002

## Hyperflows: Pruning Reveals the Importance of Weights

#### **Anonymous Authors**<sup>1</sup>

#### Abstract

Network pruning is used to reduce inference latency and power consumption in large neural networks. However, most existing methods struggle to accurately assess the importance of individual weights due to their inherent interrelatedness, leading to poor performance, especially at extreme sparsity levels. We introduce Hyperflows, a dynamic pruning approach that estimates each weight's importance by observing the network's gradient response to the weight's removal. A global pressure term continuously drives all weights toward pruning, with those critical for accuracy being automatically regrown based on their flow, the aggregated gradient signal when they are absent. We explore the relationship between final sparsity and pressure, deriving power-law equations similar to those found in neural scaling laws. Empirically, we demonstrate state-of-the-art results with ResNet-50 and VGG-19 on CIFAR-10 and CIFAR-100.

#### 1. Introduction

Overparameterization has become the norm in modern deep learning to achieve state-of-the-art performance (Neyshabur et al., 2019; Allen-Zhu et al., 2019; Li et al., 2018). Despite clear benefits for training, this practice also increases computational and memory costs, complicating deployment on resource-constrained devices such as edge hardware, IoT platforms, and autonomous robots (Shi et al., 2016; Li et al., 2019). Recent theoretical and empirical findings suggest that sparse subnetworks extracted from large dense models can match or exceed the accuracy of their dense counterparts (Frankle & Carbin, 2019; Zhou et al., 2019; Ma et al., 2021; Lee et al., 2019; De Jorge et al., 2021; Cho et al., 2023; Yite et al., 2023; Frantar et al., 2024; Wang et al., 2023) and even outperform smaller dense models of equal size (Ramanujan et al., 2020; Li et al., 2020; Zhu & Gupta, 2018). These results have created interest in network pruning as a strategy to identify minimal, high-performing subnetworks.

Pruning has a rich history (LeCun et al., 1989; Mozer & Smolensky, 1988; Thimm & Hoppe, 1995) and continues to prove valuable for real-time applications (Han et al., 2016; Jongsoo et al., 2017; Wang et al., 2019). Recent methods have significantly advanced the field by resorting to a variety of strategies, from heuristics, gradient methods and Hessian-based criteria (Han et al., 2015; 2016; LeCun et al., 1992; Singh & Alistarh, 2020; Bellec et al., 2018) to dynamic pruning approaches (Liu et al., 2020; Cho et al., 2023; Savarese et al., 2020; Kusupati et al., 2020; Wortsman et al., 2019) or combinations thereof. However, the strong interdependencies between weights remain a challenge (Jin et al., 2020; Templeton et al., 2024; Lee et al., 2019; De Jorge et al., 2021; Louizos et al., 2017), as they complicate the task of determining each weight's absolute importance.

Given this gap, we ask: *Can we rigorously quantify a weight's importance for model accuracy, while accounting for the inherent interrelatedness among neural network weights?* 

Inspired by the well-known insight that the value of something is not truly known until it is lost, we introduce Hyperflows, a dynamic pruning method which determines weight importance by first removing it. Each weight  $\theta_i$  will be pruned if its associated flow parameter  $t_i$  is negative. The value of  $t_i$  will follow the direction of  $|\theta_i|$ , as their gradients are strongly correlated, while a global *pressure* term  $L_{-\infty}$ will push all t values towards  $-\infty$ . When an important weight  $\theta_i$  is pruned, the network will attempt to increase  $t_i$ . If the aggregated gradient over multiple iterations of  $t_i$ , which we call *flow*, is larger than the aggregated pressure, the removed weight will be restored, otherwise it will remain pruned. By allowing this process to happen concurrently on all weights multiple times, the network's topology becomes noisy, disentangling the restoration process from a specific configuration and therefore providing a good approximation for weight importance.

We analyze the relationships between sparsity and pres-

<sup>&</sup>lt;sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

sure, obtaining power-law dependencies similar to those
of known scaling laws in neural networks (Hestness et al.,
2017; Kaplan et al., 2020; Henighan et al., 2020; Rosenfeld
et al., 2020; Gordon et al., 2021; Hernandez et al., 2021;
Zhai et al., 2021; Hoffmann et al., 2022).

Our method is able to improve over baseline performance for Resnet-50 and VGG-19 on CIFAR-10 for sparsities up to 98%, exceeding current state of the art method and proving the ability of Hyperflows to preserve accuracy.

065 Summarizing, our key contributions are:

- We introduce *Hyperflows*, a dynamic pruning method aiming to quantify weight importance by developing the notions of *flow* and *pressure*.
  - We set a new state-of-the-art benchmark, achieving better accuracy than existing methods in empirical validation, across several networks and datasets.
  - We explore the mathematical relationships between pressure and sparsity, finding power-laws similar to those in neural scaling laws.

### 079 **2. Related work**

066

067

068

069

070

071

074

075

076

077

078

Research on neural network pruning has a relatively old 081 history, with some methods going back decades and lay-082 ing the groundwork for modern approaches. Early tech-083 niques, such as (LeCun et al., 1989) and (LeCun et al., 1992), utilized Hessian-based techniques and Taylor expan-085 sions to identify and remove unimportant specific weights, while Mozer & Smolensky (1988) employed derivatives to 087 remove whole units, an early form of structured pruning. 088 These initial studies demonstrated the feasibility of reduc-089 ing network complexity without significantly compromising 090 performance. An influential overview (Thimm & Hoppe, 091 1995) concluded that magnitude pruning was particularly ef-092 fective, a paradigm that since then has been widely adopted 093 (Han et al., 2016; Frankle & Carbin, 2019; Zhou et al., 2019; 094 Evci et al., 2020; Kusupati et al., 2020; Han et al., 2015). 095

096 The existence of highly effective subnetworks builds 097 upon these foundational theoretical studies, with the Lottery 098 Ticket Hypothesis (Frankle & Carbin, 2019) being a good 099 example. Magnitude pruning is used to demonstrate that 100 there exists a mask which, if applied at the start of training, produces a sparse subnetwork capable of matching the performance of the original dense network after training, if the initialization is kept unmodified. Subsequent research 104 has further validated this concept by showing that these 105 subnetworks produced by masks, even without any training, 106 achieve significantly higher accuracy than random chance (Zhou et al., 2019), reaching up to 80% accuracy on MNIST. 108 Moreover, training these masks instead of the actual weight 109

values can result in performance comparable to the original network (Ramanujan et al., 2020; Zhou et al., 2019), suggesting that neural network training can occur through mechanisms different from weight updates, including the masking of randomly initialized weights. Other studies have attempted to identify the most trainable subnetworks at initialization. Lee et al. (2019) use gradient magnitudes as a way to identify trainable weights, while Savarese et al. (2020) employ  $L_0$  regularization along with a sigmoid function that gradually transitions into a step function during training, enabling continuous sparsification. These findings indicate that the specific values and even the existence of certain weights may be less critical than previously believed.

Dynamic pruning differs from classical heuristics by finding sparse networks during training and allowing the model to adjust itself. Some methods use learnable parameters, e.g. Kusupati et al. (2020) train magnitude thresholds for each layer in the network to determine which weights will be pruned. Other works, like that of Cho et al. (2023), do not have any learnable parameters, learning instead a weight distribution whose shape will determine which and how many weights are pruned. Yet another class of  $L_0$  regularization techniques (Savarese et al., 2020; Louizos et al., 2018) try to maximize the number of removed weights. Hyperflows aligns with the dynamic pruning paradigm by enabling continuous pruning of weights based on dynamically updated parameters. However, unlike most methods that rely on instantaneous gradients or fixed thresholds, Hyperflows introduces a novel mechanism that assesses each weight's importance through an aggregate gradient signal over multiple iterations.

Pruning based on gradient values is another prominent approach, often overlapping with dynamic methods, which enables the assessment of weight properties in relation to the loss function. Lee et al. (2019) and De Jorge et al. (2021) assess the trainability of subnetworks by analyzing initial gradient magnitudes relative to the loss function. AutoPrune (Xiao et al., 2019) introduces handcrafted gradients that influence training, while Dynamic Pruning with Feedback (Lin et al., 2020) uses gradients during backpropagation to recover pruned weights with high trainability, preserving accuracy. Evci et al. (2020) use gradient and weight magnitudes to determine which weights to prune and to regrow. Liu et al. (2022) build upon these concepts, by employing a zero-cost neuroregeneration scheme, which prunes and regrows the same number of weights, effectively keeping the sparsity constant while growing accuracy. Our method uses gradients magnitudes to approximate the importance of a weight when it is pruned, which proves to be effective, since gradients are correlated with the loss of features induced by the pruning that weight. Hyperflows distinguishes itself from other methods by utilizing gradient magnitudes to evaluate the importance of weights after the moment of

their pruning. Instead of predefining which weights are
(un)important based solely on instantaneous gradients or
single-stage evaluations, Hyperflows identifies a weight's
significance based on the aggregated impact its removal has
on the network's performance.

## **3. Hyperflows**

115

134

135

136

137

138 139 140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

118 The main idea of our method is to assess the contribution of 119 individual weights in a neural network by first pruning them 120 and evaluating the resulting impact on performance across 121 various network topologies. This is achieved by introducing 122 a learnable parameter  $t_i$  near each weight  $\theta_i$  which decides 123 if  $\theta_i$  is active or pruned. When a necessary weight is pruned, the gradient of its associated  $t_i$ ,  $\frac{\partial \mathcal{L}}{\partial t_i}$ , termed weight *flow*, increases  $t_i$ , which regrows the weight. Flow is strongly 124 125 126 correlated with the decrease in performance when the weight 127 is removed, thereby serving as a good approximation for its 128 importance. An  $L_{-\infty}$  penalty, defined below in (7), is used 129 to push t values towards pruning and control the overall 130 sparsity. As compression occurs, the remaining weights 131 will have increased importance, by capturing the features 132 lost from the permanently pruned weights, leading to larger 133 flows. We analyze this effect in Appendix B.

#### 3.1. Preliminaries

Consider a neural network defined as a function:

$$f: \mathcal{X} \times \theta \to \mathcal{Y},$$

where  $\mathcal{X}$  is the input space,  $\mathcal{Y}$  is the output space, and  $\theta \subseteq \mathbb{R}^d$  denotes the weight vector.

Given a training set  $\{(x_j, y_j)\}_{j=1}^J$ , learning the parameters  $\theta$  amounts to minimizing a loss function:

$$\min_{\theta} \sum_{j=1}^{J} \ell(f(x_j, \theta), y_j),$$

so that  $f(x_j, \theta)$  aligns with  $y_j$ .

We define the topology of the neural network  $\mathcal{T}$  as a binary vector  $\mathcal{T} \in \{0,1\}^d$  where  $\mathcal{T}^i \in \{0,1\}$  represents whether weight  $\theta_i$  is pruned or not. We denote a family of topologies as  $\{\mathcal{T}_k\}_{k=1}^K$ , with K its cardinality. Thus, the loss of a network with topology  $\mathcal{T}$  is:

$$\mathcal{L}(\mathcal{T}) = \sum_{j=1}^{J} \ell \big( f(x_j, \theta \odot \mathcal{T}), y_j \big),$$

160 where  $\odot$  is the Hadamard product. Note that  $\mathcal{L}(\mathcal{T})$  depends 161 on  $\theta$ .

<sup>162</sup> <sup>163</sup> For each parameter  $\theta_i$ , we introduce a learnable scaler  $t_i$ to which we refer as *flow parameter*. We denote with *t* the vector of flow parameters. Vector t is used to generate the topology  $\mathcal{T}$  with  $\mathcal{T}^i = H(t_i)$ , where:

$$H(t_i) = \begin{cases} 1 & \text{if } t_i > 0, \\ 0 & \text{if } t_i \le 0. \end{cases}$$

Thus, if  $t_i > 0$  then  $\theta_i$  is active, otherwise  $(t_i \leq 0)$ ,  $\theta_i$  is pruned. We denote with  $\mathcal{T} \setminus {\theta_i}$  a topology  $\mathcal{T}$  for which we set  $\mathcal{T}^i = 0$  and define the change in  $\mathcal{L}(\mathcal{T})$  when  $\theta_i$  is removed as:

$$\Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_i\}) = \mathcal{L}(\mathcal{T} \setminus \{\theta_i\}) - \mathcal{L}(\mathcal{T}).$$

We use a global penalty term  $L_{-\infty}$  to push all  $t_i$  values towards  $-\infty$ , which we discuss in detail in Section 3.2. Our goal is to find a topology  $\mathcal{T}_f$  and set of weights  $\theta$  such that the following loss is minimal:

$$\mathcal{J}(\mathcal{T}) = \mathcal{L}(\mathcal{T}) + L_{-\infty}(t). \tag{1}$$

#### 3.2. Weight Flow

Since the optimal topology  $\mathcal{T}^*$  is initially unknown, any metric for the importance of  $\theta_i$  measured on the initial topology  $\mathcal{T}_0$  might not be relevant for  $\mathcal{T}^*$ . For this reason, weight importance is evaluated multiple times during training. We present an importance metric, called *flow*, tied to a specific topology  $\mathcal{T}$ , and then extend it to a family of topologies  $\{\mathcal{T}_k\}_{k=1}^K$ . We begin by defining the gradient:

$$\mathcal{G}(\theta_i, \mathcal{T}) = \frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_i}, \forall t_i \in \mathbb{R}.$$
 (2)

Importantly, the sign of the gradient  $\mathcal{G}(\theta_i, \mathcal{T})$  always follows the direction of  $|\theta_i|$  (proof in Appendix A.2). If  $\theta_i$  increases or decreases in magnitude, then  $t_i$  will correspondingly increase or decrease.

In our method,  $\mathcal{G}(\theta_i, \mathcal{T})$  takes two different meanings based on whether  $t_i > 0$  or  $t_i \leq 0$ . We first define the meaning in the case  $t_i \leq 0$  as *flow* for one topology:

$$\mathcal{F}(\theta_i, \mathcal{T}) = \begin{cases} \mathcal{G}(\theta_i, \mathcal{T}) & t_i \le 0, \\ 0 & t_i > 0. \end{cases}$$
(3)

When  $\theta_i$  is pruned  $(t_i \leq 0)$  from topology  $\mathcal{T}$ , a corresponding  $\Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_i\})$  will occur. If  $\Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_i\}) > 0$ , then the pruning of  $\theta_i$  leads to an increase in the loss, and increasing  $|\theta_i|$  from 0 back to the original will reduce the loss again. Otherwise, if  $\Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_i\}) \leq 0$ , increasing  $|\theta_i|$  will not reduce the loss. Therefore, any parameter whose gradient depends on changes in  $|\theta_i|$  will have its value increased only if the pruned weight was important. Overall,  $\mathcal{F}(\theta_i, \mathcal{T})$  will create large positive changes in  $t_i$  for important pruned weights, regrowing them (proof in Appendix A.1).

165 For the case  $t_i > 0$ , we define:

$$\mathcal{M}(\theta_i, \mathcal{T}) = \begin{cases} 0 & t_i \le 0, \\ \mathcal{G}(\theta_i, \mathcal{T}) & t_i > 0. \end{cases}$$
(4)

170 Minimizing the loss function does not inherently correlate 171 with weight magnitude increases. Therefore  $\mathcal{M}(\theta_i, \mathcal{T}_k)$  will 172 represent changes in magnitude over training. Pruning will 173 be encouraged for weights whose magnitude decreases, and 174 be resisted for those whose magnitudes are increasing (proof 175 in Appendix A.2).

176 All the functions involved in backpropagation must be differ-177 entiable, but *H* is not. Since  $\mathcal{F}(\theta_i, \mathcal{T})$  should only depend 178 on the importance of  $\theta_i$  and not the value of  $t_i$ , we choose a 179 straight-through estimator for the gradient  $\frac{\partial H}{\partial t_i} = 1$ , which 180 is also used for  $\mathcal{M}(\theta_i, \mathcal{T})$  and therefore for the overall gra-181 dient  $\mathcal{G}(\theta_i, \mathcal{T})$ .

In practice, we analyze weight behaviour over several
topologies. Extending equations (3) and (4) to a family
of topologies we obtain the aggregated flow and the average
change of the weight magnitude respectively:

$$\mathcal{F}(\theta_i, \{\mathcal{T}_k\}_{k=1}^K) = \frac{1}{K} \cdot \sum_{k=1}^K \mathcal{F}(\theta_i, \mathcal{T}_k),$$
(5)

$$\mathcal{M}(\theta_i, \{\mathcal{T}_k\}_{k=1}^K) = \frac{1}{K} \cdot \sum_{k=1}^K \mathcal{M}(\theta_i, \mathcal{T}_k).$$
(6)

To drive t values towards  $-\infty$ , we employ an " $L_{-\infty}$ " loss called *pressure*, formulated as:

196 197

198

199

200

204

205

206

208 209 210

187 188

189 190

167 168

169

$$L_{-\infty}(t) = \frac{1}{d} \cdot \gamma \cdot \sum_{i=1}^{d} t_i, \tag{7}$$

where  $\gamma$  is a scalar used to control sparsity and d the number of weights in the network. From this point forward, any reference about an increase or decrease in pressure will refer to an increase or decrease in  $\gamma$ .

It is important to explore how  $\mathcal{L}(\mathcal{T})$  and  $L_{-\infty}(t)$  interact with t values in (1). For multiple iterations R of training, we get:

$$\sum_{r=1}^{R} \frac{\partial (\mathcal{L}(\mathcal{T}_r) + L_{-\infty}(t))}{\partial t_i} = \sum_{r=1}^{R} (\mathcal{G}(\theta_i, \mathcal{T}_r) + \frac{\gamma}{d}), \quad (8)$$

211 where  $\mathcal{T}_r$  is the topology at iteration r. At each iteration, 212 a weight can either be pruned or active, therefore, we can 213 partition a weight's state during training between pruned 214 and active stages. A stage  $S_f, f \in \{1, \ldots, F\}$  is a series of 215 consecutive iterations for which our weight is in the same 216 state. Each  $S_f$  has a duration of  $D_f$ , starting at iteration 217  $s_f$  and ending at  $e_f$ . We denote by  $S_f^+$  the stages when a 218 weight is present and  $S_f^-$  those when the weight is pruned. We define gradients taking place during a stage  $S_f^+$  and  $S_f^-$  as:

$$\nabla S_f^+ = \sum_{r=s_f}^{e_f} \left( \mathcal{M}(\theta_i, \mathcal{T}_r) + \frac{\gamma}{d} \right), \tag{9}$$

$$\nabla S_f^- = \sum_{r=s_f}^{e_f} \left( \mathcal{F}(\theta_i, \mathcal{T}_r) + \frac{\gamma}{d} \right).$$
(10)

Note that it is not possible to have two consecutive stages with the same weight state and all weights are present at the start of training. We partition the set of all stages into two  $\{S_1^+, S_3^+, ...\}$  and  $\{S_2^-, S_4^-, ...\}$ . We refer to the transition between stages as *implicit regrowth*. We do not know in which partition the final stage  $S_F$  will be until training ends. Equation (8) becomes:

$$\sum_{r=1}^{R} (\mathcal{G}(\theta_i, \mathcal{T}_r) + \frac{\gamma}{d}) = \nabla S_1^+ + \nabla S_2^- + \dots + \nabla S_F^{\{+\text{or}-\}}.$$
 (11)



Figure 1. A weight's state can be partitioned into pruned  $(t_i > 0)$  and active  $(t_i \le 0)$  stages. We represent with blue arrows the flow, which appears only for pruned stages, and with red lines the pressure, which appears for all stages.

The partition in stages is illustrated in Figure 1. Analyzing what happens at an individual level for each  $\nabla S_i$ , by rearranging (9) and (10), we get:

$$\nabla S_f^+ = D_f \cdot \left( \mathcal{M}(\theta_i, \{\mathcal{T}_r\}_{r=s_f}^{e_f}) + \frac{\gamma}{d} \right), \qquad (12)$$

$$\nabla S_f^- = D_f \cdot \left( \mathcal{F}(\theta_i, \{\mathcal{T}_r\}_{r=s_f}^{e_f}) + \frac{\gamma}{d} \right).$$
(13)

For  $S_f^-$ , if  $\mathcal{F}(\theta_i, \{\mathcal{T}_r\}_{r=s_f}^{e_f}) + \frac{\gamma}{d} > 0$ , then the overall change in  $t_i$  will be negative, keeping the weight pruned. Otherwise, if  $\mathcal{F}(\theta_i, \{\mathcal{T}_r\}_{r=s_f}^{e_f}) + \frac{\gamma}{d} < 0$  then  $t_i$  will be increased. In other words, all weights will be pushed towards pruning by the pressure and regrown if the flow is greater than the pressure.

Weights can be deemed unimportant in  $T_r$  and become relevant later in  $T_{r+q}$ . However, if a weight is pruned at  $T_r$ , the

pressure will continue to push its t value towards  $-\infty$  for q+1 iterations, making it difficult to regrow. To control this effect, in practice we apply the loss only on t values which are above a certain threshold T.

$$\widehat{L}_{-\infty}(t) = \gamma \cdot \sum_{i=1}^{d} t_i \cdot H(T - t_i).$$
(14)

#### 3.3. Neural pruning laws

222

223

224

225

227

228

229 230

231

233

234

235

264

265

266

267

269

270

271

272

273

274

We investigate how the pruning pressure scaler  $\gamma$ , the number of training epochs, and the network architecture shape the evolution of sparsity. These insights lay the foundation for our  $\gamma$  scheduler, introduced later, that can reach a target sparsity in any desired training time.

(0) Sparsity Convergence for a Fixed  $\gamma$ . As sparsity in-236 creases, the overall flow of the weights will become larger. 237 We ask the following question: Given a fixed  $\gamma$ , will the 238 network converge to a final sparsity s? Moreover, does this 239 mapping from  $\gamma$  to s follow any relationship? In Figure 3, 240 we test the existence of convergence empirically by run-241 ning LeNet-300 on MNIST and ResNet-50 on Cifar-10. We 242 allow each network to train for 300 to 1000 epochs with 243 a constant  $\gamma$  pressure and observe the results. We do this 244 with two different optimizers for t values, SGD and Adam. 245 Our findings suggest that there is no one curve that fits the 246 decrease in parameters for both optimizers, but the final 247 convergence point is the same regardless of the optimizer 248 used. 249



Figure 2. Convergence for fixed  $\gamma$ . We can observe that for each case there is a certain point at which weights are not pruned anymore or offer an extremely high amount of resistance

An important observation is that the final convergence point  $s_c$  is influenced by the  $\theta$  learning rate  $\eta$ . If  $\eta$  is high, convergence happens in a larger number of epochs (1000 in our experiments), at a higher sparsity. If  $\eta$  is low, convergence happens sooner 300 epochs, to a lower sparsity.

(1) Relationship Between  $\gamma$  and Final Sparsity. Assuming that all networks have a sparsity they converge to for a fixed  $\gamma$ , is there a relationship between  $\gamma$  and its associated final sparsity? Can we predict for a new  $\gamma$  the final sparsity a network will converge to? We modify the previous experiment, to run the networks 300 epochs for several values of  $\gamma$  between  $2^{-15}$  and  $2^{10}$ . Our empirical results suggest a power-law relationship:

$$\ln(s) = \ln(c) - \alpha_0 \cdot \ln(\gamma) - \alpha_1 \cdot (\ln(\gamma))^2, \qquad (15)$$

where constants, c,  $\alpha_0$ ,  $\alpha_1$  depend on dataset and network architecture.



Figure 3. Relation between  $\gamma$  and final sparsity, showing several curves that can be fitted by our power-law formula. Notice how different optimzers converge to the same points. One particularity of SGD is that for low values of pressure the networks takes longer to converge, which is why a few outliers appear in the top part diagram.

#### 3.4. Pressure Scheduler

Our findings from Section 3.3 suggest that for any sparsity we desire, there will be a certain fixed  $\gamma$  which produces that sparsity after a fixed number of epochs. However, in practical applications, this  $\gamma$  cannot be known from the start, since it would require running the method several times to find out the power-law curve. To solve this issue, we propose a dynamic scheduler that adjusts  $\gamma$  at each epoch, driving the network along a desired sparsity trajectory.

The goal of our scheduler, given a function that maps each iteration to a sparsity  $f(e) : [0, R] \rightarrow [0, 100]$ , is to adjust  $\gamma$  such that the sparsity of the network after e epochs s(e) will be equal to the desired sparsity s(e) = f(e). Since epochs increase linearly while the relationship between  $\gamma$  and sparsity is non-linear, we need our adjustments to  $\gamma$  at each epoch to also be non-linear.

We choose  $\gamma = p^{\alpha}$  as our  $\gamma$  function, where p is adjusted as described in Algorithm 1,  $\alpha$  is chosen as a hyperparameter

275 and u is a constant. We find values of  $\alpha \in [1.5, 2.0]$  to 276 be well suited for both stable and accurate pruning. We 277 use inertia terms  $p_+$  and  $p_-$  to account for the need of 278 potentially larger changes in  $\gamma$  for small  $\alpha$  values. Our 279 scheduler is able to reach the desired sparsity within a 10%280 margin of error. Ideally, the training time should be infinite and the changes in  $\gamma$  as small as possible, to allow for more 281 282 controlled pruning. In practice, we find optimal training 283 time for pruning to be somewhere between T/2 and 2T, 284 where T is the original training time needed for the network 285 to converge.

Algorithm 1 Pressure Scheduler

1: **Input:** Epoch *e*, sparsity curve f(e), current sparsity s(e), update *u* (constant),  $\alpha$ 

2: **Internals:** positive inertia  $p_+$ , negative inertia  $p_-$ , both initialized with 0.

3: if f(e) < s then

4:  $p \leftarrow p + u + p_+$ 294 5:  $p_+ \leftarrow p_+ + \frac{u}{4}$ 296 6:  $p_- \leftarrow 0$ 297 7: else if f(e) > s then

8:  $p \leftarrow p - u - p_-$ 

9:  $p_- \leftarrow p_- + \frac{u}{4}$ 

10:  $p_+ \leftarrow 0$ 

11: end if 12: Return: pressure  $\gamma = p^{\alpha}$ 

301 302

303 304

306

299

300

286

287

288

289

290

291

292

293

#### 3.5. Regrowth Stage

One of the main features of our method is the noise created 307 by removal and regrowth of weights, which leads to the 308 disentangling of weight from specific topology. However, 309 this noise is harmful for convergence. For this reason, we 310 introduce a regrowth-only stage at the end, whose purpose 311 is to allow the weights to converge as well as to stabilize the 312 network topology. Specifically, we eliminate the regulariza-313 tion term, setting  $\gamma$  to 0 and therefore allowing only weight 314 regrowth. In order to limit the number of regrown weights 315 and add only the most important lost weights, we introduce 316 a decay of  $d^e$  for the learning rate  $\eta_t$ , where d represents a 317 constant and e represents the epoch number. 318

319 Despite the number of parameters regrown being hard to 320 control, we can adjust the flow parameters learning rate 321  $\eta_t$ , as well as the decay d to constraint or promote the re-322 activation of weights. Although the absolute number of 323 parameters being regrown is small, the gains in accuracy are 324 significant. For example, we can gain up to 8% accuracy on 325 Imagenet dataset with a net increase in remaining parameters of 0.5%, from 4.0% to 4.5%. Generally, the regrowth 327 phase should be scheduled for between one-quarter to one-328 third of the total training time to allow adequate reactivation 329

and stabilization of essential weights.

#### **4. Experimental Results**

We conducted experiments with *Hyperflows* to demonstrate its effectiveness in achieving high sparsity levels while maintaining accuracy across various neural network architectures and datasets. Since *Hyperflows* quantifies the importance of each weight, we evaluate it for post-training pruning scenario, as weights need to hold significance within the network before starting the pruning process.

We compare *Hyperflow* with state-of-the-art pruning methods such as GraNet (Liu et al., 2022), RigL (Evci et al., 2020), GMP (Trevor Gale, 2019) and Synflow (Hidenori et al., 2020). We run GraNet and GMP individually using two setups. In the first setup, we use their reported best configurations for training networks from scratch. In the second setup, we initialize them with the same baseline network as ours and determine the optimal learning rate for this specific setting. We generally observe that their reported learning rates are also optimal for the post-training scenario. We use the same training budget of 160 epochs and keep all other configurations intact, to ensure no unintentional degradation occurs. We mark all the methods run individually with \*. For the methods we do not run, we use the results reported in (Liu et al., 2022) and (Kusupati et al., 2020).

We evaluate Hyperflows on the following combinations of networks and datasets: LeNet-300 on MNIST, ResNet-50 on CIFAR-10/100, VGG19 on CIFAR-10/100 and ResNet-50 on ImageNet-1K. Details on the training setups, architectures and datasets are summarized in Appendix D. Unless otherwise stated, all experiments were conducted three times, with results expressed as mean  $\pm$  standard deviation. The experiments were conducted on a system equipped with 3 RTX 4090 GPUs. Ablation studies are presented in Appendix B. Our findings indicate that Hyperflows consistently outperforms state-of-the-art pruning methods, achieving higher sparsity levels with comparable or better accuracy across multiple datasets and architectures.

#### 4.1. CIFAR-10 / 100

We evaluate the performance of *Hyperflows* on CIFAR-10 and CIFAR-100 using ResNet-50 and VGG-19 architectures. CIFAR-10 is a simpler benchmark with fewer classes, making it a smaller challenge compared to CIFAR-100, which has a larger number of classes and fewer images per class. This makes CIFAR-100 more prone to instability during training and a more rigorous test for pruning methods. Results are presented in Table 1.

A key feature of *Hyperflows* is the intentional introduction of noise during pruning. While this leads to performance fluctuations, it enhances resilience to pruning and supports

331	Table 1. Comparison of classification accuracy (%) on CIFAR-10 and CIFAR-100 datasets at different pruning ratios (90.0%, 95.0%)
222	98.0%). Results are reported for VGG-19 and ResNet-50 architectures using various pruning methods, including Hyperflows. Bold values
332	represent the best performance for each setting. We consider significant results to have at least 0.25% accuracy difference from the other
333	methods. Otherwise, we will report multiple best performing methods if it is the case.

Dataset	CIFAR-10			CIFAR-100			
Pruning ratio	90.0%	95.0%	98.0%	90.0%	95.0%	98.0%	
VGG-19 (Dense)		$\textbf{93.85} \pm \textbf{0.06}$			$\textbf{73.44} \pm \textbf{0.09}$		
SNIP	93.63	93.43	92.05	72.84	71.83	58.46	
GraSP	93.30	93.04	92.19	71.95	71.23	68.90	
STR	93.73	93.27	92.21	71.93	71.14	69.89	
SIS	93.99	93.31	93.16	72.06	71.85	71.17	
SynFlow	93.35	93.45	92.24	71.77	71.72	70.94	
RigL	93.38±0.11	$93.06 {\pm} 0.09$	$91.98 {\pm} 0.09$	$73.13 {\pm} 0.28$	$72.14{\pm}0.15$	$69.82{\pm}0.09$	
GMP*	$93.82\pm0.15$	$93.84\pm0.14$	$92.34\pm0.13$	$73.57\pm0.20$	$73.39\pm0.11$	$72.78\pm0.07$	
$\operatorname{GraNet}^*(s_i = 0)$	$93.87 \pm 0.05$	$93.84\pm0.16$	$\textbf{93.87} \pm \textbf{0.11}$	$74.08\pm0.10$	$73.86\pm0.04$	$\textbf{73.00} \pm \textbf{0.18}$	
Hyperflows (ours)	$\textbf{94.05} \pm \textbf{0.17}$	$\textbf{94.15} \pm \textbf{0.14}$	$\textbf{93.95} \pm \textbf{0.18}$	$\textbf{74.37} \pm \textbf{0.21}$	$\textbf{74.18} \pm \textbf{0.15}$	$\textbf{72.9} \pm \textbf{0.05}$	
ResNet-50 (Dense)		$\textbf{94.72} \pm \textbf{0.05}$			$\textbf{78.32} \pm \textbf{0.08}$		
SNIP	92.65	90.86	87.21	73.14	69.25	58.43	
GraSP	92.47	91.32	88.77	73.28	70.29	62.12	
STR	92.59	91.35	88.75	73.45	70.45	62.34	
SIS	92.81	91.69	90.11	73.81	70.62	62.75	
SynFlow	92.49	91.22	88.82	73.37	70.37	62.17	
RigL	$94.45 {\pm} 0.43$	$93.86 {\pm} 0.25$	$93.26 {\pm} 0.22$	$76.50 {\pm} 0.33$	$76.03 {\pm} 0.34$	$75.06 {\pm} 0.27$	
GMP*	$94.81\pm0.05$	$94.89\pm0.1$	$94.52\pm0.12$	$78.39\pm0.18$	$78.38 \pm 0.43$	$77.16\pm0.25$	
$\operatorname{GraNet}^*(s_i = 0)$	$94.69 \pm 0.08$	$94.44\pm0.01$	$94.34\pm0.17$	$79.09\pm0.23$	$78.71\pm0.16$	$\textbf{78.01} \pm \textbf{0.20}$	
Hyperflows (ours)	$\textbf{95.41} \pm \textbf{0.12}$	$\textbf{95.15} \pm \textbf{0.11}$	$\textbf{95.26} \pm \textbf{0.13}$	$\textbf{79.58} \pm \textbf{0.18}$	$\textbf{79.23} \pm \textbf{0.16}$	$77.7\pm0.08$	

362

363

an effective regrowth phase, achieving state-of-the-art accuracy. Further analysis of these fluctuations is provided in Appendix B.

On CIFAR-10, Hyperflows achieves above-baseline perfor-365 mance for all sparsity levels (10%, 5%, and 2%) for both 366 VGG-19 and ResNet-50, making it the only method to do 367 so. Accuracy differences between *Hyperflows* and the next 368 best method are generally within 1% or less. For VGG-369 19, Hyperflows outperforms GraNet\* at 90% sparsity by 370 0.18% and GMP\* by 0.23%. At 98% sparsity, the differ-371 ence with GMP\* increases to 1.61%. For ResNet-50, Hy-372 perflows maintains a consistent 0.7% accuracy advantage 373 over GraNet<sup>\*</sup> across all sparsity levels. Additionally, we 374 study layer-wise sparsity for ResNet-50 on CIFAR-10 at 375 extreme sparsity levels (99.74%, 99.01%, 98.13%) and ana-376 lyze weight distributions under extreme compression. We 377 observe that exreme sparsity significantly alters the distribu-378 tion of weights, results are detailed in Appendix C.1. 379

On CIFAR-100, *Hyperflows* achieves the best performance
in 4 out of 6 benchmarks. In the remaining 2 cases, it is
slightly behind GraNet\*, with differences of 0.1% and 0.3%.
Notably, GraNet\* benefits significantly from loading the

baseline network first, gaining nearly 2% accuracy points for ResNet-50 compared to their reported results. Hyperflows outperforms all methods, including GraNet<sup>\*</sup>, at 90% and 95% sparsity by 0.5% accuracy points. Other methods, such as SIS (Verma & Pesquet, 2021), STR (Kusupati et al., 2020), GraSP (Chaoqi Wang & Grosse, 2020), and SynFlow (Hidenori et al., 2020), suffer significant accuracy degradation at 98% sparsity, dropping to  $70\% \pm 2$  for VGG-19 and  $62\% \pm 0.5$  for ResNet-50.

#### 4.2. ImageNet-2012

To evaluate the scalability and effectiveness of *Hyperflows* in large-scale settings, we conducted experiments on the ImageNet-2012 dataset using the ResNet-50 architecture. The complexity and size of this dataset provide a rigorous test for pruning methods, especially at high sparsity levels. Table 2 summarizes the comparison between *Hyperflows* and other prominent pruning techniques, including GraNet (Liu et al., 2022), RigL (Evci et al., 2020), GMP (Trevor Gale, 2019), and STR (Kusupati et al., 2020). Notably, *Hyperflows* maintains higher accuracy across extreme sparsity levels, demonstrating its robustness and effective-

ness in preserving model performance despite significant 385 386 pruning. 387

388 Table 2. Comparison of Top-1 accuracy (%), number of param-389 eters (Params), and sparsity levels (%) for ResNet-50 on the 390 ImageNet-2012 dataset using various pruning methods, including GMP, DNW, RigL, GraNet, STR, and Hyperflows. Results are reported at sparsity levels of 90%, 95%, and 96.5%.

Method	Top-1 Acc (%)	Params	Sparsity (%)
ResNet-50	77.01	25.6M	0.00
GMP	73.91	2.56M	90.00
DNW	74.00	2.56M	90.00
RigL	73.00	2.56M	90.00
GraNet	74.50	2.56M	90.00
STR	74.31	2.49M	90.23
Hyperflows	74.40	2.54M	90.11
GMP	70.59	1.28M	95.00
DNW	68.30	1.28M	95.00
GraNet	72.30	1.28M	95.00
RigL	70.00	1.28M	95.00
STR	70.40	1.27M	95.03
Hyperflows	72.44	1.13M	95.58
RigL	67.20	0.90M	96.50
STR	67.22	0.88M	96.53
GraNet	70.5	0.90M	96.50
Hyperflows	70.91	0.92M	96.42

415 At 96.5% sparsity, Hyperflows achieves a Top-1 accuracy of 416 70.91%, surpassing RigL and STR by 3% accuracy, while 417 GraNet has similar performance at 70.51%. At 95% spar-418 sity, Hyperflows achieves a Top-1 accuracy of 72.14%, sig-419 nificantly outperforming GMP (70.59%), DNW (68.30%), 420 RigL (70.00%), and STR (70.40%), while GraNet falls 421 behind by 0.14%. At 90% sparsity, Hyperflows achieves 422 74.40% accuracy closely matching GraNet (74.50%) and 423 STR (74.31%). 424

Furthermore, we conducted an analysis on the weight his-425 tograms of ResNet-50 on ImageNet to study the difference 426 in weight distribution under two settings measured at same 427 sparsities during pruning and regrowth stages. The results 428 are shown in Appendix C.1, where we observe a shift in 429 the weight distributions and a decrease in the number of 430 non-zero weights during the pruning phase. 431

#### 4.3. LeNet-300

4

4 4 4

414

432

433

434 We utilize LeNet-300 for our experiments due to its sim-435 plicity and manageability, which make it an ideal choice for 436 examining various aspects of our method. By applying our 437 flow metric to eliminate non-essential weights, we progres-438 sively increase the sparsity of LeNet-300 up to 99.85% and 439

investigate which remaining weights in the first layer are crucial for classification. As shown in Figure 4, the areas essential for classification become evident, with the margin being pruned across all levels of sparsity. The dataset uses only normalization, with no additional transformations applied. Additionally, in Appendix C.2, we track the number of weight flips per iteration during training, and visualize how gradient flow contributes to reactivating important weights.



Figure 4. The pixels utilized during inference for various sparsified LeNet-300 networks. The first column shows the masks, with white pixels indicating those used at inference, while the first row presents the original, unmasked images.

#### 5. Discussion & Conclusion

We introduced Hyperflows, a theoretical framework around the idea of weight importance along with the notions of pressure and flow. We studied the relationships between flow, pressure and the final sparsity of a neural network, which we termed neural pruning laws. Based on these laws, we developed a pressure scheduler, allowing us to indirectly control sparsity, as opposed to pruning or regrowing a fixed number of weights. Furthermore, we achieved state-of-the-art results on benchmarks such as CIFAR-10 and CIFAR-100, overall demonstrating the potential of Hyperflows, from both an empirical and theoretical perspective. In future work, we aim to explore whether the constants in the pruning laws exhibit any properties that hold across multiple networks and datasets. Furthermore, it would be interesting to adapt Hyperflows to other problems and architectures, such as Reinforcement Learning and Large Language Models.

#### **Impact Statement**

This work's dynamic pruning approach can significantly reduce the computational and energy costs of deep learning models, making large networks more efficient and accessible for a wider range of applications. By achieving extreme sparsity with minimal accuracy loss, it could enable realtime or low-resource usage in domains like healthcare or
edge AI, while also lowering the overall environmental impact of machine learning.

#### References

444

445

449

450

451

452 453

454

455

456

471

472

473

- Allen-Zhu, Z., Li, Y., and Song, Z. A convergence theory for
  deep learning via over-parameterization. In *International Conference on Machine Learning*, 2019.
  - Bellec, G., Kappel, D., Maass, W., and Legenstein, R. Deep rewiring: Training very sparse deep networks. In *International Conference on Learning Representations*, 2018.
  - Chaoqi Wang, G. Z. and Grosse, R. Picking winning tickets before training by preserving gradient flow. In *International Conference on Learning Representations*, 2020.
- Cho, M., Adya, S., and Naik, D. PDP: Parameter-free differentiable pruning is all you need. In *International Conference on Neural Information Processing Systems*, 2023.
- 462 De Jorge, P., Sanyal, A., Behl, H. S., Torr, P. H., Rogez, G.,
  463 and Dokania, P. K. Progressive skeletonization: Trim464 ming more fat from a network at initialization. In *Inter-*465 *national Conference on Learning Representations*, 2021.
- Evci, U., Gale, Trevor, Menick, Jacob, Castro, Samuel, P.,
  and Elsen, E. Rigging the lottery: Making all tickets winners. In *International Conference on Machine Learning*,
  2020.
  - Frankle, J. and Carbin, M. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019.
- Frantar, E., Riquelme Ruiz, C., Houlsby, N., Alistarh, D.,
  and Evci, U. Scaling laws for sparsely-connected foundation models. In *International Conference on Learning Representations*, 2024.
- 480 Gordon, M. A., Duh, K., and Kaplan, J. Data and parameter
  481 scaling laws for neural machine translation. In *Proceed-*482 *ings of the 2021 Conference on Empirical Methods in*483 *Natural Language Processing*. Association for Computa484 tional Linguistics, 2021.
- Han, S., Pool, J., Tran, J., and Dally, W. Learning both
  weights and connections for efficient neural network. In *International Conference on Neural Information Process- ing Systems*, 2015.
- Han, S., Mao, H., and Dally, W. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. In *International Conference on Learning Representations*, 2016.

- Henighan, T., Kaplan, J., Katz, M., Chen, M., Hesse, C., Jackson, J., Jun, H., Brown, T. B., Dhariwal, P., and Gray, Scott, e. a. Scaling laws for autoregressive generative modeling. *arXiv preprint arXiv:2010.14701*, 2020.
- Hernandez, D., Kaplan, J., Henighan, T., and McCandlish, S. Scaling laws for transfer. *arXiv preprint arXiv:2102.01293*, 2021.
- Hestness, J., Narang, S., Ardalani, N., Diamos, G., Jun, H., Kianinejad, H., Patwary, M., Ali, M., Yang, Y., and Zhou, Y. Deep learning scaling is predictable, empirically. *arXiv* preprint arXiv:1712.00409, 2017.
- Hidenori, Tanakaa nd Daniel, K., Daniel, Yamins, and Ganguli, S. Pruning neural networks without any data by iteratively conserving synaptic flow. In *Advances in Neural Information Processing Systems*, 2020.
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., de Las Casas, D., Hendricks, L. A., Welbl, J., and Clark, Aidan, e. a. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- Jin, G., Yi, X., Zhang, L., Zhang, L., Schewe, S., and Huang, X. How does weight correlation affect the generalisation ability of deep neural networks? In *Advances in Neural Information Processing Systems*, 2020.
- Jongsoo, P., Sheng, L., Wei, W., Ping, T., Hai, L., Yiran, C., and Pradeep, D. Faster cnns with direct sparse convolutions and guided pruning. In *International Conference on Learning Representations*, 2017.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv* preprint arXiv:2001.08361, 2020.
- Kusupati, A., Ramanujan, V., Somani, R., Wortsman, M., Jain, P., Kakade, S., and Farhadi, A. Soft threshold weight reparameterization for learnable sparsity. In *International Conference on Machine Learning*, 2020.
- LeCun, Y., Denker, J. S., and Solla, S. A. Optimal brain damage. In *International Conference on Neural Information Processing Systems*, pp. 598–605, 1989.
- LeCun, Y., Denker, J. S., and Solla, S. A. Second order derivatives for network pruning: Optimal brain surgeon. In *International Conference on Neural Information Processing Systems*, 1992.
- Lee, J., Gao, J., Hsieh, C.-J., and Hassner, T. Snip: Singleshot network pruning based on connection sensitivity. In *International Conference on Learning Representations*, 2019.

- Li, C., Farkhoor, H., Liu, R., and Yosinski, J. Measuring the intrinsic dimension of objective landscapes. In *International Conference on Learning Representations*, 2018.
- Li, E., Zeng, L., Zhou, Z., and Chen, X. Edge ai: Ondemand accelerating deep neural network inference via
  edge computing. *IEEE Transactions on Wireless Commu- nications*, 2019.
- Li, Z., Wallace, E., Shen, S., Lin, K., Keutzer, K., Klein,
  D., and Gonzalez, J. E. Train big, then compress: Rethinking model size for efficient training and inference of
  transformers. In *International Conference on Machine Learning*, 2020.
- Lin, T., Stich, S. U., Barba, L., Dmitriev, D., and Jaggi, M.
  Dynamic model pruning with feedback. In *Proceedings* of the International Conference on Learning Representations (ICLR), 2020.
- Liu, J., Xu, Z., Shi, R., Cheung, R., and So, H. Dynamic sparse training: Find efficient sparse network from scratch with trainable masked layers. *arXiv:2005.06870*, 2020.
- Liu, S., Chen, T., Chen, X., Atashgahi, Z., Yin, L., Kou, H.,
  Shen, L., Pechenizkiy, M., and Wang, Z. Sparse training via boosting pruning plasticity with neuroregeneration. In *International Conference on Neural Information Processing Systems*, 2022.
- Louizos, C., Ullrich, K., and Welling, M. Bayesian com pression for deep learning. In *31st Conference on Neural Information Processing Systems*, 2017.
- Louizos, C., Welling, M., and Kingma, D. P. Learning
  sparse neural networks through *l*<sub>0</sub> regularization. In *International Conference on Learning Representations*, 2018.
- Ma, X., Yuan, G., Shen, X., Chen, T., Chen, X., Chen, X.,
  Liu, N., Qin, M., Liu, S., Wang, Z., and Wang, Y. Sanity
  checks for lottery tickets: Does your winning ticket really
  win the jackpot? In *International Conference on Neural Information Processing Systems*, 2021.
- Mozer, M. C. and Smolensky, P. Skeletonization: A technique for trimming the fat from a network via relevance
  assessment. In *International Conference on Neural Infor- mation Processing Systems*, 1988.
- Neyshabur, B., Li, Z., Bhojanapalli, S., LeCun, Y., and
  Srebro, N. The role of over-parametrization in generalization of neural networks. In *International Conference on Learning Representations*, 2019.

548

549

Ramanujan, V., Wortsman, M., Kembhavi, A., Farhadi, A., and Rastegari, M. What's hidden in a randomly weighted neural network? In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020.

- Rosenfeld, J. S., Rosenfeld, A., Belinkov, Y., and Shavit, N. A constructive prediction of the generalization error across scales. In *International Conference on Learning Representations*, 2020.
- Savarese, P., Silva, H., and Maire, M. Winning the lottery with continuous sparsification. In *International Confer*ence on Neural Information Processing Systems, 2020.
- Shi, W., Cao, J., Zhang, Q., Li, Y., and Xu, L. Edge computing: Vision and challenges. *IEEE Internet of Things Journal*, 2016.
- Singh, S. P. and Alistarh, D. Efficient second order derivatives for network compression. In *International Conference on Neural Information Processing Systems*, 2020.
- Templeton, A., Conerly, T., Marcus, J., Lindsey, J., Bricken, T., Chen, B., Pearce, A., Citro, C., Ameisen, E., Jones, A., Cunningham, H., Turner, N. L., McDougall, C., MacDiarmid, M., Tamkin, A., Durmus, E., Hume, T., Mosconi, F., Freeman, C. D., Sumers, T. R., Rees, E., Batson, J., Jermyn, A., Carter, S., Olah, C., and Henighan, T. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. Technical report, Anthropic, 2024.
- Thimm, S. and Hoppe, H. Evaluating pruning methods. In *Proceedings of the International Joint Conference on Artificial Intelligence*, 1995.
- Trevor Gale, Erich Elsen, S. H. The state of sparsity in deep neural networks. *arXiv:1902.09574*, 2019.
- Verma, S. and Pesquet, J.-C. Sparsifying networks via subdifferential inclusion. In *Proceedings of the 38th International Conference on Machine Learning*, 2021.
- Wang, K., Liu, Z., Lin, Y., Lin, J., and Han, S. Haq: Hardware-aware automated quantization with mixed precision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- Wang, Q., Dun, C., Liao, F., Jermaine, C., and Kyrillidis, A. LOFT: Finding lottery tickets through filter-wise training. In *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics*, 2023.
- Wortsman, M., Farhadi, A., and Rastegari, M. Discovering neural wirings. In *International Conference on Neural Information Processing Systems*, 2019.
- Xiao, X., Wang, Z., and Rajasekaran, S. Autoprune: Automatic network pruning by regularizing auxiliary parameters. In *International Conference on Neural Information Processing Systems*, 2019.

550 551 552	Yite, W., Dawei, L., and Ruoyu, S. Ntk-sap: Improving neural network pruning by aligning training dynamics. <i>arXiv:2304.02840</i> , 2023.
553 554 555 556	Zhai, X., Kolesnikov, A., Houlsby, N., and Beyer, L. Scaling vision transformers. <i>arXiv preprint arXiv:2106.04560</i> , 2021.
557 558 559 560 561	Zhou, H., Lan, J., Liu, R., and Yosinski, J. Deconstructing lottery tickets: Zeros, signs, and the supermask. In <i>Inter-</i> <i>national Conference on Neural Information Processing</i> <i>Systems</i> , 2019.
562 563 564 565	Zhu, M. and Gupta, S. To prune, or not to prune: Explor- ing the efficacy of pruning for model compression. In <i>International Conference on Learning Representations</i> , 2018.
566 567 568 569	
570 571 572 573	
574 575 576 577	
578 579 580	
581 582 583 584	
585 586 587 588	
589 590 591 592	
593 594 595 596	
596 597 598 599	
600 601 602 603 604	

#### 605 A. Analysis

#### 607 A.1. The Role of Weight Importance in Generating Stronger Flows

<sup>608</sup> In this section, we provide theorethical groundings for the correlation between weight importance and the magnitude of the <sup>609</sup> corresponding flow within the Hyperflows framework. To this end, i.e. we analyse why important weights induce stronger <sup>610</sup> flows within the Hyperflows framework, we analyze the relationship between weight importance and the resulting gradient <sup>611</sup> feedback. Specifically, we show that if pruning a weight  $\theta_i$  leads to a larger increase in the loss function compared to pruning <sup>612</sup> another weight  $\theta_j$ , then the flow of  $\theta_i$  is larger than that of  $\theta_j$ .

<sup>613</sup> <sup>614</sup> To assess the significance of a specific weight, we measure the change in loss induced by pruning the weight. Formally, a <sup>615</sup> weight  $\theta_i$  is considered more important than a weight  $\theta_j$  if:

 $\Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_i\}) > \Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_j\}), \quad \text{where} \quad \Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_i\}) = \mathcal{L}\big(\mathcal{T} \setminus \{\theta_i\}\big) - \mathcal{L}(\mathcal{T}).$ 

Note that we make the assumption that all weights have a positive importance value, i.e.  $\Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_i\}) > 0$ .

$$\mathcal{F}(\theta_i, \mathcal{T}) = \begin{cases} \frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_i} & \text{if } t_i \le 0, \\ 0 & \text{otherwise.} \end{cases}$$
(16)

**Proposition A.1.** If  $\Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_i\}) > \Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_j\})$ , then the flow  $\mathcal{F}(\theta_i, \mathcal{T})$  exceeds  $\mathcal{F}(\theta_j, \mathcal{T})$  in magnitude:

$$\left|\frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_i}\right| > \left|\frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_j}\right|$$

*Proof.* Consider two weights  $\theta_i$  and  $\theta_j$  with flow parameters  $t_i$  and  $t_j$ , respectively. Let  $\delta > 0$  be a small perturbation applied to the flow parameters such that pruning a weight involves decrementing its parameter by  $\delta$ , ensuring  $t_k - \delta \le 0$  for  $k \in \{i, j\}$ . This results in the topology changes  $\mathcal{T} \setminus \{\theta_i\}$  and  $\mathcal{T} \setminus \{\theta_j\}$ .

We analyze how the pruning of each of the two weights affects the loss function  $\mathcal{L}(\mathcal{T})$ :

1. **Pruning**  $\theta_i$ : When  $\theta_i$  is pruned, the topology changes to  $\mathcal{T} \setminus {\theta_i}$ . Using a first-order Taylor expansion around  $\mathcal{T}$ , the change in loss due to pruning  $\theta_i$  is approximated by:

$$\mathcal{L}(\mathcal{T} \setminus \{\theta_i\}) - \mathcal{L}(\mathcal{T}) \approx -\delta \cdot \frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_i} + \mathcal{O}(\delta^2).$$

2. **Pruning**  $\theta_j$ : Similarly, pruning  $\theta_j$  results in the topology  $\mathcal{T} \setminus \{\theta_j\}$  and the corresponding change in loss due to pruning  $\theta_j$  is approximated by:

$$\mathcal{L}(\mathcal{T} \setminus \{\theta_j\}) - \mathcal{L}(\mathcal{T}) \approx -\delta \cdot \frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_j} + \mathcal{O}(\delta^2).$$

Given the assumption that pruning  $\theta_i$  induces a larger increase in loss compared to pruning  $\theta_i$ :

$$\Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_i\}) > \Delta \mathcal{L}(\mathcal{T} \setminus \{\theta_j\}),$$

substituting the approximations yields:

$$\left| -\delta \cdot \frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_i} \right| > \left| -\delta \cdot \frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_j} \right|.$$

651 Since  $\delta > 0$  is a constant, this simplifies to:

$$\left|\frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_i}\right| > \left|\frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_j}\right|$$

Thus, the above inequality directly translates to:

 $\left|\mathcal{F}(\theta_i, \mathcal{T})\right| > \left|\mathcal{F}(\theta_j, \mathcal{T})\right|.$ 

#### A.2. Influence of Weight Magnitude and Direction on Flow Parameters

To observe the relationship between the magnitude and direction of  $\theta_i$  and its flow parameter  $t_i$ ,  $t_i > 0$ , we analyze the partial derivatives of the loss function  $\mathcal{L}$  with respect to both  $\theta_i$  and  $t_i$ . We consider the following derivatives:

$$\frac{\partial \mathcal{L}}{\partial t_i} = \frac{\partial \mathcal{L}}{\partial \alpha} \cdot \frac{\partial \alpha}{\partial H} \cdot \frac{\partial H}{\partial t_i} = \frac{\partial \mathcal{L}}{\partial \alpha} \cdot \theta_i \cdot \mathcal{I} \cdot \mathbf{1}$$
$$\frac{\partial \mathcal{L}}{\partial \theta_i} = \frac{\partial \mathcal{L}}{\partial \alpha} \cdot \frac{\partial \alpha}{\partial \theta_i} = \frac{\partial \mathcal{L}}{\partial \alpha} \cdot H(t_i) \cdot \mathcal{I}$$

where  $\alpha = \theta_i \cdot H(t_i) \cdot \mathcal{I}$  and  $H(t_i) = 1$ .

**Proposition A.2.** The sign of  $\frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_i}$  is influenced by both the gradient of the loss with respect to  $\theta_i$  and the value of  $\theta_i$ . Specifically, the update direction of  $t_i$  is determined by the product  $\theta_i \cdot \frac{\partial \mathcal{L}(\mathcal{T})}{\partial \theta_i}$ .

*Proof.* From the expression for  $\frac{\partial \mathcal{L}(\mathcal{T})}{\partial t_i}$ , we observe that:

$$\frac{\partial \mathcal{L}}{\partial t_i} = \theta_i \cdot \frac{\partial \mathcal{L}}{\partial \theta_i}, \forall t_i > 0.$$

Table 3 showcases the implications based on the sign and magnitude of  $\theta_i$ .

Table 3. Implications Based on the Sign of  $\theta_i$ 

Weight Type	<b>Derivative</b> $\frac{\partial \mathcal{L}}{\partial \theta_i}$	Change in $\theta_i$	<b>Gradient</b> $\frac{\partial \mathcal{L}}{\partial t_i}$	Effect on t <sub>i</sub>	Implication
$\theta_i > 0$	Positive	$ heta_i\downarrow$	Positive	$t_i \downarrow$	Reinforces pruning
$\theta_i > 0$	Negative	$ heta_i\uparrow$	Negative	$t_i \uparrow$	Promotes regrowth
$\theta_i < 0$	Positive	$  heta_i \uparrow$	Positive	$t_i \uparrow$	Promotes regrowth
$\theta_i < 0$	Negative	$  heta_i \downarrow$	Negative	$t_i \downarrow$	Reinforces pruning

#### **B.** Ablation Studies

#### **B.1. Factors influencing flow** $\mathcal{F}(\theta_i, \mathcal{T})$

We begin by analyzing how the flow value  $\mathcal{F}(\theta_i, \mathcal{T})$  is influenced by factors other than the learning rate  $\eta_t$ . Our findings from Section 3.3, suggest that weight learning affects the behavior of flow, by changing the final convergence point a network will reach for the same constant pressure  $\gamma$ . We study this effect in the case of LeNet-300. We run the network for 1000 epochs for three different learning rates of 0.005, 0.0005 and 0.00005, with no schedulers used and the same constant  $\gamma$ .

Our findings are summarized in Figure 5, which shows that increasing weight learning rate  $\eta_{\theta}$  leads to smaller flows and convergence at higher sparsities.

Given the impact of  $\theta$  learning rate on network convergence, we study the influence of high and low learning rates on our pruning and regrowth phases. In our experiments, we study three setups on ResNet-50 CIFAR-10. In the first two experiments, we study how constant learning rates across the entire pruning and regrowth process affect sparsity and regrowth. We choose a high learning rate of 0.01 and a low learning rate of 0.0001. For our third experiment, we start with the high learning rate which is then decayed using cosine annealing to a low learning rate until the end of regrowth. For all three studies we let our scheduler guide the network towards the same sparsity rate of 1%. However, we observe significant differences in the regrowth stage. For the first experiment, regrowth does not occur at all, with more weights being pruned even after the pressure is set to 0, while 



Figure 5. MNIST convergence for constant  $\gamma = 1$  for different learning rates

for the low learning rate, the performance initially degrades, but is followed by a substantial regrowth stage where the number of remaining parameters increases by 60%. For the third experiment performance does not degrade as much as for the low learning rate and the regrowth is done in a more controlled way, experiencing an increase in remaining parameters of 35%. The results are illustrated in Figure 6.



Figure 6. The impacts of weights learning rate on pruning and sparsity

Lastly, we study how weight flow is affected by weight decay. Being directly applied on the weights, weight decay acts on both pruned and present weights. If a weight has been pruned in the first epochs on the training, weight decay will keep making it smaller and smaller, in this way diminishing its flow. We run similar experiments to the ones before, with a learning rate of 0.01, decayed during training to 0.0001, both with and without the standard weight decay. As expected, we observe in Figure 7 that regrowth without weight decay if more ample. We run this experiment five times, and note that each time the pattern illustrated in the figure remain consistent.



Figure 7. Effect of weight decay on the regrowth process

### **B.2.** Weights $\eta_{\theta}$ and pruning

Given the large impact  $\eta_{\theta}$  has on flow, we explore its implications for producing an optimal pruning setup for Hyperflows. We run three experimental setups on ResNet-50 CIFAR-10 similar to the ones before. For each one of them, we select a starting learning rate, which is then decayed during training to 0.0001 to ensure convergence. For this setup, we run experiments using  $\eta_{\theta} = 0.1, 0.01, 0.0001$ . We analyze the results from the perspective of accuracy after pruning, noise, regrowth and final accuracy. We find that the third setup is the most effective for Hyperflows.

We observe that each of the four studied aspects has a relationship with the learning rate. The noise is increased as initial learning rate increases, accuracy at the end of pruning is decreased the most for low learning rates and the highest for large



Figure 8. Training evolution for different learning rate configurations.

learning rates. We obtain the highest final accuracy for higher learning rates and the regrowth phase is diminished the higher the learning rate. These relationships hold and can be easily seen in Figure 8

# **B.3.** $\eta_t$ values and regrowth

We analyze regrowth behavior for several values of  $\eta_t$ . At regrowth stage, we scale  $\eta_t$  with 5, 10, 20, 30 for VGG-19 CIFAR-100 to observe the behavior of regrowth stage. Our findings are summarized in Figure 9. As  $\eta_t$  increases so does the number of regrown weights. However, we note that after a point, generally about an increase of 50% in remaining parameters, the effects of regrowth start to be diminished and starts introducing noise in the performance, while also regrowing more weights.



Figure 9. How differently scaled  $\eta_t$  affect regrowth

## <sup>815</sup> C. Extended experiments

#### 817 C.1. Layerwise sparsity levels & Weight Histograms

In this section, we examine the layer-wise sparsity observed for ResNet-50 on CIFAR-10 across the following pruning rates: 99.75%, 99.01%, and 98.13%. As illustrated in Figure 13, the overall sparsity hierarchy is maintained, displaying a decreasing trend in sparsity from the initial layers down to the final layer, where this pattern is interrupted. We hypothesize that earlier layers retain more weights due to their critical role in feature extraction, while deeper layers can sustain higher levels of pruning without significantly impacting overall performance. Notably, the penultimate layer experiences the highest degree of pruning, which means that it contains higher redundancy or less critical weights for performance. Furthermore, by analyzing the weight histograms for ResNet-50 with sparsity levels of 99.01% and 99.74% in Figure 11, we observe the

influence of sparsity on the weight distributions. High sparsity levels significantly alter weight distributions, demonstrating
 that extreme pruning not only reduces the number of active weights but also changes the underlying weight dynamics within
 the network.

828 the n 829

838

839

872

873

874

875 876

877

878

879

The histograms in Figure 12 illustrates the differences in weight distributions between the pruning and regrowth stages 830 on ImageNet with ResNet-50 at approximately 4.23% remaining weights. In the pruning stage, weights are more evenly 831 distributed across the range of [-0.4, 0.4], with a noticeable dip near zero, reflecting the removal of low-magnitude weights. 832 In contrast, during regrowth stage the weight distribution shifts significantly, showing a sharp clustering of weights around 833 zero, indicating the reactivation of low-magnitude weights during this process. This change in distribution correlates with a 834 notable performance gap: the regrowth stage achieves 72.4% accuracy, while the pruning stage reaches only 66.13%, we 835 consider the cause of this to be the fact that during the pruning process the small magnitude weights are pruned and during 836 the regrowth phase we recover from these weights the ones that improve performance the most. 837

#### C.2. Weight flips & Implicit regrowth

Implicit regrowth serves as the main source of noise in our network, promoting diverse topologies throughout the training
process. In Figure 10, we identify patterns in flip frequency, such as the lower number of flips at the start of training. This
behavior is anticipated, as pruning a critical weight early on allows its features to be more readily absorbed by other weights.
Around iteration 14, we notice a plateau followed by a brief decline in weight flips, which we attribute to the network
stabilizing during this phase.

As training progresses and the number of parameters declines, the per-weight flip frequency continues to increase, while the overall flip frequency remains relatively steady, resulting in a continue increase of the per-weight flip frequency. The regrowth phase is marked by a sharp decrease in the total number of flips as the network stabilizes and the learning rate of flow parameters diminishes toward zero. This pattern is visible between iterations 70 and 130, alongside a gradual increase in the number of parameters.



*Figure 10.* Frequency of Flips: The blue histogram represents the percentage of remaining parameters on a logarithmic scale, while the orange histogram illustrates the ratio of parameter flips per iteration relative to the total number of network parameters, also on a logarithmic scale. In our figure, one iteration is equivalent to the aggregation of 100 actual training iterations. We aggregate iterations to present the flip data in a more manageable way.

In Figure 10 we can observe the behavior of *flow* in relation to the gradients of t values. Two specific type of weights emerge, as we stated in the methodology Section 3.2. Note that negative values of the gradients translate into positive updates for t values and vice-versa. The first type of weight can be seen in the top-left and bottom-right diagrams in Figure 14, where



*Figure 11.* Weight values histograms of ResNet-50 on CIFAR-10 at Different Sparsity Levels. Top 99.75% sparsity, bottom 99.1% sparsity. We can observe a reshape of weight distributions



*Figure 12.* Weight Histograms: The upper figure depicts ResNet-50 during the pruning phase, achieving an accuracy of 66.13%. In contrast, the lower figure shows ResNet-50 in the regrowth phase, attaining an accuracy of 70.51%. Both phases maintain approximately 99.56% sparsity on ImageNet.



#### Hyperflows: Pruning Reveals the Importance of Weights

1040 *Figure 14.* Gradient values over time corresponding to four remaining weight of a pruned network. The blue values represent gradients 1041 while  $t_i > 0$ , while red values represent gradients for  $t_i \le 0$ . We can observe that red gradients, if they exist for that weight, have an 1042 average with very high magnitude, which is the flow  $\mathcal{F}(\theta_i, \mathcal{T})$ , while positive gradients  $\mathcal{M}(\theta_i, \mathcal{T})$  are much smaller, but in some cases 1043 big enough to oppose pressure for several iterations.

the gradient  $\mathcal{M}(\theta_i, \mathcal{T})$  does not oppose significant pressure for  $t_i > 0$ . This leads to the weight being pruned multiple times, which coincides, with large negative values in the gradient, which push  $t_i$  back over 0. The second type of weight, as common as the first one, does not get pruned at all. In this case,  $\mathcal{M}(\theta_i, \mathcal{T})$  averaged over several iterations, attempts to increase the magnitude of the weight, therefore increasing  $t_i$  at the same time, which leads to the weight not being pruned at all. We can see that in this case the overall magnitude of the gradients is below -1.5, which in our experiment was enough to resist pressure.

### <sup>2</sup> D. Training setup and reproducibility

In this section, we summarize the details regarding hyperparameters, optimizers and initializers used in our experiments.

Table 4. Comparison of experimental settings and results across various datasets (CIFAR-10, CIFAR-100, MNIST, and ImageNet-1K) using different neural network architectures. We make our notation as follows: <sup>w</sup> represent relation to weights, <sup>t</sup> represents relation to pruning values, p,  $p_s$  and  $p_e$  represent pruning stage, pruning start and end, r,  $r_s$  and  $r_e$  regrowth stage, regrowth stage start and end.  $\eta$ is learning rate, S scheduler and O optimizer. For example  $\eta_{ps}^w$  represents the learning rate of  $\theta$  values and the beginning of pruning state, and  $\eta_p^t$  represents the learning of t values during pruning stage. For flow params, we use a constant  $\eta^t$  for pruning stage, while for regrowth we employ a exponential decay given by  $\lambda_r^t$ . For weights, we use a cosine annealing decay from the start to the end of pruning and naother cosine decay from the start of regrowth to the end of it. We find it useful to have a discontinuity when transitioning from pruning to regrowth, as it helps with training.

Dataset	CIFAR-10		CIFAR-100		MNIST	ImageNet-1K
Network	<b>ResNet-50</b>	<b>VGG19</b>	<b>ResNet-50</b>	<b>VGG19</b>	<b>LeNet-300</b>	<b>ResNet-50</b>
Acc (%)	93.0 ± 0.5	94.0 ± 0.4	93.0 ± 0.6	72.0 ± 1.2	75.0 ± 1.1	72.0 ± 1.3
Batch size	128	128	128	128	128	1024
$\eta^w_{ps} \ \eta^w_{pe} \ \eta^w_{rs} \ \eta^w_{re}$	0.1	0.1	0.1	0.1	0.001	0.1
	0.003	0.003	0.003	0.003	0.001	0.003
	0.001	0.001	0.001	0.001	0.001	0.001
	0.0001	0.0001	0.0001	0.0001	0.00001	0.0001
$\mathcal{O}^w \ \mathcal{S}^w_\mathbf{p} \ \mathcal{S}^w_\mathbf{r}$	SGD	SGD	SGD	SGD	ADAM	SGD
	Cosine	Cosine	Cosine	Cosine	Cosine	Cosine
	Cosine	Cosine	Cosine	Cosine	Cosine	Cosine
$\begin{array}{c} \eta_p^t \\ \eta_{rs}^t \end{array}$	0.001	0.001	0.001	0.001	0.001	0.001
	0.01	0.01	0.01	0.01	0.001	0.01
$egin{array}{c} \lambda_r^t \ \mathcal{S}_{\mathbf{r}}^t \end{array}$	0.75	0.75	0.75	0.75	0.75	0.75
	LambdaLR	LambdaLR	LambdaLR	LambdaLR	LambdaLR	LambdaLR
$\mathcal{O}^t$	ADAM	ADAM	ADAM	ADAM	ADAM	ADAM
Initialization	Kaiming	Kaiming	Kaiming	Kaiming	Xavier	Kaiming
<b>Epochs</b>	160	160	160	160	60	120
#Prune end	100	100	100	100	30	90