51	88.03	$51.83_{\pm 0.14}$	$44.46_{\pm 0.15}$	49.99 ± 0.10
	0	0	0	73.31 [†]	/	/
	0.1	9.84	9.92	73.53 ± 0.10	$73.86_{\pm 0.05}$	$74.04_{\pm 0.05}$
	0.3	29.15	29.26	$72.50_{\pm 0.17}$	$73.11_{\pm 0.06}$	$73.31_{\pm 0.08}$
ResNet-34	0.5	48.41	48.12	$71.27_{\pm 0.03}$	$\overline{71.71_{\pm 0.06}}$	$71.85_{\pm 0.07}$
	0.7	67.95	67.63	$68.69_{\pm 0.10}$	$\overline{68.90_{\pm 0.08}}$	$69.33_{\pm 0.04}$
	0.9	87.26	86.97	62.08 ± 0.12	$\overline{60.34_{\pm 0.03}}$	$62.26_{\pm 0.06}$
	0.95	92.16	91.69	$57.21_{\pm 0.15}$	$52.83_{\pm 0.11}$	56.69 ± 0.23
	0	0	0	70.37^{\dagger}	/	/
	0.1	9.19	17.78	$68.51_{\pm 0.04}$	$71.45_{\pm 0.07}$	$71.74_{\pm 0.04}$
	0.3	26.80	47.63	$66.60_{\pm 0.11}$	$\overline{70.00_{\pm 0.03}}$	$70.53_{\pm 0.05}$
VGG11_BN	0.5	43.86	70.56	$65.85_{\pm 0.13}$	$\overline{67.34_{\pm 0.05}}$	$68.01_{\pm 0.10}$
	0.7	60.54	87.03	61.56 ± 0.06	$\overline{62.14_{\pm 0.09}}$	$63.14_{\pm 0.08}$
	0.9	76.50	96.49	$48.34_{\pm 0.12}$	$\overline{45.11_{\pm 0.07}}$	$48.52_{\pm 0.06}$
	0.95	80.49	97.76	$35.47_{\pm 0.09}$	$33.50_{\pm 0.10}$	$38.47_{\pm 0.13}$

Table 2: Top-1 accuracy comparison between scratch training ("Scratch") and L_1 -norm pruning (Li et al., 2017) on ImageNet. "PR" means pruning ratio. [†]We adopt the official torchvision models as unpruned models. "Finetuned-1" and "Finetuned-2" refers to two finetuning LR schedules ("Finetuned-1": 90 epochs, initial 10^{-3} , decay 30/60/75; "Finetuned-2": 90 epochs, initial 10^{-2} , decay 45/68). Best results are in **bold**, second best underlined, each averaged by 3 random runs.

the empirical studies in (Crowley et al., 2018; Liu et al., 2019) to show that the "no value of inheriting weights" argument is actually inaccurate owing to improper finetuning LR schedules.

Reexamination of (Liu et al., 2019). Before presenting results, here are some important comparison setting changes worth particular attention.

(1) In (Liu et al., 2019), they compare training from scratch with *six* pruning methods (five structured pruning methods (Li et al., 2017; Luo et al., 2017; Liu et al., 2017; He et al., 2017; Huang & Wang, 2018) and one unstructured pruning method (Han et al., 2015)). Here, we only focus on the L_1 -norm pruning (Li et al., 2017) on ImageNet. The main reason is that, L_1 -norm pruning is well-known a very basic filter pruning method. If we can show it outperforms training from scratch already, it will be no surprise to see other more advanced pruning methods also outperform training from scratch. In this sense, L_1 -norm pruning is the most representative method here for our investigation.

(2) In (Liu et al., 2019), they have two variants for the number of epochs in scratch training, "Scratch-E" and "Scratch-B". For the former, different small models are trained for a fixed number of epochs; for the latter, *smaller* models are trained for *more* epochs to maintain the same computation budget (Scratch-B was shown to be better than Scratch-E in (Liu et al., 2019)). Also, they decay LR only to 10^{-3} following the official PyTorch ImageNet example¹. Here, we simply train all the networks for the same number of epochs but ensure the epochs are abundant (120 epochs) and decay LR to a very small amount (10^{-5}). These two changes are to make sure the networks are trained to *full convergence*. As we will show, one primary cause possibly leading (Liu et al., 2019) to an inaccurate conclusion is exactly that the pruned networks are *not* fully converged (see Tab. 1).

With the LR schedule changes, we rerun the experiments using the released code of (Liu et al., 2019). Results are presented in Tab. 1. In the implementations of (Liu et al., 2019), the finetuned model is outperformed by the scratch training one, hence their "no value of inheriting weights" argument. We also reproduce their settings (the two rows of "20 epochs, initial 10^{-3} , fixed" in "Our rerun") for confirming their argument. However, the finetuning LR schedule "20 epochs, initial 10^{-3} , fixed" is actually sub-optimal; the network is *not* fully converged. Using the proper ones ("90 epochs, initial 10^{-3} , decay" or "90 epochs, initial 10^{-2} , decay"), pruning *outperforms* training from scratch for both ResNet-34-A and ResNet-34-B. We note the pruned models even outperform the

¹https://github.com/pytorch/examples/tree/master/imagenet

Table 3: Comparison between L_1 -norm pruning and training from scratch with **ResNet34** on **ImageNet100** at Fair-setup-1 and Fair-setup-2. Each result is averaged by 3 random runs. The learning rate schedule of scratch training is: Initial LR 0.1, decayed at epoch 30/60/90/105 by multiplier 0.1, total: 120 epochs. "Prune30, Ft90, 1e-1" means the model is pruned at epoch 30 and finetuned for another 90 epochs with initial finetune LR 1e-1. The other can be inferred likewise.

Method	0.1	0.3	0.5	0.7	0.9	0.95
Scratch	$83.68_{\pm0.38}$	$83.31_{\pm0.13}$	$82.90_{\pm0.16}$	$82.45_{\pm 0.13}$	$79.37_{\pm0.76}$	$76.67_{\pm 0.90}$
L ₁ (Fair-setup-2: Prune30, Ft90, 1e-2)	84.13 ± 0.11	83.44 ± 0.03	83.51 ± 0.23	83.48 ± 0.19	80.03 ± 0.14	77.00 ± 0.34
L ₁ (Fair-setup-1: Prune30, Ft90, 1e-1)	$\textbf{85.61}_{\pm 0.18}$	$\textbf{85.82}_{\pm 0.30}$	$\textbf{85.54}_{\pm 0.40}$	$\textbf{84.61}_{\pm 0.41}$	$\textbf{82.07}_{\pm 0.11}$	$78.40_{\pm 0.23}$

original models. This is probably because pruning reduces the network redundancy, thus curbing overfitting. This phenomenon is also widely observed in past works (Han et al., 2016; Wen et al., 2016; He et al., 2017) especially with small pruning ratios as the cases in Tab. 1.

Tab. 1 only presents two ResNet models and their speedups are actually quite small. To see if the finetuning LR effect still holds across the full spectrum of pruning ratios and on other types of networks, we vary the pruning ratios from 0.1 to 0.95 and include experiments on VGG11_BN (Simonyan & Zisserman, 2015).

Results are presented in Tab. 2. With a more proper finetune LR scheme (column "Pruned-Fintuned-2" vs. "Pruned-Fintuned-1"), the performance can be improved *significantly*. A clear pattern is, the larger the pruning ratio, the more of the improvement. Now, comparing the results of "Pruned-Fintuned-2" to those of "Scratch", we can see pruning *outperforms* scratch-training in most cases. Exceptions appear on ResNet-34/18 under extreme pruning ratios (90% and 95%). Despite them, we believe it is fair to say **inheriting weights** *has* **value given the fact that 17/20 experiments in Tabs. 1 and 2 show pruning is better than training from scratch**, especially under the pruning ratios of practical interest (*i.e.*, non-extreme pruning ratios). Retrospectively, (Liu et al., 2019) concluded oppositely because they re-implemented the L_1 -norm pruning method *exactly following the description in the original paper* (Li et al., 2017): fixed LR 10⁻³, 20 epochs, which turns out *far from optimal* as we know now.

Stricter comparison settings. In Tab. 1 and Tab. 2, although the two schemes (pruning-finetuning vs. scratch) are trained to their full convergence (which is the most common comparison setup in the pruning literature), they are *not strictly fair* if we consider the total epochs including pretraining and finetuning – pruning-finetuning takes 90+90=180 epochs, while training from scratch only takes 120 epochs. It is of deep interest how the comparison would turn *if we keep the total epochs exactly the same*. Here we present more results to address this. Two stricter comparison setups are introduced:

- Fair-setup-1: The total epochs are kept the same while the specific LR schedules are allowed to be different.
- Fair-setup-2: The total epochs and LR schedules are both kept the same.

Clearly, Fair-setup-2 is even stricter than Fair-setup-1. Then at these two setups, we reconduct the experiments under different pruning ratios (0.1 to 0.95) on *ImageNet100* (a randomly drawn 100-class subset of ImageNet-1K to obtain results more quickly). The results are presented in Tab. 3. Even at these stricter setups, L_1 -norm pruning still *outperforms* scratch training consistently.

Reexamination of (Crowley et al., 2018). Coincidentally, (Crowley et al., 2018) adopted a *very similar* finetuning LR scheme to (Liu et al., 2019): They finetuned the pruned network with the lowest LR (8×10^{-4} , close to 10^{-3} in (Liu et al., 2019)) during scratch training and also *fixed*. Like the empirical study above, we reproduce the experiments of (Crowley et al., 2018) and rerun them with a larger initial LR (10^{-2}) and decay it during finetuning.

Detailed results are deferred to the Appendix (Tab. 8) due to the limited length here. We summarize the observation here – Exactly the same as the case in (Liu et al., 2019), when the proper finetuning LR is used, pruning actually *outperforms* the best scratch training scheme consistently.

Retrospective remarks. Simply put, results above show that the debate about the value of pruning is largely attributed to sup-optimal finetuning settings. It is worthwhile to ponder at this point why this simple reason was not spotted for years. In fact, the problem is not so straightforward to see,

Table 4: Mean JSV of the first 10 epochs under different finetuning settings. Epoch 0 refers to the model just pruned, before any finetuning. Pruning ratio is 0.9. Note, with OrthP, the mean JSV is 1 because OrthP can achieve exact isometry.

Epoch	0	1	2	3	4	5	6	7	8	9	10
$LR=10^{-2}$, w/o OrthP	0.0004	0.6557	0.8946	1.0191	0.9826	1.0965	1.1253	1.2595	1.3298	1.2940	1.4238
$LR=10^{-6}$, w/o OrthP $LR=10^{-2}$, w/ OrthP	1.00004	1.2318	1.4144	1.4277	1.4017	0.2765	1.5171	1.5551	1.6082	1.6538	1.6648
LR= 10^{-3} , w/ OrthP	1.0000	1.5135	1.6630	1.7449	1.8250	1.8720	1.9193	1.9556	1.9943	2.0084	2.0409

Table 5: Test accuracy (%) of the first 10 epochs *corresponding to Tab.* **4** under different finetuning settings. Epoch 0 refers to the model just pruned, before any finetuning. Pruning ratio is 0.9.

- <u></u>			Jaar	r,			2				•
Epoch	0	1	2	3	4	5	6	7	8	9	10
LR= 10^{-2} , w/o OrthP	9.74	63.86	79.96	79.74	80.06	85.79	85.82	86.11	86.45	86.53	85.95
LR=10 ⁻³ , w/o OrthP	9.74	9.74	9.74	12.09	21.74	27.95	33.55	35.92	49.19	65.50	69.90
$LR=10^{-2}$, w/ OrthP	9.74	91.05	91.39	91.33	91.37	91.74	91.69	90.74	91.39	91.58	91.44
LR= 10^{-3} , w/ OrthP	9.74	90.81	91.59	91.77	91.85	92.04	92.12	92.22	92.12	92.33	92.25

because it has been broadly believed that inherited weights of a pruned model already pose it at a location close to the final solution in the loss landscape, hence no need to finetune the model for *many* epochs (typically much fewer than the number of epochs to train the model from scratch). This is probably why (Li et al., 2017) finetuned the model for merely 20 epochs on ImageNet, far from the best setting. The reality, however, turns out to be the opposite way: The pruned models actually demand *more* finetuning epochs because pruning hurts the dynamical isometry, slowing down the convergence. This gap between our presumption and the reality, not noticed by previous works, finally leads to the debate about the value of pruning. Now, after discovering the role of dynamical isometry in pruning, our paper clears much of the mystery. The value of pruning is also justified.

Up to now, the results above have shown that the "no value of inheriting weights" argument in previous works is largely attributed to sup-optimal finetuning settings. A larger LR (*e.g.*, 10^{-2}) can significantly improve the finetuning performance than a small one (*e.g.*, 10^{-3}). In fact, we are not the only one to discover this. Previous works (Renda et al., 2020; Le & Hua, 2021) also reported similar observation. Nevertheless, they do not link the phenomenon with the "value of inheriting weights" argument and do not conduct systematical empirical studies as we do here. *More importantly, neither of them presented theoretical explanations about its occurrence* – next, we are about to bridge this gap. We present a faithful theoretical explanation through the lens of dynamical isometry.

3 INITIALIZATION IS THE KEY

3.1 PREREQUISITES: DYNAMICAL ISOMETRY

Dynamical isometry (DI) is studied under the topic of trainability of deep neural networks. It was first brought up in (Saxe et al., 2014). Specifically, the dynamical isometry is defined as *the singular values of the Jacobian matrix being around 1* (Saxe et al., 2014). It is easy to see that the networks with dynamical isometry is easy to train, because JSVs around 1 imply the gradient signals will not be amplified or attenuated seriously during propagation, preventing the network from gradient exploding or vanishing, which are well-known the main difficulties in deep network training (Glorot & Bengio, 2010; Sutskever et al., 2013). For linear networks, dynamical isometry can be achieved *exactly* by the orthogonal initialization proposed in (Saxe et al., 2014); while for neural networks with non-linearity (like ReLU (Nair & Hinton, 2010)) and convolution, it can only be approximated up to date (see Tab. 11 in Appendix).

Mean Jacobian singular values. Since dynamical isometry is measured by the Jacobian singular values (JSV's), we adopt *the mean of Jacobian singular values* (denoted by \overline{S}) as a scalar metric for analysis. Specifically, for a Jacobian $\mathbf{J} \in \mathbb{R}^{C \times D_{\text{in}}}$ (*C* stands for the output dimension, *i.e.*, the number of classes, D_{in} for the input dimension), apply singular value decomposition (Trefethen & Bau III, 1997) to it,

$$U, \Sigma, V = \operatorname{svd}(\mathbf{J}), \bar{S} = \frac{1}{K} \sum_{i=1}^{K} \Sigma_{ii},$$
(1)

where Σ is the singular value matrix and $K = \min(C, D_{in})$.

Properly look at mean JSV. Notably, mean JSV equal to (or close to) 1 for dynamical isometry is an ideal case. In practice, this barely holds (unless manually set). When analyzing practical networks, we need to be exact about how close is the so-called "close to" 1. In Appendix A.4, we provide an empirical study to answer this question, which can help us look at the mean JSV metric at a proper scale. Readers are encouraged to take a look.

Pruning as poor initialization for trainability. We investigate how the pruning affects \overline{S} . The results are shown in Fig. 1. As seen, the mean JSV is consistently *damaged* by pruning; and a larger pruning ratio, more decrease of the mean JSV. This means, pruning actually servers as a very poor initialization scheme for the subsequent finetuning.

In stark contrast to the broad awareness that initialization is rather critical to neural network training (Glorot & Bengio, 2010; Sutskever et al., 2013; Mishkin & Matas, 2016; Krähenbühl et al., 2016; He et al., 2015), the dynamical isometry before finetuning.



Figure 1: Mean JSV (Eq. (1)) of pruned networks w.r.t. different pruning ratios on MNIST. Note, with a larger pruning ratio, mean JSV is hurt more. initialization role of pruning has received negligible research attention, however. As far as we know, no prior works have noted this issue when pruning a pretrained network or tried to recover the broken

Simply put, our goal next is to show it is this broken dynamical isometry that answers for the performance gap between LR 10^{-2} and 10^{-3} . To validate this explanation, we need a method to *fully recover* dynamical isometry in filter pruning, as is introduced next.

3.2 DYNAMICAL ISOMETRY RECOVERY IN FILTER PRUNING

In (Saxe et al., 2014), they propose a weight orthogonalization scheme to achieve dynamical isometry for neural network *initialization*, namely, the initial weights are *randomly sampled*. Different from their case, here the initial weights are inherited from a pretrained model by pruning. Therefore, we need to adapt it to our application.

For a FC layer parameterized by a matrix $W_0 \in \mathbb{R}^{J \times K}$ (for a Conv layer parameterized by a 4-d tensor of shape $\mathbb{R}^{N \times C \times H \times W}$, it can be reshaped to a matrix of shape $\mathbb{R}^{N \times CHW}$), it reduces to matrix W of size $\mathbb{R}^{J_1 \times M_1}$ $(J_1 \leq J, K_1 \leq K)$ after structured pruning. Then, we apply the weight orthogonalization technique (Mezzadri, 2006) based on QR decomposition (Trefethen & Bau III, 1997) to W.

$$Q, R = \operatorname{qrd}(W),$$

$$W^* = Q \odot \operatorname{sign}(\operatorname{diag}(R)),$$
(2)

where qrd(·) means QR decomposition; Q is an orthogonal matrix of the same size as $W(\mathbb{R}^{J_1 \times K_1})$; R is an an upper triangular matrix of size $\mathbb{R}^{K_1 \times K_1}$; sign(·) is the sign function which returns the positive or negative sign of its argument; . represents the Hadamard (element-wise) product aligned to the last axis (since Q and sign(diag(**R**)) share the same dimension at the last axis).

As an orthogonalized version of W, W^* recovers the dynamical isometry damaged by pruning. Therefore, we propose to employ W^* instead of the original W as the initialization weights for later finetuning. We dub this weight orthogonalization method for *pruned* models as **OrthP**. With this method, we can fully recover broken dynamical isometry and continue our analysis as follows.

3.3 ANALYSIS WITH MLP-7-LINEAR ON MNIST

Evaluated network. The network for analysis is a 7-layer linear MLP. We are aware that this toy network has little practical meaning, but it is very appropriate here for two reasons. First, as mentioned above, in our analysis we need a method to recover DI exactly. Up to date, this can only be achieved on linear networks (see Tab. 11). Second, the linear MLP network is free from the intervention of modern CNN features (*e.g.*, BN (loffe & Szegedy, 2015), residual (He et al., 2016)). By our observation, these features will make the problem complex and prevent us from seeing consistent results at the early analysis stage.

LR schedule setup. When we set different LR schedules with different initial LRs, we will (1) keep the total number of epochs the same, (2) keep the last LR the same, (3) intentionally use *prolonged* training epochs (like 90 or even 900 epochs on MNIST dataset. Typically, it only needs 30 epochs to reach convergence on MNIST). All of these are to ensure the networks are fully converged, helping us render a faithful conclusion. Step decay LR schedule is employed given its broad use.

Proposed explanation and its deducted hypotheses. We list the first 10-epoch mean JSV and test accuracy of pruning MLP-7-Linear (at pruning ratio 0.9) under different finetuning setups in Tab. 4 and Tab. 5, respectively. The following observations are worth our attention.

(1) In Tab. 4, one important fact is that, the mean JSV can recover itself without any extra help during finetuning, regardless of different setups.

(2) In Tab. 4, without OrthP, the mean JSV of LR 10^{-2} arises *much faster* than that of LR 10^{-3} . By our analysis in the Appendix A.4), a larger mean JSV implies better dynamical isometry, which further implies easier optimization. Easier optimization finally leads to the better test accuracy of LR 10^{-2} against 10^{-3} when the training epoch is insufficient, as shown in Tab. 5.

(3) In Tab. 4, with OrthP, the broken dynamical isometry is exactly recovered (note at epoch 0, the mean JSV is 1). Then the accuracy advantage of LR 10^{-2} over 10^{-3} is not significant anymore, *e.g.*, in Tab. 5, at epoch 10, LR 10^{-2} is better than LR 10^{-3} by 16.05% without OrthP; while it is not better (actually worse) than LR 10^{-3} with OrthP.

(4) Particularly note how the mean JSV trend in Tab. 4 correlates well with the test accuracy trend in Tab. 5. This resonates with the possibility that it is dynamical isometry that answers for the accuracy gap fundamentally.

These observations inspire us to the following plausible explanation about why a larger finetuning LR can improve performance significantly:

A larger finetuning LR helps the network update faster, thus the dynamical isometry recovers faster (and possibly better), which further leads to faster (and possibly better) optimization. Meanwhile, better optimization usually implies better generalization as we observe in practice, thus the larger finetuning LR eventually leads to the better test accuracy as we see, especially when the finetuning epochs are limited.

That is, a larger finetuning LR shows performance advantage *only if dynamical isometry is broken by pruning first*. If dynamical isometry is fully recovered before finetuning, a larger LR should *not* pose performance advantage anymore. The validation of this explanation can be specified into the following 4 deducted hypotheses:

Hypothesis 1. Given a small number of epochs, mean JSV *cannot* be fully recovered by training, then the larger LR should show a significant advantage over the smaller LR.

Hypothesis 2. With sufficient training epochs, mean JSV can be fully recovered by training. Then the larger LR should have less advantage over the smaller LR.

Hypothesis 3. If we employ OrthP to exactly recover the mean JSV, given the small number of epochs again, the larger LR should have much less advantage now.

Hypothesis 4. If we combine abundant epochs with OrthP, mean JSV will be recovered even more, then the performance gap between the larger LR and the smaller one should be even smaller.

Corresponding to these 4 hypotheses, the 8 finetuning LR settings are summarized in Tab. 6. The unpruned MLP model is trained with LR schedule "90 epochs, initial 10^{-2} , decay 30/60". For filter pruning, we employ L_1 -norm pruning (Li et al., 2017) throughout this paper. Specifically, it sorts the neurons (or filters) by their L_1 norms in ascending order and prunes those with the least norms by a predefined pruning ratio r.

The final accuracy results are shown in Tab. 7. We first analyze the results of pruning ratio 0.8 in Tab. 7. As seen, when finetuned for 90 epochs, LR 10^{-2} shows an advantage over LR 10^{-3} by

Table 6: Summary of finetuning LR setups corresponding to the 4 proposed hypotheses in Sec. 3.3.

	Initial LR 10^{-2}	Initial LR 10^{-3}
For Hypothesis 1	90 epochs, Initial 10^{-2} , decay 30/60	90 epochs, Initial 10^{-3} , decay 45
For Hypothesis 2	900 epochs, Initial 10^{-2} , decay 300/600	900 epochs, Initial 10^{-3} , decay 450
For Hypothesis 3	OrthP, 90 epochs, Initial 10^{-2} , decay 30/60	OrthP, 90 epochs, Initial 10^{-3} , decay 45
For Hypothesis 4	OrthP, 900 epochs, Initial 10^{-2} , decay 300/600	OrthP, 900 epochs, Initial 10^{-3} , decay 450

Table 7: Test accuracies (%) of the 4 hypotheses in Tab. 6 with **MLP-7-Linear** network on MNIST. Accuracy of unpruned model: 92.77%. Each setting is randomly run 5 times, mean accuracy and stddev reported. "Acc. gain" refers to the mean accuracy improvement of LR 10^{-2} over 10^{-3} . Hyper-parameters: batch size 100, SGD optimizer, weight decay 0.0001, LR schedule: initial LR 0.01, decayed by factor 0.1 at 1/3 and 2/3 total epochs.

Finetuning setting	Pr	runing ratio 80%		Pruning ratio 90%			
i metaning setting	$LR \ 10^{-3}$	$LR \ 10^{-2}$	Acc. gain	$LR \ 10^{-3}$	$LR \ 10^{-2}$	Acc. gain	
90 epochs	$90.54_{\pm 0.02}$	$91.36_{\pm 0.02}$	0.82	$87.59_{\pm 0.01}$	$87.81_{\pm 0.03}$	0.22	
900 epochs	$92.54_{\pm 0.03}$	$91.64_{\pm 0.41}$	-0.90	$90.44_{\pm0.01}$	$87.83_{\pm 0.04}$	-2.61	
OrthP, 90 epochs	$92.77_{\pm 0.03}$	$92.79_{\pm 0.03}$	0.02	92.72 ± 0.03	$92.77_{\pm 0.04}$	0.05	
OrthP, 900 epochs	$92.84_{\pm 0.03}$	$92.81_{\pm 0.04}$	-0.03	$92.86_{\pm 0.03}$	$92.79_{\pm 0.03}$	-0.07	

0.82% accuracy. It is tempting to draw a conclusion based on this comparison that LR 10^{-2} is much better than LR 10^{-3} . However, this is *not* the whole story:

- With 900 epochs, LR 10^{-2} is greatly surpassed by LR 10^{-3} (91.64 vs. 92.54). The reason by our proposed explanation is that, with abundant epochs, the dynamical isometry can be recovered more completely, hence LR 10^{-2} does not show advantages anymore over 10^{-3} .
- When OrthP applied, LR 10^{-2} does not show significant advantages either, similar to the effect of increasing the number of training epochs. This is because that finetuning shares the same role of recovering dynamical isometry with OrthP. Just OrthP is more effective since it is analytically targeting exact isometry.
- When the best setting used (OrthP + 900 epochs), LR 10^{-3} is slightly better than 10^{-2} . Comparing "OrthP, 900 epochs" with "OrthP, 90 epochs", the gains are only marginal. This is because the dynamical isometry has already been fully recovered by OrthP, thus more training epochs do not show much value anymore.

A different pruning ratio 0.9 is also explored. Its results are in line with those of pruning ratio 0.8 as shown in Tab. 7. In short, these empirical observations are *fully in line with* our expectations, justifying the validity of the proposed explanation.

Explaining LR effect on ImageNet. In Tab. 2, there is an apparent trend that the *larger* pruning ratio, the *more* performance advantage of LR 10^{-2} over 10^{-3} . Using our explanation, this phenomenon can also be explained now – When the pruning ratio is greater, more dynamical isometry is damaged. LR 10^{-2} can find more use in these cases since it is faster/better to recover dynamical isometry, hence the more pronounced advantage.

4 CONCLUSION

In this work, we present extensive pieces of empirical evidence to show the "no value of inheriting weights" argument in prior works is inaccurate because of improper finetuning LR schedules. We further tap into dynamical isometry to explain why the finetuning LR has such a great impact on the final performance, through carefully designed control experiments with 4 hypotheses. We promote looking at pruning as a kind of initialization, which is a favorable perspective that can make seemingly inconsistent and random results become predictable and coherent.

Our investigation of trainability (dynamical isometry) in structured pruning justifies the value of inheriting weights, in line with the past research wisdom and our common beliefs. It also helps us towards a better understanding of pruning and possibly can inspire more advanced pruning algorithms as dynamical isometry recovery has been shown a worthy direction in the paper. In addition, the awareness of dynamical isometry in structured pruning can help us render a more faithful conclusion when comparing different pruning methods (*e.g.*, pay more attention to finetuning).

REFERENCES

- Davis Blalock, Jose Javier Gonzalez, Jonathan Frankle, and John V Guttag. What is the state of neural network pruning? In *SysML*, 2020. 3
- Jian Cheng, Pei-song Wang, Gang Li, Qing-hao Hu, and Han-qing Lu. Recent advances in efficient computation of deep convolutional neural networks. *Frontiers of Information Technology & Electronic Engineering*, 19(1):64–77, 2018a. 3
- Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. Model compression and acceleration for deep neural networks: The principles, progress, and challenges. *IEEE Signal Processing Magazine*, 35 (1):126–136, 2018b. 3
- Elliot J Crowley, Jack Turner, Amos Storkey, and Michael O'Boyle. A closer look at structured pruning for neural network compression. *arXiv preprint arXiv:1810.04622*, 2018. 1, 3, 4, 5, 12
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 1
- Lei Deng, Guoqi Li, Song Han, Luping Shi, and Yuan Xie. Model compression and hardware acceleration for neural networks: A comprehensive survey. *Proceedings of the IEEE*, 108(4): 485–532, 2020. 3
- Xiaohan Ding, Guiguang Ding, Jungong Han, and Sheng Tang. Auto-balanced filter pruning for efficient convolutional neural networks. In *AAAI*, 2018. 3
- Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *JMLR*, 20(55):1–21, 2019. 3
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *ICLR*, 2019. **3**
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis. In *ICML*, 2020. 3
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M Roy, and Michael Carbin. Pruning neural networks at initialization: Why are we missing the mark? In *ICLR*, 2021. 3
- Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. *arXiv* preprint arXiv:1902.09574, 2019. 3
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010. 6, 7
- Song Han, Jeff Pool, John Tran, and William J Dally. Learning both weights and connections for efficient neural network. In *NeurIPS*, 2015. 2, 3, 4
- Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. In *ICLR*, 2016. 2, 3, 5
- B. Hassibi and D. G. Stork. Second order derivatives for network pruning: Optimal brain surgeon. In *NeurIPS*, 1993. 3
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 8, 15
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *CVPR*, 2015. 7
- Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *ICCV*, 2017. 2, 4, 5
- Zehao Huang and Naiyan Wang. Data-driven sparse structure selection for deep neural networks. In *ECCV*, 2018. 4

- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 8
- Philipp Krähenbühl, Carl Doersch, Jeff Donahue, and Trevor Darrell. Data-dependent initializations of convolutional neural networks. In *ICLR*, 2016. 7
- Duong H Le and Binh-Son Hua. Network pruning that matters: A case study on retraining variants. In *ICLR*, 2021. 1, 6
- Y. LeCun, J. S. Denker, and S. A. Solla. Optimal brain damage. In NeurIPS, 1990. 1, 3
- Namhoon Lee, Thalaiyasingam Ajanthan, and Philip Torr. Snip: Single-shot network pruning based on connection sensitivity. In *ICLR*, 2019. 3
- Namhoon Lee, Thalaiyasingam Ajanthan, Stephen Gould, and Philip HS Torr. A signal propagation perspective for pruning neural networks at initialization. In *ICLR*, 2020. **3**, 13
- Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *ICLR*, 2017. 2, 3, 4, 5, 6, 8, 12, 13, 15
- Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. In ICCV, 2017. 4
- Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network pruning. In *ICLR*, 2019. 1, 2, 3, 4, 5, 13
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *ICLR*, 2017. 14
- Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural network compression. In *ICCV*, 2017. 4
- Huizi Mao, Song Han, Jeff Pool, Wenshuo Li, Xingyu Liu, Yu Wang, and William J Dally. Exploring the granularity of sparsity in convolutional neural networks. In *CVPR Workshop*, 2017. 2
- Francesco Mezzadri. How to generate random matrices from the classical compact groups. *arXiv* preprint math-ph/0609050, 2006. 7
- Dmytro Mishkin and Jiri Matas. All you need is a good init. In ICLR, 2016. 7
- Michael C Mozer and Paul Smolensky. Skeletonization: A technique for trimming the fat from a network via relevance assessment. In *NeurIPS*, 1989. 1
- Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010. 6
- Vivek Ramanujan, Mitchell Wortsman, Aniruddha Kembhavi, Ali Farhadi, and Mohammad Rastegari. What's hidden in a randomly weighted neural network? In CVPR, 2020. 3
- R. Reed. Pruning algorithms a survey. *IEEE Transactions on Neural Networks*, 4(5):740–747, 1993. 1, 3
- Alex Renda, Jonathan Frankle, and Michael Carbin. Comparing rewinding and fine-tuning in neural network pruning. In *ICLR*, 2020. 1, 6
- Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *ICLR*, 2014. 2, 6, 7, 13, 15
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 5
- Sidak Pal Singh and Dan Alistarh. Woodfisher: Efficient second-order approximations for model compression. In *NeurIPS*, 2020. 3
- Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *ICML*, 2013. 6, 7

- Vivienne Sze, Yu-Hsin Chen, Tien-Ju Yang, and Joel S Emer. Efficient processing of deep neural networks: A tutorial and survey. *Proceedings of the IEEE*, 105(12):2295–2329, 2017. 1, 3
- Lucas Theis, Iryna Korshunova, Alykhan Tejani, and Ferenc Huszár. Faster gaze prediction with dense networks and fisher pruning. *arXiv preprint arXiv:1801.05787*, 2018. **3**, 12
- Lloyd N Trefethen and David Bau III. Numerical linear algebra, volume 50. Siam, 1997. 6, 7
- Chaoqi Wang, Roger Grosse, Sanja Fidler, and Guodong Zhang. Eigendamage: Structured pruning in the kronecker-factored eigenbasis. In *ICML*, 2019a. 3
- Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by preserving gradient flow. In *ICLR*, 2020. 3
- Huan Wang, Qiming Zhang, Yuehai Wang, Lu Yu, and Haoji Hu. Structured pruning for efficient convnets via incremental regularization. In *IJCNN*, 2019b. 3
- Huan Wang, Can Qin, Yulun Zhang, and Yun Fu. Neural pruning via growing regularization. In *ICLR*, 2021. 3, 15
- Huan Wang, Can Qin, Yulun Zhang, and Yun Fu. Recent advances on neural network pruning at initialization. In *IJCAI*, 2022. **3**
- Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. In *NeurIPS*, 2016. 2, 5
- Barret Zoph and Quoc Le. Neural architecture search with reinforcement learning. In *ICLR*, 2017. 3

A APPENDIX

A.1 RESULTS OF REEXAMINATION OF CROWLEY ET AL. (2018)

Table 8: Test accuracy comparison between 2 pruning schemes and 4 scratch training schemes in (Crowley et al., 2018). Network: **WRN-40-2** (unpruned accuracy: 95.08%, params: 2.24M). Dataset: **CIFAR-10**. Results above the dashline are directly cited from (Crowley et al., 2018); results below the dashline are from our reproducing (with the official code of (Crowley et al., 2018) at https://github.com/BayesWatch/pytorch-prunes for fair comparison). "Rerun" means we rerun the code of (Crowley et al., 2018) as it is. "LR 10^{-2} " means we redo the finetuning for the pruned models in "Rerun" using our finetuning LR schedule (120 epochs, initial 10^{-2} , decay 60/90). Finetuning is randomly repeated for 3 times, mean (stddev) accuracies reported. The main point here is that (Crowley et al., 2018) draws the conclusion that scratch training is better than pruning *because of an improper finetuning LR scheme*. With the proper finetuning LR scheme, pruning is actually *better* than scratch training.

Mathad	Saratah	500K para	500K params budget		ns budget	1.5M params budget	
Method	Scraten:	Params (M)	Acc. (%)	Params (M)	Acc. (%)	Params (M)	Acc. (%)
L_1 -norm pruning (Li et al., 2017)	X	0.51	90.86	1.02	92.61	1.52	93.63
Fisher pruning (Theis et al., 2018)	X	0.52	92.59	1.02	93.51	1.52	94.51
Varying Depth	1	0.69	93.56	1.08	<u>94.54</u>	1.47	94.64
Varying Width	1	0.50	93.45	0.98	94.30	1.48	94.66
Varying Bottleneck	1	0.50	93.69	1.00	94.40	1.49	94.79
Fisher Scratch	1	0.52	93.72	1.02	94.65	1.52	94.86
L_1 -norm pruning (Li et al., 2017) (Rerun)	X	0.50	91.23	1.00	92.80	1.51	93.52
L_1 -norm pruning (Li et al., 2017) (LR 10 ⁻²)	X	0.50	<u>93.88</u> (0.10)	1.00	<u>94.49</u> (0.10)	1.51	<u>94.92</u> (0.16)
Fisher pruning (Theis et al., 2018) (Rerun)	X	0.52	92.17	0.98	93.57	1.48	94.67
Fisher pruning (Theis et al., 2018) (LR 10^{-2}) 🗶	0.52	94.27 (0.09)	0.98	94.80 (0.02)	1.48	95.10 (0.13)

U	2		
Dataset	MNIST	CIFAR10	ImageNet
Solver	SGD (0.9, 1e-4)	SGD (0.9, 5e-4)	SGD (0.9, 1e-4)
LR schedule (Initial 1e-2)	Multi-step (decay 30/60)	Multi-step (decay 60/90)	Multi-step (decay 30/60/75)
LR schedule (Initial 1e-3)	Multi-step (decay 45)	Multi-step (decay 80)	Multi-step (decay $45/68$)
Total epoch	90	120	90
Batch size	100	128	256
Data augmentation	None	Random crop and horizontal flip	Random crop and horizontal flip

Table 9: Training setting summary in finetuning. For the SGD solver, in the parentheses are the momentum and weight decay.

A.2 TRAINING SETTING SUMMARY

There are three datasets in our experiments in the paper: MNIST, CIFAR10, and ImageNet. Apart from some key settings stated in the paper, a more detailed training setting summary is shown as Tab. 9. We use 8 NVIDIA V100 GPUs for all our experiments.

A typical finetuning LR schedule looks like this in our paper: "90 epochs, initial 10^{-2} , decay 30/60". Its meaning is: the total number of training epochs is 90; initial LR is 10^{-2} ; at epoch 30 and 60, LR is multiplied by a default factor 0.1. The others can be inferred likewise.

A.3 WHY OUR REPRODUCED RESULTS ARE BETTER THAN THOSE REPORTED IN (LIU ET AL., 2019)

In Tab. 1, we present the results of ResNet34 on ImageNet, pruned by the L_1 -norm pruning method (Li et al., 2017). Readers may be curious that why our scratch-training results are much higher than those of Rethinking (Liu et al., 2019) (by 0.4-0.6%). This has a historical reason. (Liu et al., 2019)² refers to the official PyTorch ImageNet example³, where the LR is only decayed twice (from 1e-1 to 1e-2, then to 1e-3). However, by our empirical observation, 1e-3 is still not the lowest LR that the network converges. If the LR is decayed another time (to 1e-4), the top-1 accuracy can still bump by round 0.5-0.8% point, which is a significant improvement for ImageNet thus not negligible. Therefore, we choose to decay the LR further to 1e-4 and finally to 1e-5 to ensure the network is *fully converged*. This reiterates the comparison rule of our paper: comparing different methods *in their best shape*. Comparing scratch training to L_1 -norm pruning *before* the finetuned network finally converges may appear fair (given the same number of epochs) but is of less practical meaning and may hide the true picture.

A.4 How to *properly* look at the mean JSV metric for dynamical isometry

In the main paper, we use mean JSV as the metric to measure dynamical isometry, for the following 3 specific reasons: (1) It was used by Lee et al. (2020), which analyzes the dynamical isometry for randomly initialized network. Given its success there, it is very natural for us to also employ this metric for analyzing pretrained networks here. (2) We currently do not have a better alternative, either. (3) In practice, we find mean JSV is informative (*e.g.*, in Tab. 7, the mean JSV trend is well-correlated with the test accuracy trend), provided we see it *properly*.

This section is meant to provide more background regarding how to look at the mean JSV for dynamical isometry *properly*. It is not a new invention of our paper but a general practical guideline about the relationship between mean JSV and dynamical isometry in order to help readers better understand our paper.

DI (dynamical isometry) is defined by mean JSV close to 1 in Saxe et al. (2014). Rigorously, in Saxe et al. (2014), DI describes the distribution of *all* JSVs. Mean JSV is only an average sketch of the distribution. Nevertheless, this average approximation is accurate enough for analysis. In other words, *if a network has mean JSV close to 1, we can say this network has dynamical isometry*.

Then, a non-trivial technical question is: When we deal with practical DNNs in the real world, how close is the so-called "close to 1"? To our best knowledge, there is no outstanding theory to

²https://github.com/Eric-mingjie/rethinking-network-pruning/tree/master/imagenet/l1-norm-pruning ³https://github.com/pytorch/examples/tree/master/imagenet

quantify this, so we resort to empirical analysis, specifically on the MLP-7-Linear network used in the main paper.

Table 10: Mean JSV and test accuracies (%) of MLP-7-Linear on MNIST under different pruning ratios. Each result is randomly run for 3 times. We report the mean accuracy and (std). "ft." is short for finetuning.

Pruning ratio	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
mean JSV	2.4987	1.7132	0.9993	0.5325	0.2711	0.1180	0.0452	0.0151	0.0040	0.0004
Acc. before ft.	92.77	91.35	78.88	62.21	32.14	11.47	9.74	9.74	9.74	9.74
Acc. after ft.	/	92.82 (0.05)	92.80 (0.04)	92.80 (0.01)	92.77 (0.01)	92.77 (0.02)	92.77 (0.00)	92.78 (0.02)	91.37 (0.03)	87.82 (0.03)

As seen, there is a clear trend in Tab. 10: *larger* pruning ratio, *smaller* mean JSV, and *lower* test accuracy (either before or after finetuning).

Particularly note the mean JSV range where the pruned network can be *finetuned back to the original accuracy* (92.77%), which is ≥ 0.0151 . This means, for networks with mean JSV greater than 0.0151, in spite that their immediate accuracies (without finetuning) can be distinct (*e.g.*, 91.35% at PR 0.1 vs. 9.74% at PR 0.7), intrinsically, they are *equivalently potential* after finetuning.

DI theory suggests that mean JSV equal to 1 is the best case. Then we would ask, how about mean JSV equal to 2? 1.5? 0.8? Are they really worse than the ideal value 1? Tab. 10 above shows *not necessarily*, because a network with mean JSV 0.2711 can reach comparable accuracy to the network with mean JSV 1. Only when the mean JSV is smaller than some threshold (in this case, 0.0151), it leads to an irrecoverable damage to the network optimization, which eventually results in lower generalization ability (*i.e.*, test accuracy).

By this "trainable with equivalent potential" rule, if a mean JSV lies in the range of >= 0.0151, we can regard it as "close to 1" because they can do *just as well as* 1.

Here, we only mention the *lower* bound, how about the *upper* bound? Can mean JSV 10,000 also be regarded as "close to 1"? This is a good question worth further dedicated investigation. Yet, for now, we do not need to worry much about it, because in practice, a *normally* trained network *rarely* presents a very large mean JSV by our empirical observation, but is quite likely to have a very small mean JSV (0.0004 at pruning ratio 0.9 in Tab. 10 is a concrete example).

Other metrics for dynamical isometry? Given dynamical isometry is defined as mean JSV equal to 1, a seemingly plausible metric for dynamical isometry is the mean deviation-from-unity: $\frac{1}{K}\sum_{i=1}^{K} (\Sigma_{ii} - 1)^2$, where Σ is the singular value matrix. We need to point out that this metric is not a good option in fact.

Here is a simple example to see. Consider we have two initialization schemes for a neural network, A: mean JSV = 0 (this can be done by removing one internal layer; since the network is disconnected, its mean JSV is definitely 0), B: mean JSV = 3 (this can be done by re-scaling the weights in each layer by $3^{1/7}$ after using the orthogonal initialization, for the MLP-7-Linear network). Then the former has mean deviation-from-unity of 1, the latter has mean deviation-from-unity of 4. By the metric, the former should be better, but clearly it is not since the network is not trainable at all while the latter is trainable – we empirically confirmed this with PyTorch on the MLP-7-Linear network. This is why we think the metric may not be better than mean JSV.

A.5 MORE RESULTS

Effect of different LR schedules. In the main paper, we mention the finetuning LR effect (*i.e.*, a larger initial finetuning LR makes the final performance better) also appear using a different LR schedule, rather than the common step decay LR schedule. Here is the example.

We consider *Cosine Annealing LR schedule* Loshchilov & Hutter (2017) here, referring to the official PyTorch Cosine LR implementation⁴. When we switch from Step LR schedule to Cosine, the initial LR and minimum LR are kept the same (namely, the start point and end point of LR are the same; the only difference is the scheduling in between). The scratch model is trained for 200 epochs, initial

⁴https://pytorch.org/docs/stable/generated/torch.optim.lr_scheduler.CosineAnnealingLR.html

Table 11: JSVs (Jacobian singular values) of orthogonal initialization (Saxe et al., 2014) on diffe	er-
ent types of neural networks on MNIST dataset. Note, only the linear MLP network can achie	ve
dynamical isometry <i>exactly</i> (<i>i.e.</i> , all the JSVs equal to 1).	

Network	Mean JSV	Max JSV	Min JSV
MLP-7-Linear	$1.0000_{\pm 0.0000}$	1.0000	1.0000
MLP-7-ReLU	$1.2268_{\pm 0.5519}$	3.2772	0.2282
LeNet-5-Linear	$0.9983_{\pm 0.0842}$	1.2330	0.7896
LeNet-5-ReLU	$1.8331_{\pm 0.5731}$	3.6007	0.6151

LR 0.1, stp decay at epoch 100 and 150 by multiplier $\frac{1}{10}$ (referring to the original ResNet CIFAR10 training recipe in the ResNet paper He et al. (2016)). For finetuning of pruning methods (L_1 Li et al. (2017) and GReg-1 Wang et al. (2021) methods are considered here as they are representative works of importance-based pruning and regularization-based pruning, the two major paradigms of filter pruning), the initial LR is 0.01 or 0.001, the minimum LR 0.0001, total epochs 120.

Table 12: Effect of LR schedule of ResNet56 on CIFAR10. Baseline accuracy 93.78%, Params: 0.85M, FLOPs: 0.25G. The best in each column is highlighted in red, second best in blue. For each sub-block (separated by the horizontal dashline), the **relatively better** one is highlighted in **bold**.

Pruning ratio Sparsity/Speedup	0.3 31.14%/1.45×	0.5 49.82%/1.99×	0.7 70.57%/3.59×	0.9 90.39%/11.41×
Scratch (Step LR) Scratch (Cosine LR)	$93.16_{\pm 0.16} \\ \textbf{93.84}_{\pm 0.06}$	$\begin{array}{c} 92.78_{\pm 0.23} \\ \textbf{93.20}_{\pm 0.31} \end{array}$	$\begin{array}{c} 92.11_{\pm 0.12} \\ \textbf{92.15}_{\pm 0.21} \end{array}$	$\frac{\textbf{88.36}_{\pm 0.20}}{88.17_{\pm 0.43}}$
L_1 Li et al. (2017) (Step LR 0.001)	$93.43_{\pm 0.06}$	$93.12_{\pm 0.10}$	$91.77_{\pm 0.11}$	87.57 _{±0.09}
L_1 Li et al. (2017) (Step LR 0.01)	$93.79_{\pm 0.06}$	93.51 _{±0.07}	$92.26_{\pm 0.17}$	$86.75_{\pm 0.31}$
L_1 Li et al. (2017) (Cosine LR 0.001)	$93.48_{\pm 0.04}$	$93.11_{\pm 0.09}$	$91.65_{\pm 0.11}$	87.17 _{±0.14}
L_1 Li et al. (2017) (Cosine LR 0.01)	$93.82_{\pm 0.07}$	93.74 ±0.06	92.27 ±0.00	$86.90_{\pm 0.20}$
GReg-1 Wang et al. (2021) (Step LR 0.01)	$94.00_{\pm 0.10}$	$93.36_{\pm 0.03}$	$92.43_{\pm 0.06}$	89.59 _{±0.03}

Table 13: Effect of LR schedule of VGG19 on CIFAR100. Baseline accuracy: 74.02%, Params: 20.08M, FLOPs: 0.80G. The best in each column is highlighted in red, second best in blue. For each sub-block (separated by the horizontal dashline), the **relatively better** one is highlighted in **bold**

0.3 19.24%/1.23×	0.5 51.01%/1.97×	0.7 74.87%/3.60×	0.9 90.98%/8.84×
$72.84_{\pm 0.25}$	71.88 $_{\pm 0.14}$	70.79 $_{\pm 0.08}$	66.52 ±0.37
$73.54_{\pm 0.22}$	$71.87_{\pm 0.09}$	$70.10_{\pm 0.24}$	$65.92_{\pm 0.10}$
$73.67_{\pm 0.05}$	$72.04_{\pm 0.12}$	$70.21_{\pm 0.02}$	$64.72_{\pm 0.17}$
$74.01_{\pm 0.18}$	$73.01_{\pm 0.22}$	71.49 $_{\pm 0.14}$	$66.05_{\pm 0.04}$
$73.69_{\pm 0.08}$	$72.10_{\pm 0.08}$	$69.96_{\pm 0.09}$	$63.93_{\pm 0.15}$
74.39 _{±0.07}	$73.51_{\pm 0.18}$	71.78 ±0.21	65.70 ±0.11
$74.03_{\pm 0.11}$	$73.21_{\pm 0.08}$	$71.64_{\pm 0.08}$	$67.75_{\pm 0.16}$
	$\begin{array}{r} 0.3\\ 19.24\%/1.23\times\\ \hline 72.84_{\pm 0.25}\\ \textbf{73.54}_{\pm 0.22}\\ \hline \textbf{73.67}_{\pm 0.05}\\ \textbf{74.01}_{\pm 0.18}\\ \hline \textbf{73.69}_{\pm 0.08}\\ \textbf{74.39}_{\pm 0.07}\\ \hline \textbf{74.03}_{\pm 0.11}\\ \end{array}$	$\begin{array}{c ccccc} 0.3 & 0.5 \\ \hline 19.24\%/1.23\times & 51.01\%/1.97\times \\ \hline 72.84_{\pm 0.25} & \textbf{71.88}_{\pm 0.14} \\ \hline \textbf{73.54}_{\pm 0.22} & 71.87_{\pm 0.09} \\ \hline 73.67_{\pm 0.05} & 72.04_{\pm 0.12} \\ \hline \textbf{74.01}_{\pm 0.18} & \textbf{73.01}_{\pm 0.22} \\ \hline 73.69_{\pm 0.08} & 72.10_{\pm 0.08} \\ \hline \textbf{74.03}_{\pm 0.11} & \textbf{73.21}_{\pm 0.08} \\ \hline \end{array}$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$

The results of ResNet56 (on CIFAR10) and VGG19 (on CIFAR100) are presented in Tab. 12 and Tab. 13. The following observations are worth our attention:

(1) The advantage of initial LR 0.01 over 0.001 does not only appear with the Step LR schedule, but also appears with the Cosine LR schedule. This implies the finetuning LR effect is *generic*, not limited to one particular LR schedule, which further highlights the importance of the topic we have been studying in the paper.

(2) Comparing the best scratch result to the best pruning result at each sparsity, **there is no case that scratch training outperforms pruning**. In most cases, scratch training is not even the second best (highlighted in blue).

A.6 THE CONNECTION BETWEEN DYNAMICAL ISOMETRY AND GENERALIZATION

Dynamical isometry is an indicator for trainability. Rigorously, trainability only implies a easy optimization process (namely, closely related to the optimization *speed*, not necessarily implying anything about optimization *quality*). How is it related to generalization?

The reason that a better trainability typically implies a better generalization in practice is that the network has multiple local minima. Converging to the good local minima takes many epochs of training. When the number of training epochs is limited (in practice we cannot train the network for infinite epochs), the networks with a worse trainability will "lag behind" compared to those with a good trainability, which finally *pose as a worse test accuracy* when the training is finished. Note here, by "pose", we mean they are *not necessarily* worse if unlimited training epochs are given. Just allowing that long time for the network to explore a better local minimum is not practically meaningful. All the comparison in this work is still based on the practical application scenario, *i.e.*, the network is trained for *a fair number* of epochs. Exploring the extreme cases (*e.g.*, given a very long training) is *not* the interest of this paper (and most past pruning papers, either).

To understand the stance of this paper, here is a quick example. Consider a network which normally takes K epochs to convergence. Now, pruning is shown better than scratch training with K (or fairly more epochs, *e.g.*, 2*K) epochs. Meanwhile, allowing the training time to be longer, such as 10*K, scratch training can match pruning. In this case, we still consider pruning has value because at least it saves the total training epochs even if it does not show a final performance advantage.

A.7 LIMITATIONS AND POTENTIAL NEGATIVE SOCIETAL IMPACTS

Limitations. Essentially, this paper is an empirical work with connections to theoretical explanations (such the dynamical isometry part). We draw our conclusions based on the experiments we have. It is well-known that deep learning typically has heavy hyper-parameters. We are not 100% sure they are set most correctly, especially in the sparse case. This uncertainty passes on to our conclusion inevitably.

Yet, we are doing our best following the common practices in the area and repeat multiple experiments to exclude the influence of irrelevant statistical variations. Based on the empirical proofs we have, the proposed explanation is self-consistent and can translate across different networks / datasets. This said, a more rigorous understanding (in terms of mathematical proofs) is still desired.

Potential Negative Societal Impacts. Simply put, this work proposes theories and algorithms that can boost the efficiency of a network by network pruning, that is, make it smaller, faster, and possibly consume less energy in practical applications. We focus on the classification task, which is generally the foundation of the many up-stream computer vision tasks like detection and segmentation. Therefore, this work potentially has a broad application especially in the computer vision areas.

The algorithm itself has few negative societal issues, but when it makes many AI-driven technologies applicable in practice, the impact really depends on how humans put them into use in practice. This actually falls into the general ethical discussion on whether AI is good or not. Beyond this scope, this work does not have specific negative societal impacts brought by its potential application, to our best knowledge.