

# FGGP: Fixed-Rate Gradient-First Gradual Pruning

Anonymous Full Paper  
Submission 36

## Abstract

In recent years, the increasing size of deep learning models and their growing demand for computational resources have drawn significant attention to the practice of pruning neural networks, while aiming to preserve their accuracy. In unstructured gradual pruning, which sparsifies a network by gradually removing individual network parameters until a targeted network sparsity is reached, recent works show that both gradient and weight magnitudes should be considered. In this work, we show that such mechanism, e.g., the order of prioritization and selection criteria, is essential. We introduce a *gradient-first* magnitude-next strategy for choosing the parameters to prune, and show that a *fixed-rate* subselection criterion between these steps works better, in contrast to the annealing approach in the literature. We validate this on CIFAR-10 dataset, with multiple randomized initializations on both VGG-19 and ResNet-50 network backbones, for pruning targets of 90, 95, and 98% sparsity and for both initially dense and 50% sparse networks. Our proposed **fixed-rate gradient-first gradual pruning (FGGP)** approach outperforms its state-of-the-art alternatives in most of the above experimental settings, even occasionally surpassing the upperbound of corresponding dense network results, and having the highest ranking across the considered experimental settings.

## 1 Introduction

In deep-learning for a given problem setting, typically first a network architecture is engineered (hand-crafted) and then the parameters of such network are learned. However, even when the problem setting has an established solution with a known network architecture, the required number of features/filters, contextual depth, layer sizes, and other architectural settings often need to be adjusted empirically to achieve optimal results. Alternatively, overparameterized deep neural networks are used, as most state-of-the-art today, aiming to capture hidden patterns in the data without manually optimizing architectures. Overparameterization, however, comes with largely increased computational costs at both training and inference time; requiring more energy, yielding higher CO<sub>2</sub> emissions, and making models less suitable for time critical tasks and embedded/edge/mobile computing. Overparameterized

models are also more likely to overfit the training data, hence yielding reduced performance especially without suitable regularization treatments, due to suboptimal optimization approaches or insufficient time required for lengthy training.

Neural Architecture Search (NAS) [1], a type of meta-learning and a subfield of automated machine learning (AutoML), aims to find optimal NN architectures. NAS typically treats networks as black-box models updating them based solely on observed output or prediction accuracy using methods such as evolutionary search, reinforcement learning, Bayesian approximation, etc. This typically requires significant amount of resources and does not seek principled update strategies that consider the intrinsic dynamics and parameters of a network.

Pruning is a network model compression technique that aims to remove the network parameters that are least important, i.e., that would change the accuracy minimally. Ideas of adapting network architectures started as early as the introduction of neural networks themselves, including seminal works such as “Optimal Brain Damage” (OBD) [2] by LeCun et al. For a while, research mostly focused on devising expressive neural representations, on efficient optimization strategies, and on solving practical large problems thanks to advances in compute capability. Recently, pruning has gained popularity again with an increasing focus on model compression for highly complex problems and edge computing.

Our main contributions in this paper include: (1) We provide a clear and transparent definition and review of multi-step top-K selection processes in gradual pruning. (2) We show the order prioritization and selection criteria both being essential and inter-related for a successful gradual pruning algorithm. (3) We propose a gradient-first top-K selection criterion that performs well with a fixed-rate selection quota. (4) We set the new state-of-the-art in gradual pruning for CIFAR-10 dataset.

## 2 Background

Network pruning was shown by Frankle and Carbin [3] and Liu et al. [4] to achieve similar or even better classification performance than corresponding dense models, with less than half of the original parameters. These helped draw further attention to the redundancy of state-of-the-art deep neural networks. *Structured pruning* removes neurons in

096 fully-connected (FC) or convolutional (conv) layers  
097 (the latter also known as *channel pruning*), hence  
098 not changing the layer’s original structural property,  
099 i.e., yielding a respective FC or conv layer. *Unstruc-*  
100 *tured pruning* removes weights (connections), which  
101 then typically makes the layer (and the network)  
102 *unstructured*, i.e., not a conventional FC or conv  
103 anymore.

104 **Structured pruning** aims to remove redundant  
105 channels such as based on LASSO regression [5] and  
106 discrimination-aware loss [6], or to remove redundant  
107 neurons or convolutional filters such as based on  
108 mean activation magnitude [7]. Instead of such  
109 (structured) pruning, it is more natural to decide on  
110 the (unstructured) pruning of network parameters  
111 since these are the entities that are determined via  
112 optimization during training.

113 **Scheduling.** Some methods update the neurons  
114 (filters) using a single one-shot approach, either at  
115 the very beginning (following initialization) or at the  
116 very end (after training to convergence). Pruning at  
117 start, also known as *foresight pruning*, assumes that  
118 an optimal subnetwork, which is capable of achiev-  
119 ing success at convergence, is already identifiable at  
120 initialization. For instance, SNIP [8] approximates  
121 synapse sensitivity after initialization by estimating  
122 the change in loss with respect to the removal of each  
123 parameter. To avoid numerous forward passes by  
124 removing each parameter individually, SNIP instead  
125 makes an infinitesimal (multiplicative) approxima-  
126 tion to removal that can be computed in a single  
127 forward-backward pass, and then pre-prunes the  
128 parameters with small magnitude  $|\theta_i|$  and small gra-  
129 dents  $|g_i|$  at initial state. Gradient Signal Preser-  
130 vation (GraSP) [9] revisits SNIP by considering ex-  
131 pected subsequent gradient flow using a second-order  
132 term  $|Hg|$ , which avoids explicit Hessian computa-  
133 tions; however, this leads to results not substantially  
134 different than SNIP. Despite the simplicity and at-  
135 traction of one-shot methods, superior results are  
136 often achieved using sequential pruning approaches,  
137 indicating that fixing the network structure once is  
138 not an optimal strategy.

139 **Iterative pruning** is applied repeatedly over  
140 multiple rounds, following complete convergence af-  
141 ter each round, which is hence computationally very  
142 costly. Lottery Ticket Hypothesis (LTH) [3] assumes  
143 that a pruning-candidate subnetwork has the *win-*  
144 *ning ticket* primarily thanks to its random initial-  
145 ization of parameters. LTH then trains a network,  
146 prunes the weights with smallest magnitudes, and  
147 then retrains the pruned network starting from the  
148 *same* initial random parameters, and repeats this  
149 process iteratively until a desired sparsity level is  
150 reached. Later, Liu et al. [4] confute LTH by showing  
151 that any arbitrary initialization with such pruned  
152 network achieve similar results, hence showing that  
153 the key is the architecture, not the initialization.

**Fixed-sparsity pruning** initializes a network at  
the target low sparsity and then trains this network  
to convergence while keeping its sparsity constant,  
such that the total training cost can be kept lower  
than a dense network. RigL [10] is such an exam-  
ple, which during training first selects a subset of  
weights with the smallest magnitudes to prune, and  
then momentarily sets all missing weights to zero  
to compute their gradient with a backprop, to de-  
termine the highest gradient-magnitude weights to  
add (grow) to keep the sparsity constant.

**Gradual pruning** prunes the network while it is  
being trained, slowly changing the network sparsity  
to a final targeted value, e.g. at regular iteration or  
epoch intervals some parameters are pruned based on  
a priority criterion and mechanism. This was shown  
to achieve comparable performance to iterative prun-  
ing, while incurring much lower computational costs.  
The main challenge here is that the network is not  
in a converged state during the pruning decisions.  
Medeiros et al. [11] prune connections with lower  
correlations between the errors within a layer and  
those backpropagated to the preceding layer, which  
they call the MAXCORE principle. Dynamic Net-  
work Surgery [12] employs gradual pruning with a  
binary mask for pruned/spliced connections, while  
updating both the pruned and the remaining param-  
eters. Zhu et al. [13] gradually change the network  
sparsity based on a cubic scheduling function, while  
pruning the weights with smallest magnitudes – al-  
though the weights alone are not sufficiently infor-  
mative in an uncovered network state. Dettmers  
et al. [14] utilize exponentially smoothed gradients  
(momentum) to identify layer and parameter con-  
tributions to error reduction, while both pruning  
and regrowing the connections based on momenta.  
GraNet [15] combines the pruning schedule of [13]  
with the pruning criterion of RigL [10], achieving  
the state-of-the-art results in unstructured gradual  
pruning. Note that although GraNet calls the subset  
selection process as weight “addition” (where the  
second stage is explained as if adding [back] high-  
gradient parameters), this is somewhat a misnomer  
as GraNet does not aim and cannot grow synapses  
inexistent at the beginning of pruning. In this pa-  
per, we describe GraNet with a literature-consistent  
terminology, which helps to better contrast it with  
our proposed method.

### 3 Methods

For a neural network  $f$  parametrized by  $\Theta =$   
 $\{\theta_1, \theta_2, \dots, \theta_w\}$  with  $w$  parameters, the goal of train-  
ing on a dataset  $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_K, y_K)\}$   
with  $K$  input-groundtruth pairs  $(x_i, y_i)$  is

$$\min_{\Theta} L = \sum_{i=1}^K l(f(x_i; \Theta), y_i), \quad (1)$$

where  $l(\cdot, \cdot)$  is the penalty/loss function for the distance between  $y_i$  and the network prediction  $f(x_i; \Theta)$ . The optimization problem is then solved iteratively via backpropagation (training), which is typically stabilized by using a regularization of parameters as an additional objective.

### 3.1 Pruning schedule

Let sparsity  $s$  define the number of parameters,  $N$ , in a pruned network with respect to its dense equivalent  $N^*$  as  $N = (1 - s)N^*$ . Gradual pruning reduces the network parameters slowly over training time based on a desired decay pattern (also called "schedule"). To prune a network from an initial sparsity  $s_{\text{ini}}$  at iteration  $t_{\text{ini}}$  to an intended target sparsity  $s_{\text{fin}}$  at iteration  $t_{\text{fin}}$ , we employ cubic sparsity scheduling [13] where sparsity  $s_t$  at iteration  $t$  is given by:

$$s_t = s_{\text{fin}} + (s_{\text{ini}} - s_{\text{fin}}) \left(1 - \frac{t - t_{\text{ini}}}{t_{\text{fin}} - t_{\text{ini}}}\right)^3, \quad (2)$$

which then defines the desired number of parameters at any iteration  $t$  as  $N_t = (1 - s_t)N^*$ .

Pruning events can either be applied regularly during training, e.g., every  $\Delta t$  epochs or iterations, or be at instances sampled randomly from a probability distribution. Each event will then prune  $N_p$  network parameters to reduce their number to that desired (scheduled) at that instance, i.e.,  $N_p = N_{t-\Delta t} - N_t$  assuming pruning events with  $\Delta t$  iteration interval.

### 3.2 Pruning strategy

Parameter selection criteria and mechanism have the utmost importance that can affect the outcome significantly. We motivate our choice based on the framework of OBD [2], which approximates the sensitivity of loss  $L$  to individual network parameters  $\theta_i$  using a second-order Taylor-series expansion as:

$$\delta L = \underbrace{\sum_i g_i \delta \theta_i}_{\text{1st term}} + \underbrace{\frac{1}{2} \sum_i h_{ii} \delta \theta_i^2}_{\text{2nd term}} + \underbrace{\frac{1}{2} \sum_{i \neq j} h_{ij} \delta \theta_i \delta \theta_j}_{\text{3rd term}} \quad (3)$$

where higher order terms are omitted,  $g_i$  is the gradient of  $L$  with respect to  $\theta_i$ , and  $h_{ij}$  are the elements of the Hessian matrix. Pruning a parameter, which nullifies its effect, causes a negative change equivalent to its value, i.e.,  $\delta \theta_i = -\theta_i$ . The 3rd term is often omitted by assuming minimal cross-parameter effect. If the network is already trained (i.e., converged at a local minimum), then the gradients diminish and the 1st term can be omitted as well. OBD then estimates the diagonal Hessian terms to prune parameters based on the 2nd term.

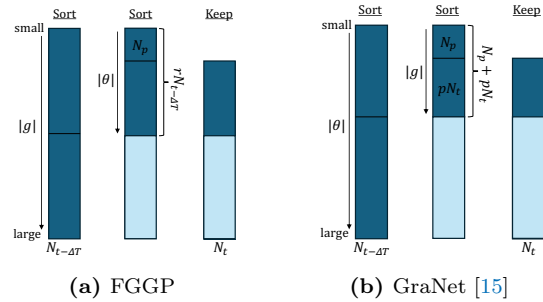
The above, however, cannot be assumed in a gradual pruning setting where the network is not converged. Assuming a simpler first-order Taylor expansion  $\delta L \approx \sum_i g_i \delta \theta_i$ , several works aim to minimize this by simply pruning parameters with **small** mag-

#### Algorithm 1 FGGP algorithmic overview

---

**Inputs:**  
1: network  $f_\Theta$ , dataset  $D$   
**Initialize**  
2: Neural network  $f_\Theta$   
3:  $s_t$  : sparsity scheduled as in (2)  
4:  $\Delta t$  : update interval  
5:  $r$  : sub-selection rate  
6: **for** each training iteration  $t$  **do**  
7:     Sample a minibatch  $B_t \sim D$   
8:     **if**  $t \equiv 0 \pmod{\Delta t}$  **then**  
9:         Sort the  $N_{t-\Delta t}$  parameters in ascending order by gradient magnitude (step 1 in Figure 1(a))  
10:         For the first  $r \cdot N_{t-\Delta t}$  parameters, sort in ascending order by weight magnitude (step 2)  
11:         Prune the first  $N_{t-\Delta t} - N_t$  parameters, so there are  $N_t$  parameters left (step 3 in Figure 1(a))  
12:     **end if**  
13:     Update parameters via backpropagation  
14: **end for**

---



**Figure 1.** Comparison of the parameter selection mechanisms between our proposed FGGP and GraNet [15].

nitudes, but this only applies if those gradients are not large. Although some recent works [8, 10, 15] consider the gradients in addition, they do this without a basis on the terms higher than the first order. In this work, we consider the gradients first, focusing on the parameters with small gradient magnitudes for which the 1st term in (3) has a basis to be omitted, and then we focus on small magnitudes that ensure the 1st term to diminish as well as the 2nd term where they appear quadratically – selecting the parameters with minimal effect on the loss.

A pseudocode of our proposed approach **fixed-rate gradient-first gradual pruning (FGGP)** is given in Algorithm 1, with the pruning criteria visualized in Figure 1(a). At every  $\Delta t$  iterations, our method chooses the parameters to prune with a two-step selection process: We first rank the parameters by their gradient magnitudes  $|g_i|$ ; we select the smallest  $rN_{t-\Delta T}$  out of these, and then rank those by their parameter magnitude  $|\theta_i|$ ; finally we select the smallest  $N_p$  of these to prune. This strategy avoids the magnitude-based selection from applying to unconverged parameters, whose values are still being changed, i.e., having large gradient magnitudes. GraNet, in contrast, applies an opposite order of

283 selection as illustrated in Figure 1(b), which may  
 284 fail by pruning important parameters if many pa-  
 285 rameters with small magnitudes are still in change  
 286 (i.e., with large gradients), since its 2nd step can  
 287 only consider those culled/selected from the 1st step.  
 288 With our subset selection order, we avoid this.

289 For pruning we consider the parameters from all  
 290 the network layers, unlike RigL [10] which does not  
 291 prune the first convolutional layer and LTH [3] which  
 292 does not prune the last fully-connected layer. We  
 293 prune the network *globally*, pooling the parameters  
 294 from all layers for their prioritization in pruning.

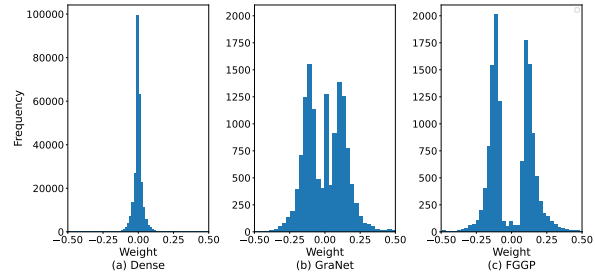
### 295 3.3 Subset selection rate

296 In the above process, an important factor is the de-  
 297 cision of which gradients to consider as large to omit  
 298 from the next steps of parameter selection. Any  
 299 fixed threshold would not be applicable, since the  
 300 parameters and gradients all have values relative to  
 301 each other, and in gradual pruning a predetermined  
 302 schedule has to be met to achieve a desired sparsity.  
 303 So, the selection needs to be based on a ratio from  
 304 a ranked (sorted) prioritization. In its magnitude-  
 305 first strategy, GraNet employs (cosine) annealing to  
 306 reduce a parameter  $p$ , seen in Figure 1(b). This re-  
 307 duces the subset considered from step 1 over training  
 308 iterations, focusing on increasingly smaller param-  
 309 eter magnitudes at later iterations. In our gradient-  
 310 first strategy, such reduction is not necessary and is  
 311 found to be counterproductive in our ablation stud-  
 312 ies. Instead we utilize a fixed rate  $r$  as the ratio of  
 313 gradient magnitudes to selected from the first step.  
 314 We herein set  $r = 0.5$  such that the parameters  
 315 with gradient magnitudes smaller than their median  
 316 value are taken into further consideration.

## 317 4 Experiments and Results

318 For evaluation we use the CIFAR-10 dataset as com-  
 319 mon in the field, allowing us to compare our re-  
 320 sults to multiple published works. We evaluate our  
 321 method based on ResNet-50 and VGG-19 architec-  
 322 tures, as were adapted for the CIFAR dataset [4].  
 323 Implementation details are given in Appendix A. We  
 324 aim for target sparsities  $s_{\text{fin}} = \{90, 95, 98\}$  in two  
 325 experimental settings of dense-to-sparse ( $s_{\text{ini}} = 0\%$ )  
 326 and sparse-to-sparse ( $s_{\text{ini}} = 50\%$ ). For initializ-  
 327 ing the parameters in sparse networks, we employ  
 328 Erdős-Rényi Kernel (ERK) [10] inline with the com-  
 329 pared state-of-the-art. See Appendix B for details.

330 In comparisons, we provide single-shot pruning  
 331 results from other methods as reference. The prun-  
 332 ing accuracies of the methods with a dense network  
 333 upperbound are seen in Table 1, where  $\pm$  results indi-  
 334 cate those from three different random **initialization**.  
 335 The single-shot methods are seen to be inferior to  
 336 the gradual methods, and among the latter the ones



**Figure 2.** Comparison of weight magnitude distribu-  
 tions at the end of training with (a) a dense network, as  
 well as pruned with (b) GraNet and (c) FGGP from 0%  
 to 95% sparsity.

that consider both gradient and parameter magni- 337  
 tudes (e.g., RigL, GraNet, and our FGGP) perform 338  
 relatively better. These differences are more pro- 339  
 nounced at higher sparsity targets and for ResNet-50 340  
 architecture, indicating that these potentially rep- 341  
 resent more challenging pruning scenarios. In most 342  
 scenarios our method FGGP outperforms the other 343  
 state-of-the-art methods. For VGG-19, our method 344  
 is more successful at higher sparsity targets of 95% 345  
 and 98%, for both sparse- and dense-to-sparse train- 346  
 ing scenarios. In both scenarios for ResNet-50, two- 347  
 out-of-three sparsity levels our method outperforms 348  
 the others – although some results may be too close 349  
 to call for a clear winner. For an overall comparison 350  
 across all experiments, we employ a ranking strat- 351  
 egy where the methods are ranked by their mean 352  
 accuracy in each experimental configuration (each 353  
 column and grouping). Our method is seen to lead 354  
 the overall rankings. 355

For a sparsified network, the distribution of num- 356  
 ber of parameters across layers exemplified for VGG- 357  
 19 in Appendix C indicates that the depth of this 358  
 network was potentially redundant for the given 359  
 task, as most filters in the latter half of the lay- 360  
 ers have been sparsified almost completely, without 361  
 much reducing the accuracy as seen in our results. 362

To provide further insight, in Figure 2 we present 363  
 the distributions of weight magnitudes for networks 364  
 trained using a dense model as well as using GraNet 365  
 and the proposed FGGP. Although both pruning 366  
 methods reduce near-zero weights, which may have 367  
 less impact on the final predictions, our approach is 368  
 seen to be more effective in this regard. 369

### 370 4.1 Ablation Study

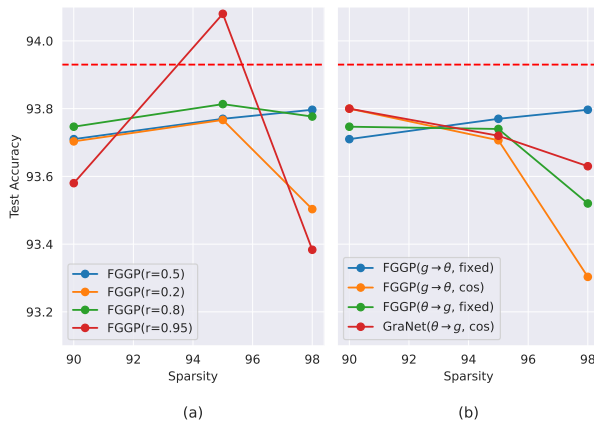
To study the effect of the proposed method compo- 371  
 nents and parameters, we conduct ablation experi- 372  
 ments for the dense-to-sparse setting with VGG-19 373  
 on CIFAR-10, repeating each experiment with three 374  
 initialization seeds and reporting the mean values 375  
 for comparison. 376

First, we evaluate the impact of the subset selec- 377  
 tion rate  $r$  by comparing results for values  $r =$  378



**Table 1.** Test accuracy of pruned VGG-19 and ResNet-50 networks on CIFAR-10 dataset, with mean $\pm$ std values from experiments with three different seeds. GraNet [15] and RigL [10] results were taken from [15], while the other results were compiled from [9, 16] for the reported settings. The bold numbers indicate the best accuracy in each given gradual pruning subcategory. The rightmost column shows the average ranking of methods within a subcategory across the six experimental settings (columns) (the stars indicate averages from VGG-19 only).

Sparsity target (N%)	VGG-19				ResNet-50			Rank
	90%	95%	98%		90%	95%	98%	
Single-shot (0% $\rightarrow$ N% at init)	SNIP [8]	93.63	93.43	92.05	92.65	90.86	87.21	2.17
	GraSP [9]	93.30	93.04	92.19	92.47	91.32	88.77	2.33
	SynFlow [17]	93.35	93.45	92.24	92.49	91.22	88.82	1.50
Sparse-to-sparse (50% $\rightarrow$ N% gradual)	Deep-R [18]	90.81	89.59	86.77	-	-	-	5.00*
	SET [19]	92.46	91.73	89.18	-	-	-	4.00*
	RigL [10]	93.38 $\pm$ 0.11	93.06 $\pm$ 0.09	91.98 $\pm$ 0.09	94.45 $\pm$ 0.43	93.86 $\pm$ 0.25	93.26 $\pm$ 0.22	3.00
	GraNet [15]	<b>93.73<math>\pm</math>0.08</b>	93.66 $\pm$ 0.07	93.38 $\pm$ 0.15	94.64 $\pm$ 0.27	<b>94.38<math>\pm</math>0.28</b>	94.01 $\pm$ 0.23	1.67
	FGGP (ours)	93.68 $\pm$ 0.04	<b>93.94<math>\pm</math>0.17</b>	<b>93.63<math>\pm</math>0.15</b>	<b>94.76<math>\pm</math>0.11</b>	94.27 $\pm$ 0.38	<b>94.22<math>\pm</math>0.24</b>	<b>1.33</b>
Dense-to-sparse (0% $\rightarrow$ N% gradual)	STR [20]	93.73	93.27	92.21	92.59	91.35	88.75	4.50
	SIS [16]	<b>93.99</b>	93.31	92.16	92.81	91.69	90.11	3.67
	GMP [21]	93.59 $\pm$ 0.10	93.58 $\pm$ 0.07	93.52 $\pm$ 0.03	94.34 $\pm$ 0.09	94.52 $\pm$ 0.08	94.19 $\pm$ 0.04	3.17
	GraNet [15]	93.80 $\pm$ 0.10	93.72 $\pm$ 0.11	93.63 $\pm$ 0.08	94.49 $\pm$ 0.08	94.44 $\pm$ 0.01	<b>94.34<math>\pm</math>0.17</b>	2.00
	FGGP (ours)	93.71 $\pm$ 0.15	<b>93.77<math>\pm</math>0.25</b>	<b>93.80<math>\pm</math>0.02</b>	<b>94.78<math>\pm</math>0.19</b>	<b>94.64<math>\pm</math>0.38</b>	94.33 $\pm$ 0.42	<b>1.67</b>
Dense (0% upperbound)	93.93 $\pm$ 0.35				94.73 $\pm$ 0.06			



**Figure 3.** (a) Comparison of FGGP for four different subset selection ratios  $r$ . (b) Ablations of FGGP indicated as  $(\cdot, \cdot)$  where the first indicated the pruning criteria order ( $g \rightarrow \theta$ : gradient-first &  $\theta \rightarrow g$ : magnitude-first) and the second the subset selection strategy (with the rate *fixed* or varying as *cosine annealed*). GraNet’s proposed method choices are also indicated in the same notation for clarity. Note that  $FGGP(g \rightarrow \theta, fixed)$  is our proposed mechanism with gradient-first fixed-rate subset selection. Experiments are reported for dense-to-sparse pruning of VGG-19 for target sparsities of  $\{90,95,98\}\%$ .

379  $\{0.20, 0.50, 0.80, 0.95\}$ . Note that in the extreme  
380 case of  $r = 1$ , the second stage would be the  
381 sole determining criterion — effectively reducing  
382 the method similar to gradual magnitude pruning  
383 (GMP). The results are depicted in Figure 3 (a).  
384 Although settings differ in their performance (with  
385  $r = 0.95$  even surpassing the upper bound at 95%  
386 sparsity), overall  $r = 0.5$  and  $r = 0.8$  consistently  
387 perform well. We chose  $r = 0.5$  for all the experi-  
388 ments given its superior trend with higher sparsity,  
389 as a primary target of network compression.

390 Second, we ablate different parts of our method

391 FGGP, mimicking the GraNet behaviour to assess  
392 the components with positive contribution. For the  
393 subset selection, we use our fixed rate as well as the  
394 varying (cosine annealing) rate change from GraNet.  
395 We also test the change of order from gradient-first  
396 to magnitude-first, as well as the combinations of  
397 this with the rate choice above. The results are seen  
398 in Figure 3(b). As the sparsity level gets higher,  
399 our proposed method FGGP with the gradient-first  
400 and fixed-rate settings are seen to achieve the best  
401 performance. Note that for all these methods the  
402 cubic scheduling function already reduces  $N_p$  over  
403 time, so the results may indicate that an additional  
404 reduction of the subset selection rate is redundant  
405 and detrimental.

## 5 Discussion and Conclusion 406

407 In this paper, we consider both gradient and weight  
408 magnitudes in the unstructured gradual pruning  
409 of parameters to sparsify networks. We herein argue  
410 that the criteria and the mechanism (the order,  
411 thresholds, etc) used in the prioritization of pruned  
412 parameters are essential. We propose to use a fixed-  
413 ratio of parameter gradient magnitudes as a first  
414 decision criteria for pruning, and experimentally  
415 validate this in a variety of settings. Pruning has  
416 the potential to substantially reduce computational  
417 costs in deep learning, thereby contributing to lower  
418 energy consumption and carbon emissions without  
419 sacrificing ultimate performance. Lower energy use  
420 can enable novel end-user experience, such as rendering  
421 IoT devices feasible. Having smaller models  
422 with similar performances could allow the deployment  
423 of complex models on smaller hardware and  
424 of very large/deep models that would otherwise be  
425 infeasible.

426 Note that Liu et al. [15] explain their method  
427 GraNet as being able to regenerate/add ( $pN_t$ ) con-  
428 nections/parameters to a network during the prun-  
429 ing operations. However, following their descriptions  
430 and pseudocode it becomes evident that they apply  
431 a two-step strategy which first ranks the parameters  
432 by their magnitude  $|\theta|$  to select the subset of small-  
433 est  $N_p + pN_t$  for further ranking gradient magnitude  
434  $|g_i|$  to select the smallest  $N_p$  of them to prune, as  
435 demonstrated in Figure 1(b). We believe this dem-  
436 ystification of such state-of-the-art is a further  
437 minor contribution of our work, as it also enabled  
438 us herein to technically compare our methods and  
439 experimentally design comparative ablation experi-  
440 ments. Note that at any given training step, some  
441 parameters may momentarily be zero (e.g., while  
442 changing sign) despite having nonzero gradients. If  
443 the number of such parameters exceeds  $N_p + pN_t$ ,  
444 the second step of GraNet may inadvertently prune  
445 essential parameters, especially for larger  $p$  at earlier  
446 iterations. In contrast, our approach prioritizes gra-  
447 dient magnitudes in the first ranking step, ensuring  
448 that only relatively stable parameters (those that  
449 have locally converged) are considered for pruning.  
450 Note that our above criterion is more conservative  
451 than GraNet, and it may disregard some good pa-  
452 rameter candidates. Also, if all gradients are large,  
453 changing parameters may still be considered erro-  
454 neously. Nevertheless, the results show that our  
455 strategy is superior to that of the earlier state-of-  
456 the-art. In the future, including  $|g\theta|$  and/or Hessian  
457 approximations can potentially improve the results  
458 further.

## 459 References

- 460 [1] C. White, M. Safari, R. Sukthanker, B. Ru, T.  
461 Elsken, A. Zela, D. Dey, and F. Hutter. *Neural*  
462 *Architecture Search: Insights from 1000 Papers*.  
463 2023. arXiv: [2301.08727](https://arxiv.org/abs/2301.08727).
- 464 [2] Y. LeCun, J. Denker, and S. Solla. “Optimal  
465 Brain Damage”. In: *Advances in Neural Infor-*  
466 *mation Processing Systems (NeurIPS)*. Vol. 2.  
467 1989.
- 468 [3] J. Frankle and M. Carbin. “The lottery ticket  
469 hypothesis: Finding sparse, trainable neural  
470 networks”. In: *Int Conf on Learning Representations (ICLR)*. 2019.  
471
- 472 [4] Z. Liu, M. Sun, T. Zhou, G. Huang, and T.  
473 Darrell. “Rethinking the Value of Network  
474 Pruning”. In: *International Conference on*  
475 *Learning Representations*. 2019.
- 476 [5] Y. He, X. Zhang, and J. Sun. “Channel prun-  
477 ing for accelerating very deep neural networks”.  
478 In: *IEEE Int Conf on computer vision (ICCV)*.  
479 2017, pp. 1389–1397.
- [6] Z. Zhuang, M. Tan, B. Zhuang, J. Liu, Y. Guo,  
480 Q. Wu, J. Huang, and J. Zhu. “Discrimination-  
481 aware channel pruning for deep neural net-  
482 works”. In: *Advances in neural information*  
483 *processing systems (NeurIPS)* 31 (2018).  
484
- [7] N. K. Dinsdale, M. Jenkinson, and A. I. Nam-  
485 burete. “STAMP: Simultaneous Training and  
486 Model Pruning for low data regimes in med-  
487 ical image segmentation”. In: *Medical Image*  
488 *Analysis* 81 (2022), p. 102583.  
489
- [8] N. Lee, T. Ajanthan, and P. H. Torr. “Snip:  
490 Single-shot network pruning based on connec-  
491 tion sensitivity”. In: *Int Conf on Learning*  
492 *Representations (ICLR)*. 2019.  
493
- [9] C. Wang, G. Zhang, and R. Grosse. “Picking  
494 winning tickets before training by preserving  
495 gradient flow”. In: *Int Conf on Learning Rep-*  
496 *resentations (ICLR)*. 2020.  
497
- [10] U. Evci, T. Gale, J. Menick, P. S. Castro,  
498 and E. Elsen. “Rigging the lottery: Making all  
499 tickets winners”. In: *International Conference*  
500 *on Machine Learning*. PMLR. 2020, pp. 2943–  
501 2952.  
502
- [11] C. M. Medeiros and G. A. Barreto. “A novel  
503 weight pruning method for MLP classifiers  
504 based on the MAXCORE principle”. In: *Neur-*  
505 *al Computing and Applications* 22 (2013),  
506 pp. 71–84.  
507
- [12] Y. Guo, A. Yao, and Y. Chen. “Dynamic  
508 network surgery for efficient dnns”. In: *Ad-*  
509 *vances in neural information processing sys-*  
510 *tems (NeurIPS)* 29 (2016).  
511
- [13] M. Zhu and S. Gupta. “To prune, or not to  
512 prune: exploring the efficacy of pruning for  
513 model compression”. In: *ICLR workshop*. 2018.  
514
- [14] T. Dettmers and L. Zettlemoyer. *Sparse net-*  
515 *works from scratch: Faster training without*  
516 *losing performance*. 2019. arXiv: [1907.04840](https://arxiv.org/abs/1907.04840).  
517
- [15] S. Liu, T. Chen, X. Chen, Z. Atashgahi, L. Yin,  
518 H. Kou, L. Shen, M. Pechenizkiy, Z. Wang, and  
519 D. C. Mocanu. “Sparse Training via Boosting  
520 Pruning Plasticity with Neuroregeneration”.  
521 In: *Advances in Neural Information Processing*  
522 *Systems (NeurIPS)*. 2021.  
523
- [16] S. Verma and J.-C. Pesquet. “Sparsifying net-  
524 works via subdifferential inclusion”. In: *Int*  
525 *Conf on Machine Learning (ICML)*. 2021,  
526 pp. 10542–52.  
527
- [17] H. Tanaka, D. Kunin, D. L. Yamins, and S.  
528 Ganguli. “Pruning neural networks without  
529 any data by iteratively conserving synaptic  
530 flow”. In: *Advances in neural information pro-*  
531 *cessing systems* 33 (2020), pp. 6377–6389.  
532

- 533 [18] G. Bellec, D. Kappel, W. Maass, and R. Leg-  
534 enstein. “Deep rewiring: Training very sparse  
535 deep networks”. In: *Int Conf on Learning Rep-*  
536 *resentations (ICLR)*. 2018.
- 537 [19] D. C. Mocanu, E. Mocanu, P. Stone, P. H.  
538 Nguyen, M. Gibescu, and A. Liotta. “Scal-  
539 able training of artificial neural networks with  
540 adaptive sparse connectivity inspired by net-  
541 work science”. In: *Nature communications* 9.1  
542 (2018), p. 2383.
- 543 [20] A. Kusupati, V. Ramanujan, R. Somani, M.  
544 Wortsman, P. Jain, S. Kakade, and A. Farhadi.  
545 “Soft threshold weight reparameterization for  
546 learnable sparsity”. In: *Int Conf on Machine*  
547 *Learning (ICML)*. PMLR. 2020, pp. 5544–55.
- 548 [21] T. Gale, E. Elsen, and S. Hooker. *The state of*  
549 *sparsity in deep neural networks*. 2019. arXiv:  
550 [1902.09574](https://arxiv.org/abs/1902.09574).

## A Experiment details

Our proposed method FGGP is implemented in Pytorch 2.0. We use cross-entropy loss for classification and the Stochastic Gradient Descent (SGD) for optimization in all experiments. Training hyperparameters are tabulated in Table A.1. We set  $\Delta t = 1000$  iterations. Pruning is stopped after 80% of the targeted epochs (i.e.,  $t_{fin}$  is set to the iteration number forecast for the 148th epoch), hence leaving the remaining 20% of training to fine-tune the model parameters without any disruption from architectural changes – a strategy also common in other gradual pruning approaches. CIFAR-10 dataset contains 10 different classes representing airplanes, cars, birds, cats, deer, dogs, frogs, horses, ships, and trucks. It consists of 60k  $32 \times 32$  color images with 6k images for each class, with a split of 50k training and 10k testing images. We use data augmentation with random crops with padding of 4 and horizontal flips.

## B Erdős–Rényi initialization

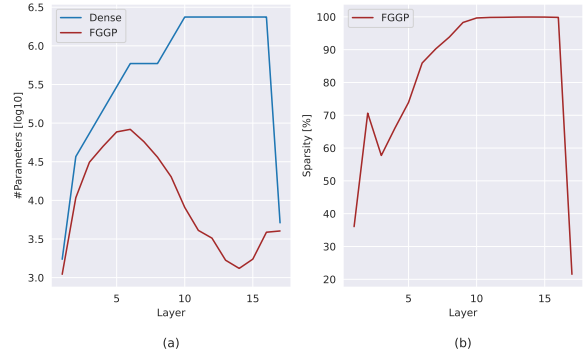
Erdős–Rényi [19] is a strategy for initializing the parameters in fully-connected layers. Erdős–Rényi Kernel (ERK) [10] offers an extension of this to convolutional layers. Such initialization was shown to perform superior to networks initialized randomly [10]. ERK [10] determines a factor  $f_l$  to scale the initialization of the parameters in the convolutional kernel  $l$  with width  $w_l$  and height  $h_l$  as:

$$f_l = 1 - \frac{n_{l-1} + n_l + w_l + h_l}{n_{l-1} \times n_l \times w_l \times h_l}, \quad (4)$$

where  $n_l$  is the total number of parameters in that convolutional kernel.

## C Sparsity of pruned network

To give an insight of where the parameters are pruned the most and how the pruned networks look like, in this section we show the sparsity distribution of networks after applying FGGP. The number of parameters of a dense and an FGGP pruned VGG-19



**Figure C.1.** (a) Number of parameters per layer in a dense and FGGP-pruned VGG-19 network, shown in logarithmic scale. (b) Sparsity of each layer after pruning. The results are shown for a sample experiment. Note that the final layer is fully-connected while the others are convolutional.

are shown in Figure C.1(a), with the resulting layer sparsity plotted in Figure C.1(b). As can be seen, the second half of the network is almost completely sparsified, likely leaving one or a few unit filters to simply forward propagate the information extracted by the initial convolutional layers for the final the prediction in the last fully-connected layer (which hence could not sparsify much). This observation suggests that the depth of VGG-19 may be highly redundant for the task of CIFAR-10 classification, which could potentially be tackled with a half-the-depth network. Future studies shall investigate this aspect, potentially using pruning as an automatic tool to determine optimal network shape and size of traditional handcrafted architectures.

**Table A.1.** Training hyperparameters.

Data	CIFAR-10
Model	VGG-19 / ResNet-50
Epochs	160
Batch Size	128
LR	0.1
LR Decay Epoch	[80, 120]
LR Decay Factor	0.1
Weight Decay (L2)	0.0005
Momentum	0.9