

---

# Neural expressiveness for beyond importance model compression

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Neural Network Pruning has been established as driving force in the exploration of  
2 memory and energy efficient solutions with high throughput both during training  
3 and at test time. In this paper, we introduce a novel criterion for model com-  
4 pression, named “Expressiveness”. Unlike existing pruning methods that rely  
5 on the inherent “Importance” of neurons’ and filters’ weights, “Expressiveness”  
6 emphasizes a neuron’s or group of neurons ability to redistribute informational  
7 resources effectively, based on the overlap of activations. This characteristic is  
8 strongly correlated to a network’s initialization state, establishing criterion auton-  
9 omy from the learning state (*stateless*) and thus setting a new fundamental basis  
10 for the expansion of compression strategies in regards to the “When to Prune”  
11 question. We show that expressiveness is effectively approximated with arbitrary  
12 data or limited dataset’s representative samples, making ground for the exploration  
13 of *Data-Agnostic strategies*. Our work also facilitates a “hybrid” formulation of  
14 expressiveness and importance-based pruning strategies, illustrating their com-  
15plementary benefits and delivering up to  $10\times$  extra gains w.r.t. weight-based  
16 approaches in parameter compression ratios, with an average of 1% in performance  
17 degradation. We also show that employing expressiveness (independently) for  
18 pruning leads to an improvement over top-performing and foundational methods in  
19 terms of compression efficiency. Finally, on YOLOv8, we achieve a 46.1% MACs  
20 reduction by removing 55.4% of the parameters, with an increase of 3% in the  
21 mean Absolute Precision ( $mAP_{50-95}$ ) for object detection on COCO dataset.

## 22 1 Introduction

23 To address the computational constraints of existing models, Model Compression [7] has emerged as  
24 a prominent solution in exploring models that achieve comparable performance, but with reduced  
25 computational complexity [52]. Within this scope, Floating Point Operations (*FLOPs*) are used to  
26 estimate a model’s computational complexity, by measuring the arithmetic operations required for a  
27 forward pass, while parameters (*params*) are associated with a model’s size in terms of memory space  
28 [48] and their reduction can be a precursor towards more energy efficient solutions [5]. Although  
29 *FLOPs* and *params* often correlate, their relationship isn’t strictly linear. For instance, VGG16 [43]  
30 has  $17\times$  more parameters than ResNet-56 [17] but only  $3\times$  more *FLOPs*, largely due to VGG16’s  
31 extensive use of fully connected layers. At first sight, this can be attributed to the differences in  
32 network topologies. From a deeper perspective, the intricacies of various operations at handling  
33 computational workloads, such as residual structures [17, 55], depthwise separable convolutions [19],  
34 inverted residual modules [18], channel shuffle operations [59] and shift operations [53], coupled  
35 with their interplay, may significantly affect the relationship between *FLOPs* and *params* in a neural  
36 network. In a nutshell, besides the use of more computationally efficient operations as above-

37 mentioned, Model Compression aims to maintain model performance while optimizing the two  
38 aforementioned metrics via tensor decomposition, data quantization, and network sparsification [7].

39 In this paper we emphasize on the sparsification strategy of pruning [49], which we use as a basis  
40 framework to introduce “**Expressiveness**” as a new criterion for compressing neural networks.  
41 Existing pruning methods focus on removing redundant network elements – be they weights, neurons,  
42 or structures of neurons – in ways that minimally affect the overall performance of a network, based  
43 on the criterion of “**Importance**”, e.g. [38, 58, 20, 30]. Importance-based methods address questions  
44 like “How much does the removal of a network’s element cost in terms of performance degradation?”  
45 and “How much information does a network element contain?” in various ways. More specifically,  
46 they are motivated by the information inherent in network elements, such as the magnitude of weights  
47 [15, 28], similarity of weights or weight matrices [29, 60] ; and their sensitivity to the network’s loss  
48 function, such as the magnitude of gradients [38] and more [49, 3]. Such dependencies on weights’  
49 distributions constitute the aforementioned pruning methods to be “data-aware” since they intrinsically  
50 rely on the input data and the information state of the model, making the importance estimation  
51 of the network’s elements challenging and often costly due to factors like i) the stochasticity from  
52 training with minibatches, ii) the presence of plateau areas in the optimization space, and iii) the  
53 complexity introduced by nonlinearities [38]. Liu et al. [36] have also discussed limitations in the  
54 perception of importance within trained models, i.e. the authors criticize the ability of network’s  
55 elements importance to generalize to pruned derivatives, while also questioning the necessity of  
56 training large-scale models prior pruning.

57 Inspired by the concepts of “Information Plasticity” [2] and the “Lottery Ticket Hypothesis” (LTH)  
58 [12], we aim to address the limitations of previous importance-based methods through elaborating  
59 the “Expressiveness” criterion in model compression. In contrast to “Importance”, we focus on  
60 understanding the capability of network elements to redistribute informational resources to subsequent  
61 network elements. We define “Expressiveness” as - “*A neuron’s or group’s of neurons potential*  
62 *(when a network is not fully trained) or ability (when it is trained) to extract features that maximally*  
63 *separate different samples*”. As derived by [2], the early training phase of a model is crucial in  
64 shaping its expressiveness, with the formation of critical paths —strong connections that determine  
65 the “workload distribution”— being particularly significant during these initial stages. It’s essential to  
66 note that the network’s initialization state influences the formation of those paths, which interestingly  
67 enables “Expressiveness” to be a fit criterion for compression during all time instances of a networks’  
68 convergence [12], setting a baseline for answering the question of “When to prune?” [42]. *Our*  
69 *proposed pruning metric centers on measuring the overlap of activations between datapoints of the*  
70 *feature space*. In that way, expressiveness is based on effectively evaluating the inherent ability of the  
71 network’s neurons to differentiate sub-spaces within the feature space. We experimentally show that  
72 utilizing either small sets of arbitrary data points from the feature space or stratified sampling [34]  
73 from each class yields consistent estimations of expressiveness. Finally, we propose and implement a  
74 new “hybrid” pruning optimization strategy that cooperatively searches, exploits and characterizes  
75 the complementary benefits between “Importance” and “Expressiveness” for model compression.  
76 In summary, this work offers the following four-fold contribution: (i) we propose Expressiveness,  
77 a novel criterion based on the overlap of activations for model compression; (ii) we provide an  
78 in-depth theoretical analysis of both the fundamental principles and the technical intricacies of the  
79 proposed criterion; (iii) we validate the hypothesis that Expressiveness can be approximated with  
80 little to none input data, opening the road for data-agnostic pruning strategies; and (iv) through  
81 extensive experimentation we offer a thorough comparison w.r.t to both foundational and state-of-  
82 the-art methods demonstrating the efficiency and effectiveness of the proposed technique in model  
83 compression, while also examining the feasibility and effectiveness of a “hybrid” expressiveness-  
84 importance pruning strategy.

85 Specifically, we validate “Expressiveness” on the CIFAR-10 [24] and ImageNet [40] datasets using  
86 a variety of models with different design characteristics [44, 17, 45, 21, 19]. We demonstrate the  
87 superiority of our novel criterion over existing solutions, including many top performing structural  
88 pruning methods [31, 61, 58, 32, 23, 46, 11], and show significant params reduction while maintaining  
89 comparable performance. We experimentally explore and analyze the complementary nature of  
90 expressiveness and importance, showing that summary numeric evaluation provides up to 10×  
91 additional parameter compression ratio gains, with an average of 1% loss decrease w.r.t group  $\ell_1$ -  
92 norm [28]. Finally, we experiment on the current state-of-the-art computer vision model (YOLOv8  
93 [9, 22]), showcasing notable compression rates of 53.9% together with performance gains of 3% on

94 the COCO dataset [33], and highlighting the ability of more expressive neurons to better recover lost  
95 information from the pruning operation.

## 96 2 Related Work

97 **Weight (Non-Structural) Importance.** Han et al. [15, 14] and Guo et al. [13] approached the  
98 importance of weights based on their magnitude, removing connections below given thresholds.  
99 However, earlier works [25, 16] emphasized on the Hessian of the loss and have questioned whether  
100 magnitude is a reliable indicator of weight’s importance, as small weights can be necessary for  
101 low error. In this direction, several studies [4, 47, 41, 8] have proposed strategies of iterative  
102 magnitude pruning, in the form of “adaptive weight importance”, where weights are ranked based on  
103 their sensitivity to the loss. From a different perspective, Yang et al. [56] address the limitations of  
104 individual weight’s saliency that fail to account for their collective influence and provide a formulation  
105 of weight’s importance based on the error minimization of the output feature maps. Expanding on this  
106 concept, Xu et al. [54] propose a layer-adaptive pruning scheme that encapsulates the intra-relation  
107 of weights between layers, focusing on minimizing the output distortion of the network. Amongst  
108 other factors and limitations (as also discussed in 1), weight importance is very expensive to measure,  
109 mainly because of the increased complexity induced by the mutual influences of the weights among  
110 interconnected neurons. This, coupled with the requirement for specialized hardware to manage the  
111 irregular sparsity patterns resulting from weight pruning [57], has shifted research focus towards  
112 structural pruning [28], where neurons or entire filters are removed.

113 **Neuron and Filter (Structural) Importance.** Many were driven by the success of Iterative  
114 Shrinkage and Thresholding Algorithms (ISTA) [6] in non-structural sparse pruning and proposed  
115 filter-level adaptations [28, 29, 32, 26], based on the relaxation ( $\ell_1$  and  $\ell_2$ ) of  $\ell_0$  norm minimization.  
116 However, the loss of universality of such magnitude-based methods remains a limitation in the  
117 approximation of importance even in the structural scope. Yu et al. [58] further elaborate on the  
118 idea of error propagation ignorance, where the analysis is limited to the statistical properties of a  
119 single [28, 29] or two consecutive layers [37]. The authors suggest that the importance of neurons  
120 is better approximated from the minimization of the reconstruction error in the final response layer  
121 from which it is propagated to previous layers. In contrast to this view, Zhuang et al. [61] emphasize  
122 on the discriminative power of a filter as a more effective measure of importance and highlight that  
123 this aspect is not effectively assessed by the minimization of the reconstruction error. In a manner  
124 that reflects the progression of weight importance, Molchanov et al. [38] define “adaptive filter  
125 importance” as the squared change in loss and apply first and second-order Taylor expansions to  
126 accelerate importance’s computations. Predominantly, the data-awareness imposed by most pruning  
127 strategies is added to their already high-complexity – i.e. mostly non-convex, NP-Hard problems  
128 that require combinatorial searches. This renders the estimation of importance both computationally  
129 expensive and labor-intensive, similarly to non-structural approaches. Notably, Lin et al. [30] propose  
130 a less data-dependent solution based on the observation that the average rank of multiple feature maps  
131 generated by a single filter remains constant. HRank [30], alongside several other feature-guided  
132 filter pruning approaches, are valuable indicators towards data independence. Such works form a  
133 principle that pruning elements are better evaluated in the activation phase, where the importance of  
134 information and the richness of characteristics for both input data and filters are better reflected. In  
135 this work, we expand on this belief and we through extensive experimental analysis, we demonstrate  
136 that neither the information state nor the input data is required for the discriminative characterization  
137 of an element.

## 138 3 Neural Expressiveness

### 139 3.1 Weights and Activations: Importance vs Expressiveness

140 Neurons are the main constituent element of a neural network. Given a neural network  $\mathcal{N}$ , we  
141 denote neurons by  $a_i^{(l)}$ , where  $l \in L$  is indicative of the neuron’s layer in a network with  $L =$   
142  $\{l_0, \dots, l_l, \dots, l_{|L|}\}$  layers and  $i$  of its position in the given layer  $l = \{a_0, \dots, a_i, \dots, a_{|l|}\}$ . Another  
143 important element are the learning parameters of the network. Otherwise the weights represent the  
144 strength of connections between neurons in adjacent layers and are denoted by  $w_{ij}^{(l)}$ , where  $i$  and  $j$   
145 index the neurons in the current and previous layers. In that manner, neuron’s can be perceived as

146 switches that allow or block information from propagating through-out a network. The activation  
 147 (or not) of a neuron  $a_i^{(l)}$  depends on the output value of its activation function  $\sigma(\cdot)$ , where there are  
 148 many popular options for the definition of  $\sigma$ , e.g., sigmoid, tanh, and ReLU functions. Specifically, a  
 149 neuron’s output is defined as follows,

$$a_i^{(l)} = \sigma \left( \sum_j w_{ij}^{(l)} a_j^{(l-1)} + b_i^{(l)} \right) \quad (1)$$

150 where  $b_i^{(l)}$  denotes the bias term. From eq. 1, we observe that a neuron’s activation is affected by  
 151 the activation of the previous layers, hence affecting in the same way the consecutive layers. This  
 152 interdependence between activations  $a^{(l)}$ , for a given layer  $l$  defines a recurrent form that can be  
 153 generalized as follows,

$$a^{(l)} = \sigma \left( W^{(l)} f \left( a^{(l-2)}, \dots, a^{(1)} \right) + b^{(l)} \right). \quad (2)$$

154 On the other hand, weights are a more static representation of information as they modulate how  
 155 much influence one neuron’s activation has on another’s, compared to activations that control the  
 156 flow of information in a network. This differentiation has motivated us to define two axes of study in  
 157 the categorisation of pruning criteria, one based on the weights (“importance”) and one based on the  
 158 activation phase (“expressiveness”).

159 **Generalization of concepts in a structural level.** The aforementioned principles extend to the  
 160 structural representations of weights and activations, the most common being Convolutional Neural  
 161 Networks (CNNs). For a CNN model with a set of  $K$  convolutional layers, where  $C^l$  is the  $l$ -th  
 162 convolutional layer. We denote filters (weight maps) and feature maps (activation maps) as  $F_k^l$  and  
 163  $C_k^l$  respectively, where  $k$  is the index within a layer. Given filter with dimensions  $m \times n$ , eq. 1 is  
 164 adapted as follows,

$$C_k^{(l)}(x, y) = \sigma \left( \sum_{i=1}^m \sum_{j=1}^n F_{ij}^{(l,k)} a_{x+i-1, y+j-1}^{(l-1)} + b_k^{(l)} \right) \quad (3)$$

165 where  $(i, j)$  and  $(x, y)$  are the coordinates of weights and output activations within the filter and the  
 166 output activation map respectively. Similarly, a convolution layer  $l$  can be analytically expressed as  
 167 follows,

$$C^{(l)} = \begin{cases} \sigma \left( \bigoplus_{k=1}^{K^{(1)}} F^{(1,k)} * X + B^{(1)} \right) & \text{if } l = 1 \\ \sigma \left( \bigoplus_{k=1}^{K^{(l)}} F^{(l,k)} * C^{(l-1)} + B^{(l)} \right) & \text{if } l > 1 \end{cases} \quad (4)$$

168 with  $X$  being the input to the first layer of the network, and where symbol  $*$  denotes convolution  
 169 operation and  $\bigoplus$  denotes the concatenation operation. Within this context<sup>1</sup>, eq. 2 is generalized as  
 170 follows,

$$C^{(l)} = \sigma \left( \bigoplus_{k=1}^{K^{(l)}} F^{(l,k)} * f \left( C^{(l-2)}, \dots, C^{(1)} \right) + B^{(l)} \right). \quad (5)$$

171 **Conceptualization of information propagation.** Consider a task with  $X = \{x_i\}_{i=1}^{|D|}$  denoting  
 172 dataset samples, where  $|D|$  is the size of the dataset. Given the information state (weight state) of  
 173 a CNN model with  $K$  convolutional layers at a given time  $t_i$ ,  $X$  is mapped through the network as  
 174  $f(X, W_{t_i})$ , where  $W_{t_i} = \{F_{t_i}^1, \dots, F_{t_i}^l, \dots, F_{t_i}^{K^{(l)}}\}$  and  $F_{t_i}^l = \{F_{t_i}^{(l,1)}, \dots, F_{t_i}^{(l,k)}, \dots, F_{t_i}^{(l,K^{(l)})}\}$ ,  
 175 with  $K^{(l)}$  being the amount of weight maps (filters) in a given layer  $l$ . This process can be further  
 176 analyzed as follows,

$$f(X, \mathbf{W}_{t_i}) = \mathcal{F}_{|K|}(\mathcal{F}_{|K|-1}(\dots \mathcal{F}_1(X; \mathbf{F}_{t_i}^1); \mathbf{F}_{t_i}^2); \dots; \mathbf{F}_{t_i}^{K^{(l)}}), \quad (6)$$

177 where  $\mathcal{F}_l$  represents the mapping operation of convolutional layer  $l$ .

178 Based on eq. 2 and eq. 5, the equivalent of the previous based on the activations of the layers can be  
 179 expressed as,

$$f(X, \mathbf{W}_{t_i}) = C^{(K^{(l)})} \left( \dots \left( C^{(2)} \left( C^{(1)}(X, \mathbf{F}_{t_i}^1), \mathbf{F}_{t_i}^2 \right) \dots \right), \mathbf{F}_{t_i}^{K^{(l)}} \right). \quad (7)$$

<sup>1</sup>We do not include pooling and batch normalization layers in the formulations; however, the equations can be expanded to incorporate them as intermediate steps based on each architecture.

180 Here,  $C^{(l)}$  represents the activation map of the  $l$ -th layer, where  $C^{(l)} = \mathcal{F}_i(C^{(l-1)}; \mathbf{F}_{t_i}^l)$  aligns  
 181 with the structure defined in eq. 4. In this formulation,  $C^{(1)}$  is the activation map of the first layer,  
 182 computed using the input  $X$  and the first layer’s filters  $\mathbf{F}_{t_i}^1$ . Subsequent layers’ activation maps  
 183  $C^{(l)}$  are derived from the previous layer’s output  $C^{(l-1)}$  and their respective filters  $\mathbf{F}_{t_i}^l$ . Assuming a  
 184 classification task, the final layer  $C^{(|K|)}$  is considered the classification layer, effectively summarizing  
 185 the hierarchical feature extraction and transformation process across all convolutional layers.

### 186 3.2 Mathematical Foundation of Neural Expressiveness.

187 We observe that the training parameters of the model, in this case  $W_{t_i}^2$ , are responsible  
 188 for transforming the original input feature space  $X$  into a sequence of intermediate feature  
 189 spaces  $\{C^{(1)}, \dots, C^{(|K|-1)}\}$ , progressing towards the final prediction formulated by the prediction  
 190 layer  $C^{(|K|)}$ .

191 Based on this intrinsic characteristic of neural networks and inspired by the goal of optimizing  
 192 feature discrimination, akin to the entropy reduction strategy in decision trees [51], we assess network  
 193 elements ability, in this scenario filters, to extract features, i.e., activation patterns, that maximally  
 194 separate different input samples  $x_i$ . In other words, we score the expressiveness of the filters within  
 195  $W_{t_i}$ , based on the discriminative quality of the intermediate feature spaces they generate, where the  
 196 feature space generated by a filter  $F_k^l$ , is denoted as  $C_k^l$ .

197 **Neural Expressiveness foundational concept.** When assessing the expressiveness of an element  
 198 within  $W_{t_i}$  based on its generated feature spaces, e.g.,  $NEXP(F_{t_i}^l; C^l)$ , we cooperatively evaluate  
 199 all of its preceding elements, as derived from eq. 5. This can be formulated as,

$$NEXP(F_{t_i}^l; C^l) = NEXP(F_{t_i}^l; (C^{(l-1)}, C^{(l-2)}, \dots, C^{(1)})), \quad (8)$$

200 which can be further extended to incorporate the inter-dependencies between the examined element  
 201 and its predecessors, in accordance with eq. 7, as detailed below:

$$\begin{aligned} NEXP\left(F_{t_i}^l; \left(C^{(l-1)}, C^{(l-2)}, \dots, C^{(1)}\right)\right) = \\ NEXP\left(F_{t_i}^l; \left((C^{(l-2)}, \mathbf{F}_{t_i}^{l-1}), (C^{(l-3)}, \mathbf{F}_{t_i}^{l-2}), \dots, (X, \mathbf{F}_{t_i}^1)\right)\right). \end{aligned} \quad (9)$$

202 The aforementioned eqs. 8 and 9 provide the *foundational concepts for utilizing the evaluation of*  
 203 *the activation phase*, in an endeavor to encourage the development of more universal solutions by  
 204 addressing the limitations of universality inherent in the assessment of the weight state alone (as also  
 205 discussed in sections 1 and 2).

206 **Formulation of Neural Expressiveness (NEXP) Score.** Diving deeper into the Neural Expressive-  
 207 ness (NEXP) scoring process, we follow eq. 9 previously and assume a mini-batch  $X' = \{x'_i\}_{i=1}^N$ ,  
 208 with  $N$  being the number of samples in it. Mapping the batch through the network, based  
 209 on eqs. 6 and 7, generates a set of sequences of feature spaces (activation maps), denoted as  
 210  $S = \{s_1, \dots, s_i, \dots, s_N\}$ , where  $s_i = \{x'_i, \dots, C_i^l, \dots, C_i^{(|K|)}\}$  is the sequence of the activation  
 211 patterns generated from sample  $x'_i \in X'$  and  $|s_i| = |K| + 1$  is its cardinality, including the feature  
 212 space of sample  $x'_i$ . To evaluate a specific filter  $k$  in layer  $l$ , denoted as  $F_k^l$ , we utilize the retrieved  
 213 activation patterns from that filter, denoted as  $\{s_{i,k}^l\}_{i=1}^N$ , where  $s_{i,k}^l = C_{i,k}^l$  is the activation pattern  
 214 retrieved from filter  $k$  in layer  $l$ .

215 To score the Neural Expressiveness of  $F_k^l$ , we first construct a  $N \times N$  matrix that expresses all  
 216 possible combinations of the activation patterns derived from the different input samples. This table  
 217 can be visualised as follows,

$$\begin{pmatrix} s_{(1,1),k}^l & s_{(1,2),k}^l & \cdots & s_{(1,N),k}^l \\ s_{(2,1),k}^l & s_{(2,2),k}^l & \cdots & s_{(2,N),k}^l \\ \vdots & \vdots & \ddots & \vdots \\ s_{(N,1),k}^l & s_{(N,2),k}^l & \cdots & s_{(N,N),k}^l \end{pmatrix}. \quad (10)$$

<sup>2</sup>Bias terms are excluded for simplicity.

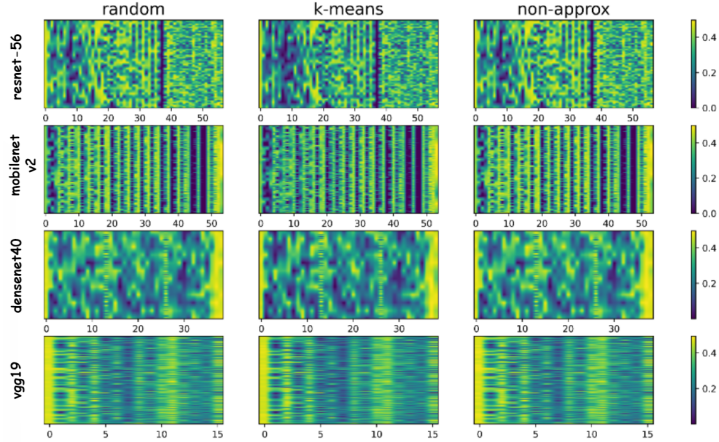


Figure 1: **Expressiveness statistics of feature maps from different convolutional layers and architectures on CIFAR-10.**

218 where  $s_{(i,j),k}^l$  denotes the dissimilarity of activations patterns between the  $i$ -th and the  $j$ -th sample  
 219 of the batch. In other words, the matrix in eq. 10 represents all the possible combinations of NEXP  
 220 calculations, where each element  $s_{(i,j),k}^l$  derives from  $f(s_{i,k}^l, s_{j,k}^l)$ , with  $f$  being any dissimilarity  
 221 function. Without loss of generality, for the rest of the study, we use the Hamming distance as the  
 222 operator implementing dissimilarity function. Activations are first binarized (values greater than 0  
 223 become 1, and the rest become 0), i.e. enabling to evaluate the degree of overlap between the binary  
 224 activation patterns using  $f$ .

225 We note that the matrix's diagonal, where  $i$  equals  $j$ , along with the elements below the diagonal,  
 226 where  $i$  is greater than  $j$ , do not contribute additional value to quantifying the discriminative ability  
 227 of an element. The diagonal elements represent comparisons of the same sample's activation patterns,  
 228 rendering them redundant. Meanwhile, the lower triangular elements are considered duplicates  
 229 since  $s_{(i,j),k}^l$  is equal to  $s_{(j,i),k}^l$ , thereby not adding any new information. Drawing from these two  
 230 observations, we define the Neural Expressiveness score (NEXP) as follows,

$$NEXP(F_k^l) = \frac{1}{\frac{N(N-1)}{2}} \sum_{i=1}^N \sum_{j=i+1}^N f(s_{i,k}^l, s_{j,k}^l) \quad (11)$$

231 The **more similar** the activation patterns derived from an element are, the **less expressive** it is  
 232 declared to be. In eq. 11, we also normalize the score w.r.t the total amount of combinations  
 233 ( $\frac{N(N-1)}{2}$ ), thereby deriving the average expressiveness score. This average score is then utilized to  
 234 characterize the discriminative capability/capacity of the examined network element. In this study,  
 235 we used the mean operation, however, we note that alternate statistical measures, e.g., minimum,  
 236 maximum, median, etc., could feasibly be applied in the computation of the overall score.

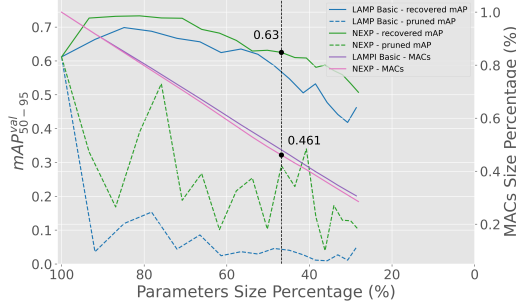
### 237 3.3 Dependency to Input Data

238 NEXP evaluates the inherent property of network elements to maximally distinguish between input  
 239 samples. We extend this line of thought and assess its sensitivity to input data  $X$  and mini-batch  
 240 size  $N$ , in order to delineate the dependence between NEXP and the input data. To achieve that,  
 241 we perform a sensitivity analysis of NEXP to the mini-batch data  $X$ , using two input sampling  
 242 strategies to assemble a batch with 60 samples, namely random sampling (denoted as 'random') and  
 243 class-representative sampling via k-means (denoted as 'k-means'). We define the true NEXP score  
 244 (denoted as 'non-approx') for each filter as the value obtained by comparing all activation patterns  
 245 across the *entire training dataset* (more info in A.1). Fig. 1 presents a detailed comparative illustration  
 246 of the results that highlight the similarities in NEXP estimations across various trained networks,  
 247 including VGGNet [44], ResNet [17], MobileNet [19] and DenseNet [21] on CIFAR-10 dataset.  
 248 Columns represent the aforementioned sampling strategies, while colors indicate expressiveness  
 249 levels, with higher values signifying greater expressiveness. In each sub-figure, the x-axis indicates

**Algorithm 1** NEXP Pruning Algorithm

**Define:**  $NEXP_{map} = \{\{NEXP(F_k^l)\}_{k=1}^{|C^l|}\}_{l=1}^{|K|}$   
**Require:** A mini-batch  $X$ , a neural network  $\mathcal{N}(W_{t_i})$ , a theoretical speed-up target, denoted  $\tau$ , and the allowed amount of pruning steps, denoted  $steps_{max}$ .  
**Ensure:**  $\frac{FLOPs(\mathcal{N})}{FLOPs(\mathcal{N}_{pruned})} \geq \tau$   
1: Initialize  $NEXP_{map} \leftarrow f(X; W_{t_i})$   
2: Initialize  $\tau_{current}$  as 1  
3: Initialize  $steps_{current}$  as 1  
4: Initialize  $\mathcal{N}_{pruned}$  as  $\mathcal{N}$   
5: **while** ( $\tau_{current} < \tau$ ) and ( $steps_{current} \leq steps_{max}$ ) **do**  
6:  $F_{to\_prune} = bottom_{\kappa}(NEXP_{map})$   
7:  $\mathcal{N}_{pruned}(W_{pruned}) = prune(\mathcal{N}_{pruned}, F_{to\_prune})$   
8:  $NEXP_{map} \leftarrow f(X; W_{pruned})$   
9:  $\tau_{current} = \frac{FLOPs(\mathcal{N})}{FLOPs(\mathcal{N}_{pruned})}$   
10:  $steps_{current} + +$   
11: **end while**  
12: **return**  $\mathcal{N}_{pruned}(W_{pruned})$

Figure 2: **Pruning YOLOv8m trained on COCO for Object Detection.** Comparative results between neural expressiveness (NEXP) and layer-adaptive magnitude-based pruning method (LAMP) [26]. More comparisons in the supplementary material.



convolutional layer indices, and the y-axis shows feature map indices per layer, standardized through pixel-wise interpolation to align with the layer having the most feature maps. Fig. 1 confirms that NEXP can be effectively estimated using random and limited data samples. Detailed results of this analysis, are presented in Appendix A. The comparative analysis reveals that a mini-batch of 60 samples (0.4% of  $D$  in this case) effectively approximates the NEXP scores calculated from the entire dataset, yielding consistent similarity scores above 99% across most similarity metrics (Table. 3).

**3.4 Pruning Process**

Alg. 1 describes the proposed NEXP-based pruning process, and it has been implemented as extension in the DepGraph pruning framework [11]. A target theoretical speed-up is specified, referred to as the Compression FLOPs Ratio ( $\downarrow$ ) and denoted by  $\tau$ . This ratio is calculated using the formula  $\frac{original\ FLOPs}{compressed\ FLOPs}$ . To achieve this target ratio, the network may undergo pruning in one or several steps, dictated by the intricacies of the pruning criterion and adjusted according to the quantity of elements removed at each step. For example, NEXP benefits from additional steps, since a filter’s score is reliant on its preceding elements (Section 3.2), and a more gradual update on the scores allows for improved pruning precision. A more in-depth analysis of Alg. 1 along with more details on the implementation options are presented in Appendix B.

**4 Experimental Evaluation**

Details on the experimental settings can be found in Appendix C, including the (a) Datasets and Models (C.1), (b) Adversaries (C.2), (c) Evaluation Metrics (C.3) and (d) Configurations (C.4).

**4.1 Comparison w.r.t. State-of-Art Model Compression Strategies**

**Image Classification on CIFAR-10 and Imagenet-1k.** We compare against a plethora of foundational and top-performing approaches, ranging from filter magnitude-based [28, 32, 29] and loss sensitivity-based [58] methods to feature-guided strategies [23, 30] and search algorithms [35, 31].

**Outcomes and Discussion.** Our findings for various target FLOPs pruning ratios are presented in Tab. 1 (and Tab.6-9 in Appendix D.2) for CIFAR-10, and in Tab. 2 for ImageNet. It is essential to acknowledge the subjectivity in reported performance metrics (accuracy), influenced by the fine-tuning process post-pruning, e.g. the authors in DCP [61] fine-tune for 400 epochs, in contrast to ours 100. We observe that our approach yields consistent improvements in params reduction compared to other methods for given FLOPs ratios, which notably scale significantly for regimes of higher target FLOPs compression ratios  $\tau$ . For example, on ResNet-56 we show  $+0.92\times$  average params reduction gains in the  $2\times-2.20\times$  FLOPs reduction regime, with  $-0.38\%$ ,  $+0.05\%$  and  $-0.37\%$  percentage difference in loss respectively to ABC [31], SCP [23] and HRank [30], while on ResNet-110 we

Table 1: Analytical Comparison of Importance-based solutions and Expressiveness on CIFAR-10 using ResNet architectures [17] - ResNet-56 (left) and ResNet-110 (right).

Method	top-1 acc		Compression Ratio ↓		Method	top-1 acc		Compression Ratio ↓	
	Base (%)	Δ (%)	#Params	#FLOPs		Base (%)	Δ (%)	#Params	#FLOPs
L1 [28]	93.06	+0.02	1.16×	1.37×	L1 [28]	93.55	+0.02	1.02×	1.19×
NEXP (Ours)	93.36	+0.05	<b>1.69</b> ×	1.53×	NEXP (Ours)	93.79	+0.66	<b>1.10</b> ×	1.20×
GAL-0.6 [32]	93.26	+0.12	1.13×	1.60×	GAL-0.1 [32]	93.50	+0.09	1.04×	1.23×
NISP-56 [58]	-	-0.03	1.74×	1.77×	HRank [30]	93.50	+0.73	1.65×	1.70×
DCP-Adapt [61]	93.80	+0.01	3.37×	1.89×	NISP-110 [58]	-	-0.18	1.76×	1.78×
HRank [30]	93.26	-0.09	1.74×	2.01×	NEXP (Ours)	93.79	+0.18	1.78×	1.80×
SCP [23]	93.69	-0.46	1.94×	2.06×	GAL-0.5 [32]	93.50	-0.76	1.81×	1.94×
NEXP (Ours)	93.36	-0.41	<b>2.87</b> ×	2.11×	HRank [30]	93.50	-0.14	2.46×	2.39×
ABC [31]	93.26	-0.03	2.18×	2.18×	NEXP (Ours)	93.79	+0.10	<b>2.72</b> ×	2.42×
NEXP (Ours)	93.36	-1.58	<b>4.3</b> ×	2.50×	ABC [31]	93.50	+0.08	3.09×	2.87×
GAL-0.8 [32]	93.26	-1.68	2.93×	2.51×	NEXP (Ours)	93.79	-0.37	<b>3.81</b> ×	3.01×
HRank [30]	93.26	-2.54	3.15×	3.86×	HRank [30]	93.50	-0.85	3.25×	3.19×
NEXP (Ours)	93.36	-5.12	<b>21.5</b> ×	5.00×	NEXP (Ours)	93.79	-0.59	<b>4.38</b> ×	3.27×

282 show  $+1.21\times$  average params reduction gains in the  $2.87\times$ - $3.27\times$  FLOPs reduction regime, with  
 283  $-0.67\%$  and  $+0.26\%$  percentage difference in loss respectively to ABC [31] and HRank [30]. Similar  
 284 observations are evident across all tables, where in certain regimes we also show notable performance  
 285 gains, up to  $+1.5\%$ , especially for VGGNet, which is more prone to params reductions due to its  
 286 plain structure.

287 **Object Detection with YOLOv8.** We evaluate expressiveness against four importance based  
 288 methods, i.e layer-adaptive magnitude-based pruning (LAMP) [26], network slimming (SLIM) [35],  
 289 Wang’s et al. proposed method (DepGraph) [11] and random pruning that serves as a generic pruning  
 290 baseline [3]. The experiments were conducted on the YOLOv8m model version [22], utilizing the  
 291 DepGraph pruning framework [11] with an iterative pruning schedule of 16 steps, where after each  
 292 pruning step the model was fine-tuned for 10 epochs using the coco128 dataset.

293 **Outcomes and Discussion.** We report the comparative pruning progress of expressiveness versus  
 294 the baseline methods, i.e. the remaining percentage of the original model in terms of *MACs* and  
 295 *params* after each pruning step, named *MACs* Size Percentage (MSP) and *Parameters* Size Percentage  
 296 (PSP) respectively, and highlight the  $mAP_{50-95}^{val}$  both after pruning (pruned mAP) and fine-tuning  
 297 (recovered mAP). We observe that expressiveness outperforms the rest of the reported methods across  
 298 the whole pruning spectrum, as shown in Fig. 2 (more in Appendix D.2), preserving the initial  
 299 performance of the model for percentage sizes that reach up to 40% ( $2.5\downarrow$ ) of that of the original  
 300 model, with less than 0.5% of recovered performance degradation. Our method even achieves a 3%  
 301 increase in recovered mAP for 46.1% MSP ( $2.17\downarrow$ ), in comparison to the baselines that showcase  
 302 weak recovery capabilities after the 60% ( $1.67\downarrow$ ) mark in both MSP and PSP. This can be attributed  
 303 to the intrinsic property of expressiveness to maintain network elements that are more robust to  
 304 information redistribution, in contrast to “important” labeled structures by other methods. In our  
 305 experimental scenario, that characteristic is further amplified by the iterative pruning format and the  
 306 higher amount of fine-tuning epochs at each step, in comparison to conventional pruning schedules  
 307 that fine-tune for 1 epoch after each iteration or perform a unified fine-tuning session after the last  
 308 pruning iteration. Interestingly, our criterion also demonstrates significant resistance to performance  
 309 loss after pruning, achieving 18% increased average performance in terms of pruned mAP compared  
 310 to the importance-based methods. We have empirically observed that expressiveness benefits from  
 311 increased cardinality in pruning granularity settings, e.g amount of intermediate steps to achieve a  
 312 given compression ratio. This stems from expressiveness interactive nature of all elements, as also  
 313 explained in Sec. 3, where smaller pruning steps combined with iterative fine-tuning, enhance pruning  
 314 precision and allow for “smoother” redistribution of information in a network, thus contributing to  
 315 the increased resistance to performance deficits after each pruning step.

## 316 4.2 Assessing Hybrid Compression space

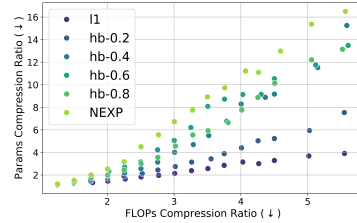
317 In this section, we assess the potential efficiency of “hybrid” pruning strategies exploiting the  
 318 cooperation between importance and expressiveness. We explore the solution space of “hybrid”  
 319 compression, using a linear combination of importance and neural expressiveness criteria. We guide  
 320 exploration through the scoring function:  $W_{imp} \cdot IMP + W_{nexp} \cdot NEXP$  and conduct experiments with



Table 2: Analytical Comparison of Importance-based solutions and Expressiveness on ImageNet-1k using ResNet-50 [17].

Method	Base (%)		$\Delta$ Acc (%)		Compression Ratio	
	top-1	top-5	top-1	top-5	#Params $\downarrow$	#FLOPs $\downarrow$
NISP-50-B [58]	-	-	-0.89	-	1.78 $\times$	1.79 $\times$
NEXP (Ours)	76.13	92.86	-1.35	-0.93	2.00 $\times$	2.02 $\times$
ThiNet [37]	72.88	91.14	-1.87	-1.12	2.06 $\times$	2.25 $\times$
DCP [61]	76.01	92.93	-1.06	-0.56	2.06 $\times$	2.25 $\times$
ABC [31]	76.01	92.96	-2.49	-1.45	2.27 $\times$	2.30 $\times$
NEXP (Ours)	76.13	92.86	-6.77	-3.43	<b>4.05<math>\times</math></b>	3.04 $\times$
GAL-1-joint [32]	76.15	92.87	-6.84	-3.75	2.50 $\times$	3.68 $\times$
Hrank [30]	76.15	92.87	-7.15	-3.29	3.08 $\times$	4.17 $\times$

Figure 3: **Linear exploration of the combinatorial space between importance and expressiveness.**



321 various weight combinations, subject to the constraint  $W_{\text{imp}} + W_{\text{nexp}} = 1$ . Given that exhaustive  
 322 search is impractical, we introduce the hyper-parameter  $\alpha \in \{0.0, 0.2, \dots, 0.8, 1.0\}$  to restrict the  
 323 set of permissible combinations, and modify the constraint to  $(1 - \alpha) \cdot W_{\text{imp}} + \alpha \cdot W_{\text{nexp}} = 1$ . We  
 324 use group L1-norm [28] as the importance criterion (IMP) and assess all permissible combinations  
 325 across a linear scale, denoted as  $\tau$ , representing the target FLOPs compression ratios that we utilized  
 326 for pruning, on ResNet-56 for CIFAR-10. The outcomes are visualized in Figure 3, which maps our  
 327 predetermined  $\tau$  values on the x-axis against the various parameter compression ratios achieved by  
 328 each combination. Regarding performance, we report the averaged percentage differences in top-1  
 329 accuracy between the baseline importance method (L1) and each hybrid format: -0.21% for hb-0.2,  
 330 -0.96% for hb-0.4, -1.55% for hb-0.6, -1.07% for hb-0.8, and -2.18% for NEXP.

331 **Observations.** A consistent pattern is observed across the values of  $\alpha$ , where larger values yield  
 332 higher params compression ratios. Notably, hybrid derivatives allow us to explore sub-spaces with  
 333 higher parameter compression ratios by sacrificing slight performance accuracy. We also observe  
 334 that the solution vectors corresponding to IMP and EXP act as extremal points in the solution space  
 335 of hybrid combinations, thus suggesting a degree of partial orthogonality between the two criteria.  
 336 Furthermore, the findings reveal a polynomial relationship between parameter compression ratios and  
 337 FLOPs reduction, with compression ratios increasing polynomially to linear increments in FLOPs  
 338 reduction, and thus enabling more efficient explorations.

### 339 4.3 Evaluating Neural Expressiveness at Initialization

340 The nature of NEXP allows to be applied in a weight agnostic manner, i.e. on untrained networks.  
 341 An extended version of the section’s 3.3 analysis, which also includes untrained models (Appendix  
 342 A), reveals that  $NEXP_{\text{map}}$ ’s obtained at initialization and after network convergence share some  
 343 expressiveness pattern similarities, particularly in the initial layers. Our numeric evaluation shows a  
 344 notable correlation between the initialization and converged states for DenseNet-40 and VGG-19,  
 345 with cosine similarities of 84.10% and 86.82%, respectively. It also indicates greater consistency in  
 346 neural expressiveness measurements for the first layers of all networks, which could be considered  
 347 important for the formation of critical paths [2]. Motivated by these observations, we also assess the  
 348 efficacy of expressiveness as criterion for Pruning at Initialization against various SOTA approaches  
 349 [27, 50, 46] (Appendix D.1). Our method consistently outperforms (in terms of top-1 acc) all other  
 350 algorithms, particularly in regimes of lower compression, up to  $10^2(\downarrow)$  with an average increase of  
 351 1.21% over SynFlow, while maintaining competitiveness at higher compression levels, above  $10^2(\downarrow)$   
 352 with an average percentage difference of 4.82%, 3.72% and -2.74%, compared to [50], [27] and [46].  
 353 In summary, under the assumption that the selection of hyperparameters remains congruent with  
 354 the initialization [12], consistent map measurements between initial and final states can effectively  
 355 evaluate NEXP’s ability to identify winning tickets. However, a robust evaluation should also consider  
 356 the initial state quality and the training process, while addressing the "When to prune" question [42].

## 357 5 Conclusions

358 In this work, we have introduced “Neural Expressiveness” as a new criterion for model compression.  
 359 In our NEXP steps, we will explore optimal solutions for the “When” and “How” to prune questions.

## References

- 360
- 361 [1] Armstrong Aboah, Bin Wang, Ulas Bagci, and Yaw Adu-Gyamfi. Real-time multi-class helmet violation  
362 detection using few-shot data sampling technique and yolov8. In *Proceedings of the IEEE/CVF Conference*  
363 *on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 5350–5358, June 2023.
- 364 [2] Alessandro Achille, Matteo Rovere, and Stefano Soatto. Critical learning periods in deep networks. In  
365 *International Conference on Learning Representations*, 2019.
- 366 [3] Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John Gutttag. What is the state of neural  
367 network pruning? In I. Dhillon, D. Papailiopoulos, and V. Sze, editors, *Proceedings of Machine Learning*  
368 *and Systems*, volume 2, pages 129–146, 2020.
- 369 [4] Miguel Á. Carreira-Perpiñán and Yerlan Idelbayev. “learning-compression” algorithms for neural net  
370 pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,  
371 June 2018.
- 372 [5] Y Chen, Tien-Ju Yang, Joel Emer, and Vivienne Sze. Understanding the limitations of existing energy-  
373 efficient design approaches for deep neural networks. *Energy*, 2(L1):L3, 2018.
- 374 [6] I. Daubechies, M. Defrise, and C. De Mol. An iterative thresholding algorithm for linear inverse problems  
375 with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57(11):1413–1457, 2004.
- 376 [7] Lei Deng, Guoqi Li, Song Han, Luping Shi, and Yuan Xie. Model compression and hardware acceleration  
377 for neural networks: A comprehensive survey. *Proceedings of the IEEE*, 108(4):485–532, 2020.
- 378 [8] Xiaohan Ding, guiguang ding, Xiangxin Zhou, Yuchen Guo, Jungong Han, and Ji Liu. Global sparse  
379 momentum sgd for pruning very deep neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer,  
380 F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,  
381 volume 32. Curran Associates, Inc., 2019.
- 382 [9] Tausif Diwan, G Anirudh, and Jitendra V Tembhurne. Object detection using yolo: Challenges, architectural  
383 successors, datasets and applications. *multimedia Tools and Applications*, 82(6):9243–9275, 2023.
- 384 [10] Andrei Dumitriu, Florin Tatui, Florin Miron, Radu Tudor Ionescu, and Radu Timofte. Rip current  
385 segmentation: A novel benchmark and yolov8 baseline results. In *Proceedings of the IEEE/CVF Conference*  
386 *on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 1261–1271, June 2023.
- 387 [11] Gongfan Fang, Xinyin Ma, Mingli Song, Michael Bi Mi, and Xinchao Wang. Depgraph: Towards  
388 any structural pruning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
389 *Recognition (CVPR)*, pages 16091–16101, June 2023.
- 390 [12] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural  
391 networks, 2019.
- 392 [13] Yiwen Guo, Anbang Yao, and Yurong Chen. Dynamic network surgery for efficient dnns. In D. Lee,  
393 M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing*  
394 *Systems*, volume 29. Curran Associates, Inc., 2016.
- 395 [14] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural networks with  
396 pruning, trained quantization and huffman coding, 2016.
- 397 [15] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient  
398 neural network. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in*  
399 *Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- 400 [16] Babak Hassibi and David Stork. Second order derivatives for network pruning: Optimal brain surgeon. In  
401 S. Hanson, J. Cowan, and C. Giles, editors, *Advances in Neural Information Processing Systems*, volume 5.  
402 Morgan-Kaufmann, 1992.
- 403 [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
404 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- 405 [18] Andrew Howard, Andrey Zhmoginov, Liang-Chieh Chen, Mark Sandler, and Menglong Zhu. Inverted  
406 residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation. In *CVPR*,  
407 2018.
- 408 [19] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand,  
409 Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile  
410 vision applications, 2017.

- 411 [20] Hengyuan Hu, Rui Peng, Yu-Wing Tai, and Chi-Keung Tang. Network trimming: A data-driven neuron  
412 pruning approach towards efficient deep architectures, 2016.
- 413 [21] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convo-  
414 lutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*  
415 (*CVPR*), July 2017.
- 416 [22] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. Ultralytics yolov8, 2023.
- 417 [23] Minsoo Kang and Bohyung Han. Operation-aware soft channel pruning using differentiable masks. In  
418 Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine*  
419 *Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 5122–5131. PMLR, 13–18  
420 Jul 2020.
- 421 [24] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical  
422 report, 2009.
- 423 [25] Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In D. Touretzky, editor, *Advances in*  
424 *Neural Information Processing Systems*, volume 2. Morgan-Kaufmann, 1989.
- 425 [26] Jaeho Lee, Sejun Park, Sangwoo Mo, Sungsoo Ahn, and Jinwoo Shin. Layer-adaptive sparsity for the  
426 magnitude-based pruning, 2021.
- 427 [27] Namhoon Lee, Thalayasingam Ajanthan, and Philip H. S. Torr. Snip: Single-shot network pruning based  
428 on connection sensitivity, 2019.
- 429 [28] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient  
430 convnets, 2017.
- 431 [29] Yuchao Li, Shaohui Lin, Baochang Zhang, Jianzhuang Liu, David Doermann, Yongjian Wu, Feiyue Huang,  
432 and Rongrong Ji. Exploiting kernel sparsity and entropy for interpretable cnn compression. In *Proceedings*  
433 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- 434 [30] Mingbao Lin, Rongrong Ji, Yan Wang, Yichen Zhang, Baochang Zhang, Yonghong Tian, and Ling Shao.  
435 Hrank: Filter pruning using high-rank feature map. In *Proceedings of the IEEE/CVF Conference on*  
436 *Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- 437 [31] Mingbao Lin, Rongrong Ji, Yuxin Zhang, Baochang Zhang, Yongjian Wu, and Yonghong Tian. Channel  
438 pruning via automatic structure search. *arXiv preprint arXiv:2001.08565*, 2020.
- 439 [32] Shaohui Lin, Rongrong Ji, Chenqian Yan, Baochang Zhang, Liujuan Cao, Qixiang Ye, Feiyue Huang,  
440 and David Doermann. Towards optimal structured cnn pruning via generative adversarial learning. In  
441 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June  
442 2019.
- 443 [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,  
444 and C. Lawrence Zitnick. Microsoft COCO: Common Objects in Context. In David Fleet, Tomas Pajdla,  
445 Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, Lecture Notes in Computer  
446 Science, pages 740–755, Cham, 2014. Springer International Publishing.
- 447 [34] Tantan Liu and Gagan Agrawal. Stratified k-means clustering over a deep web data source. In *Proceedings*  
448 *of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages  
449 1113–1121, 2012.
- 450 [35] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning  
451 efficient convolutional networks through network slimming, 2017.
- 452 [36] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network  
453 pruning, 2019.
- 454 [37] Jian-Hao Luo, Jianxin Wu, and Wei Yao Lin. Thinet: A filter level pruning method for deep neural network  
455 compression. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- 456 [38] Pavlo Molchanov, Arun Mallya, Stephen Tyree, Iuri Frosio, and Jan Kautz. Importance estimation for  
457 neural network pruning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
458 *Recognition (CVPR)*, June 2019.
- 459 [39] Sumit Pandey, Kuan-Fu Chen, and Erik B. Dam. Comprehensive multimodal segmentation in medical  
460 imaging: Combining yolov8 with sam and hq-sam models. In *Proceedings of the IEEE/CVF International*  
461 *Conference on Computer Vision (ICCV) Workshops*, pages 2592–2598, October 2023.

- 462 [40] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,  
463 Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large  
464 Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, December  
465 2015.
- 466 [41] Victor Sanh, Thomas Wolf, and Alexander Rush. Movement pruning: Adaptive sparsity by fine-tuning. In  
467 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information  
468 Processing Systems*, volume 33, pages 20378–20389. Curran Associates, Inc., 2020.
- 469 [42] Maying Shen, Pavlo Molchanov, Hongxu Yin, and Jose M. Alvarez. When to prune? a policy towards  
470 early structural pruning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
471 Recognition (CVPR)*, pages 12247–12256, June 2022.
- 472 [43] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recogni-  
473 tion, 2015.
- 474 [44] Karen Simonyan and Andrew Zisserman. Very Deep Convolutional Networks for Large-Scale Image  
475 Recognition, April 2015. arXiv:1409.1556 [cs].
- 476 [45] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru  
477 Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of  
478 the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- 479 [46] Hidenori Tanaka, Daniel Kunin, Daniel L Yamins, and Surya Ganguli. Pruning neural networks without  
480 any data by iteratively conserving synaptic flow. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan,  
481 and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6377–6389.  
482 Curran Associates, Inc., 2020.
- 483 [47] Lucas Theis, Iryna Korshunova, Alykhan Tejani, and Ferenc Huszár. Faster gaze prediction with dense  
484 networks and fisher pruning, 2018.
- 485 [48] Neil C. Thompson, Kristjan Greenewald, Keeheon Lee, and Gabriel F. Manso. The computational limits of  
486 deep learning, 2022.
- 487 [49] Sunil Vadera and Salem Ameen. Methods for pruning deep neural networks. *IEEE Access*, 10:63280–63300,  
488 2022.
- 489 [50] Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by preserving  
490 gradient flow, 2020.
- 491 [51] Qing Ren Wang and Ching Y. Suen. Analysis and design of a decision tree based on entropy reduction  
492 and its application to large character set recognition. *IEEE Transactions on Pattern Analysis and Machine  
493 Intelligence*, PAMI-6(4):406–417, 1984.
- 494 [52] Shiqiang Wang. Efficient deep learning. *Nature Computational Science*, 1(3):181–182, March 2021.  
495 Number: 3 Publisher: Nature Publishing Group.
- 496 [53] Bichen Wu, Alvin Wan, Xiangyu Yue, Peter Jin, Sicheng Zhao, Noah Golmant, Amir Gholaminejad,  
497 Joseph Gonzalez, and Kurt Keutzer. Shift: A zero flop, zero parameter alternative to spatial convolutions.  
498 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- 499 [54] Kaixin Xu, Zhe Wang, Xue Geng, Min Wu, Xiaoli Li, and Weisi Lin. Efficient joint optimization of  
500 layer-adaptive weight pruning in deep neural networks. In *Proceedings of the IEEE/CVF International  
501 Conference on Computer Vision (ICCV)*, pages 17447–17457, October 2023.
- 502 [55] Pengtao Xu, Jian Cao, Fanhua Shang, Wenyu Sun, and Pu Li. Layer pruning via fusible residual convolu-  
503 tional block for deep neural networks, 2020.
- 504 [56] Tien-Ju Yang, Yu-Hsin Chen, and Vivienne Sze. Designing energy-efficient convolutional neural networks  
505 using energy-aware pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern  
506 Recognition (CVPR)*, July 2017.
- 507 [57] Jiecao Yu, Andrew Lukefahr, David Palframan, Ganesh Dasika, Reetuparna Das, and Scott Mahlke. Scalpel:  
508 Customizing dnn pruning to the underlying hardware parallelism. In *Proceedings of the 44th Annual  
509 International Symposium on Computer Architecture, ISCA '17*, page 548–560, New York, NY, USA, 2017.  
510 Association for Computing Machinery.
- 511 [58] Ruichi Yu, Ang Li, Chun-Fu Chen, Jui-Hsin Lai, Vlad I. Morariu, Xintong Han, Mingfei Gao, Ching-  
512 Yung Lin, and Larry S. Davis. Nisp: Pruning networks using neuron importance score propagation. In  
513 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

- 514 [59] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional  
515 neural network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and*  
516 *Pattern Recognition (CVPR)*, June 2018.
- 517 [60] Zhengguang Zhou, Wengang Zhou, Houqiang Li, and Richang Hong. Online filter clustering and pruning  
518 for efficient convnets. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages  
519 11–15, 2018.
- 520 [61] Zhuangwei Zhuang, Mingkui Tan, Bohan Zhuang, Jing Liu, Yong Guo, Qingyao Wu, Junzhou Huang, and  
521 Jinhui Zhu. Discrimination-aware Channel Pruning for Deep Neural Networks. In S. Bengio, H. Wallach,  
522 H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information*  
523 *Processing Systems*, volume 31. Curran Associates, Inc., 2018.

## 524 A Duality of Independence: Data ( $X$ ) and Information State ( $W_{t_i}$ )

525 Fig. 4 presents a detailed comparative illustration that highlights the similarities in NEXP estimations  
526 across various networks, including VGGNet [44], ResNet [17], MobileNet [19] and DenseNet [21]  
527 on CIFAR-10 dataset. Specifically, for each network architecture, we showcase expressiveness  
528 distributions in both untrained (PaI) and trained (PaT) states. In each sub-figure, the x-axis indicates  
529 convolutional layer indices, and the y-axis shows feature map indices per layer, standardized through  
530 pixel-wise interpolation to align with the layer having the most feature maps. Columns represent  
531 various sampling strategies, while colors indicate expressiveness levels, with higher values signifying  
532 greater expressiveness. In other words, the figure illustrates a two-fold sensitivity analysis of NEXP  
533 to (i) the mini-batch data ( $X$ , as outlined in Alg. 1), using two input sampling strategies to assemble  
534 a batch with 60 samples, namely random sampling (denoted as ‘random’) and class-representative  
535 sampling via k-means (denoted as ‘k-means’), and (ii) the information state ( $W_{t_i}$ ), specifically  
536 comparing expressiveness at initialization (PaI) against expressiveness after training (PaT), when  
537 weights have converged.

### 538 A.1 True NEXP value (non-approx).

539 We define the true NEXP score for each filter as the value obtained by comparing all activation  
540 patterns across the entire training dataset  $D$ . In that way, the ability of each element to extract  
541 maximal features is evaluated for every data-point in the input feature space of a task at hand. In  
542 this study however, due to GPU memory constraints (limited to 12GB of GDDR6 SDRAM), we  
543 employed 25% of the total training set, ensuring class distribution is preserved, to determine these  
544 exact NEXP scores, denoted as non-approx.

### 545 A.2 Data Agnostic.

546 To evaluate NEXP’s sensitivity to input data, we conduct a similarity analysis for each row in  
547 Fig. 4. For each information state (PaI and PaT), we compare the expressiveness map ( $NEXP_{\text{map}}$ )  
548 derived from each sampling strategy against the true NEXP values (non-approx), corresponding to  
549 each respective state. For a comprehensive comparison, we utilize various similarity metrics, such  
550 as Euclidean Distance, Cosine Similarity, Pearsonr Similarity, and the Structural Similarity Index  
551 Measure (ssim\_index). Detailed results of this analysis, specific to each state, are presented in Tables  
552 3 (PaT) and 4 (PaI). The comparative analysis reveals that a mini-batch of 60 samples, with a balanced  
553 representation from each class, effectively approximates the NEXP scores calculated from the entire  
554 dataset, yielding consistent similarity scores above 99% across all similarity metrics for both PaI  
555 and PaT. Interestingly, random sampling consistently outperforms the k-means selection strategy,  
556 which involves selecting 6 representative samples per CIFAR-10 class. This is especially notable in  
557 PaT, with random sampling showing up to a 7.51% higher Pearson correlation, 5% improvement in  
558 ssim\_index, and 1.14 reduction in Euclidean distance compared to k-means. This further reinforces  
559 the statement that comparing activation patterns reflects the intrinsic ability of neural networks to  
560 distinguish various input spaces, thus effectively extending the NEXP criterion to random input data  
561 and laying the foundation for investigating *Data-Agnostic strategies*.

### 562 A.3 Weight Agnostic

563 Fig. 4 reveals that  $NEXP_{\text{map}}$ ’s obtained at initialization and after network convergence share some  
564 expressiveness pattern similarities, particularly in the initial layers. Detailed comparisons of these  
565 similarities across all layers, and specifically for the first five, are presented in Table 5, contrasting  
566 the initial maps with the true  $NEXP_{\text{map}}$  post-training. The summary of our numeric evaluation  
567 confirms a notable correlation between the initialization and converged states for DenseNet-40 and  
568 VGG-19, showing up to 84.10% and 86.82% in cosine similarity respectively. It also indicates  
569 greater consistency in neural expressiveness measurements for the first layers of all networks, which  
570 could be considered important for the formation of critical paths. In this context, the formation  
571 of the final state depends on hyperparameter choices, like weight decay and learning rate, and the  
572 stochastic nature of training, that could potentially alter the model’s progression from its initial state,  
573 as also highlighted by Frankle et al. [12]. In that manner, under the assumption that the selection  
574 of hyperparameters remains congruent with the initialization, “Expressiveness” can be considered a  
575 fit criterion for Pruning at Initialization (PaI). In summary, the consistency of map measurements

576 between initial and final states may serve as a solid metric for evaluating NEXP’s ability to identify  
 577 winning tickets. Nevertheless, a more robust process of its evaluation should also take into account  
 578 the quality of the initial state as well as the subsequent training process.

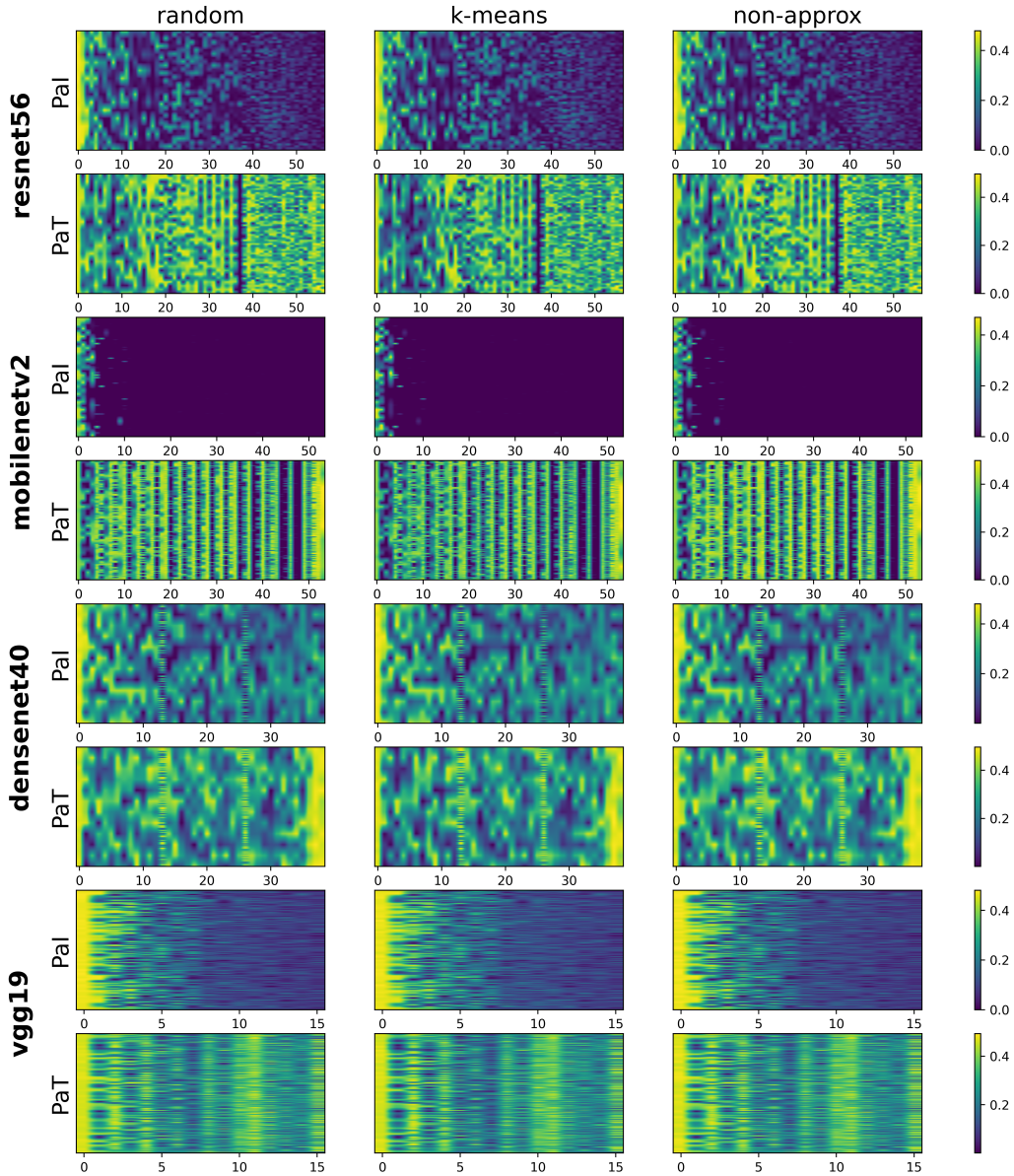


Figure 4: **Expressiveness statistics of feature maps from different convolutional layers and architectures on CIFAR-10 (Extended).** For each architecture we demonstrate the expressiveness distribution for both an untrained instance of the model (PaI), as well as a converged one (PaT). The x-axis represents the indices of convolutional layers and y-axis that of the feature maps in each layer. To maintain consistency across the y-axis, we have interpolated each layer’s feature maps (pixel-wise) to match the layer with the most feature maps. Columns denote different sampling strategies and different colors denote different expressiveness values (the higher the value, the more expressive the feature map). To approximate the expressiveness score of each element, denoted as “non-approx”, we used 25% of all dataset’s samples (not 100% due to memory limitations) maintaining the label’s distribution. As can be seen, the rank of each feature map (column of the sub-figure) is almost unchanged (the same color), regardless of the image batches. Hence, even a small number of images can effectively estimate the average rank of each feature map in different architectures.

Table 3: Sensitivity analysis of the input’s sampling strategies after training (PaT) using various similarity metrics.

Model	Sampling Strategy	Euclidean Distance	Cosine Similarity	Pearsonr Similarity	ssim_index
ResNet-56 [17]	random	<b>0.2349</b>	<b>0.9998</b>	-	<b>0.9979</b>
	k-means	1.3729	0.9949	-	0.9479
MobileNet-v2 [19]	random	<b>0.2903</b>	<b>0.9994</b>	<b>0.9810</b>	<b>0.9988</b>
	k-means	1.1197	0.9960	0.9059	0.9794
DenseNet-40 [21]	random	<b>0.2751</b>	<b>0.9997</b>	<b>0.9818</b>	<b>0.9970</b>
	k-means	1.1669	0.9959	0.9527	0.9614
VGG-19 [44]	random	<b>0.5150</b>	<b>0.9989</b>	<b>0.9814</b>	<b>0.9894</b>
	k-means	0.8438	0.9964	0.9556	0.9728

Table 4: Sensitivity analysis of the input’s sampling strategies at Initialization (PaI) using various similarity metrics.

Model	Sampling Strategy	Euclidean Distance	Cosine Similarity	Pearsonr Similarity	ssim_index
ResNet-56 [17]	random	<b>0.1333</b>	<b>0.9996</b>	<b>0.9979</b>	<b>0.9984</b>
	k-means	0.3948	0.9984	0.9868	0.9859
MobileNet-v2 [19]	random	<b>0.0340</b>	<b>0.9565</b>	-	<b>0.9994</b>
	k-means	0.2441	0.9454	-	0.9776
DenseNet-40 [21]	random	<b>0.2297</b>	<b>0.9997</b>	0.9927	<b>0.9977</b>
	k-means	0.2972	0.9994	<b>0.9941</b>	0.9955
VGG-19 [44]	random	<b>0.2688</b>	<b>0.9988</b>	0.9652	<b>0.9950</b>
	k-means	0.4882	0.9975	<b>0.9724</b>	0.9856

Table 5: Sensitivity analysis of  $NEXP_{\text{map}}$ ’s retrieved at initialization compared with the true  $NEXP_{\text{map}}$  following model convergence.

Model	Metric	random		k-means		non-approx (PaI)	
		All	first-5	All	first-5	All	first-5
ResNet-56 [17]	Euclidean Distance	9.0326	5.2005	<b>8.8029</b>	<b>5.1177</b>	8.9986	5.1850
	Cosine Similarity	0.7584	0.8765	<b>0.7677</b>	<b>0.8784</b>	0.7592	0.8751
	ssim_index	0.0194	0.3794	<b>0.0243</b>	<b>0.3990</b>	0.0206	0.3810
MobileNet-v2 [19]	Euclidean Distance	<b>10.5470</b>	<b>7.4966</b>	10.6056	8.0843	10.5492	7.5134
	Cosine Similarity	0.4645	<b>0.6478</b>	0.4910	0.5862	<b>0.6702</b>	0.6461
	ssim_index	-0.0018	<b>0.1187</b>	-0.0011	0.0942	-0.0020	0.1142
DenseNet-40 [21]	Euclidean Distance	6.1326	<b>4.6957</b>	<b>6.0594</b>	4.7157	6.1043	4.7364
	Cosine Similarity	0.8357	<b>0.8769</b>	<b>0.8410</b>	0.8762	0.8378	0.8761
	ssim_index	<b>0.0169</b>	<b>0.4552</b>	0.0101	0.4493	0.0150	0.4464
VGG-19 [44]	Euclidean Distance	6.3171	4.9532	<b>6.1194</b>	<b>4.8525</b>	6.3083	4.9810
	Cosine Similarity	0.8610	0.8979	<b>0.8682</b>	<b>0.9030</b>	0.8624	0.8972
	ssim_index	0.0808	0.3798	<b>0.0812</b>	<b>0.3844</b>	0.0808	0.3712



## 579 **B Pruning Process: An in-depth analysis**

### 580 **B.1 Global vs local -scope pruning.**

581 NEXP is used in the pruning process to evaluate and rank different network elements, guiding their  
582 subsequent removal based on their scores. In our study, we focused on the removal of filters, i.e.,  
583 Filter Pruning, where we pruned convolutional structures by removing the least expressive filters.  
584 This can be approached in two ways: (i) on a local (layer-by-layer) basis, where filters are assessed  
585 and removed according to their expressiveness relative to other filters within the same layer, e.g.,  
586 eliminating the least  $\mu$  expressive filters from each layer. (ii) On a global (network-wide) basis, where  
587 all filters across layers are normalized in terms of their scores, allowing for the removal of the least  
588  $\kappa$  expressive filters from the entire network. We experimentally observed that "Global Pruning"  
589 yields consistent results and outperforms "Local Pruning" when using the NEXP pruning criterion.  
590 Therefore, all the experiments reported in this paper were conducted using the "Global Pruning"  
591 approach.

### 592 **B.2 One-shot vs Iterative pruning.**

593 Furthermore, another design parameter to consider in the pruning process is its coordination with  
594 fine-tuning. In this context, two widely adopted strategies are: (a) "One-Shot" pruning, where pruning  
595 is completed entirely before any fine-tuning occurs, and (b) "Iterative" pruning, which involves  
596 alternating between pruning and fine-tuning via an iterative sequence. The first one (a) can be  
597 considered a more lightweight approach and allows for a more robust evaluation of the pruning metric  
598 at hand, when compared to the later one (b). This is because it has no extra dependency on the training  
599 data and its efficiency does not depend on the iterative re-calibration of the information state through  
600 the fine-tuning process. In this study, most experiments were conducted using "One-Shot" pruning,  
601 while we also explored the integration of NEXP in an "Iterative" pruning process with YOLOv8  
602 (more details on 4.1), where we noted a reduction in performance declines and an improvement in  
603 the performance recovery after each pruning step, leading to better overall results.

### 604 **B.3 Detailed description of all algorithmic steps.**

605 More in detail regarding Algorithm 1, we define the data structure  $NEXP_{\text{map}}$ , i.e., a dictionary  
606 in our implementation, to store the NEXP scores for every filter in the neural network after each  
607 iteration. Given a neural network  $\mathcal{N}$  with its current weight state  $W_{t_i}$ , we initially set up all variables  
608 required for the pruning loop (Lines 1-4). The network is then gradually pruned until one of the  
609 following conditions is met: the target ratio is achieved or the allowed number of pruning steps  
610 is exceeded (Line 5). During each pruning iteration, the  $\kappa$  least expressive filters from the current  
611 pruned state of the network are initially selected (Line 6). These filters are then removed, followed  
612 by an update to  $NEXP_{\text{map}}$  for the subsequent iteration (Lines 7-8). To obtain the NEXP scores,  
613 a forward pass  $f(X; W_{\text{pruned}})$  is conducted using a user-provided mini-batch as input. Finally, the  
614 conditions variables are updated in preparation for the next pruning iteration (Lines 9-10).

### 615 **B.4 Acceleration of NEXP computations.**

616 In Algorithm 1, Line 8 accounts for the bulk of the computational complexity. Specifically, the  
617 calculation of  $NEXP_{\text{map}}$  can be divided into two sub-processes: (i) performing a forward pass to  
618 retrieve all activation patterns, and (ii) estimating the NEXP score for each element in the map.  
619 However, performing a forward pass can be considered negligible compared to computing the NEXP  
620 score for each filter. This is because the later involves multiple comparisons between the activation  
621 patterns of all samples in the mini-batch  $X$  for every filter. Two effective ways to reduce this  
622 computational demand are: first, all operations involved in computing the NEXP score are compatible  
623 with widely-used BLAS libraries, facilitating hardware acceleration; second, the frequency of score  
624 updates can be strategically decreased under certain conditions, e.g., every  $n$  pruning iterations.

## 625 C Experimental Settings

### 626 C.1 Datasets and Models.

627 This paper explores Computer Vision tasks through extensive experiments on various datasets, such  
628 as CIFAR-10 [24] and ImageNet [40] for image classification, and COCO [33] for object detection.  
629 To demonstrate the robustness of our approach, we experiment on several popular architectures and a  
630 wide span of architectural elements, including VGGNet with a plain structure [44], ResNet with a  
631 residual structure [17], GoogLeNet with inception modules [45], MobileNet with depthwise separable  
632 convolutions [19], DenseNet with dense blocks [21] and YOLOv8 with a variety of different modules,  
633 e.g. C2f and SPPF [22].

### 634 C.2 Adversaries.

635 We assess the efficacy of expressiveness as criterion for Pruning both after Training (PaT) and at  
636 Initialization (PaI), using arbitrary (random) data-points. For PaT (4.1), we compare against a plethora  
637 of foundational and state-of-the-art approaches, ranging from filter magnitude-based [28, 32, 29]  
638 and loss sensitivity-based [58] methods to feature-guided strategies [23, 30] and search algorithms  
639 [35, 31]. Regarding PaI (4.3 and D.1), our comparison is two-fold, as we evaluate expressiveness  
640 using (i) single-shot and (ii) iterative pruning. More specifically, the adversaries for PaI include  
641 pruning with random scoring, two state-of-the-art single-shot pruning strategies, namely SNIP [27]  
642 and GraSP [50], as well as one state-of-the-art iterative pruning strategy, named SynFlow [46].

### 643 C.3 Evaluation Metrics.

644 To effectively quantify the efficiency of reported solutions, we adopt a 3-dimensional evaluation  
645 space, consisting of i) two widely-used metrics i.e. *FLOPs* and *params*, that define the 2-dimensional  
646 compression solution efficiency, alongside with ii) an NN model accuracy to assess the predictions  
647 of pruned derivatives [3]. Within the compression space, we define, (a) Compression Ratio( $\downarrow$ ) =  
648  $\frac{\text{original size}}{\text{compressed size}}$  and (b) Compressed Size Percentage (%) =  $\frac{\text{compressed size}}{\text{original size}} \cdot 100$ . To assess task-specific  
649 capabilities, we report the top-1 accuracy of pruned models for image classification on CIFAR-10  
650 [24], both top-1 and top-5 accuracies for ImageNet [40], and the mean Average Precision (mAP) over  
651 IoU (Intersection over Union) thresholds ranging from 0.5 to 0.95, denoted as  $mAP_{50-95}^{val}$ , for object  
652 detection on the COCO dataset [33].

### 653 C.4 Configurations.

654 We implement the proposed “expressiveness” pruning criterion on PyTorch, version 2.0.1+cu117, by  
655 extending the DepGraph pruning framework [11] to maintain models compatibility and to ensure  
656 structural coupling during the removal of network elements e.g., simultaneously removing any inter-  
657 dependent network elements such as kernel pairs of convolutional and batch-normalization batched  
658 layers. All experiments are conducted on a NVIDIA GeForce RTX 3060 GPU with 12GB of GDDR6  
659 SDRAM. For all experiments we use a batch of 64 random data-points to estimate expressiveness,  
660 except those that are reported for CIFAR-10 and ImageNet on 4.1, where we used K-Means to select  
661 60 samples (6 from each class). Additionally, the baseline models on CIFAR-10 were trained for 200  
662 epochs by using 128 batch size and Stochastic Gradient Descent algorithm (SGD) with an initial  
663 learning rate of 0.1 that is divided by 10 after 60 and 120 epochs respectively. For ImageNet models  
664 and YOLOv8, we utilize the available pre-trained weights on PyTorch vision library and ultralytics  
665 [22]. We fine-tune the pruned networks for 100 epochs on CIFAR-10 and for 30 epochs on ImageNet  
666 to compensate for the performance loss, using a batch size of 128 and 32 respectively.

## 667 D Supplementary Experimental Results

### 668 D.1 Neural Expressiveness at Initialization: A comparative study

669 **Adversaries.** We establish our comparative study in a two-fold manner, as we compare expressiveness  
670 against (i) single-shot and (ii) iterative pruning approaches. More specifically, the adversaries include  
671 pruning with random scoring, two state-of-the-art single-shot pruning strategies, namely SNIP [27]

672 and GraSP [50], as well as one state-of-the-art iterative pruning strategy, named SynFlow [46]. For  
 673 our approach, we implement one-shot pruning, utilizing a batch of 64 arbitrary data points for the  
 674 estimation of expressiveness.

675 **Experimental Setup.** We adopt the experimental framework of Tanaka et al. [46], who assess  
 676 algorithm performance across an exponential scale ( $10^r$ ) of parameters compression ratios  $r \in$   
 677  $\{0.00, 0.25, 0.50, 0.75, \dots\}$ . Their proposed settings also enable for the evaluation of an algorithm's  
 678 resilience to "layer collapses", typically observed at higher compression levels. **Results.** We prune  
 679 VGG-16 on CIFAR-10 and compare against the findings of [46]. We remain consistent with our  
 680 adversaries and train the model for 160 epochs, using a batch size of 128 and an initial learning rate  
 681 of 0.1, which is reduced by a factor of 10 after 60 and 120 epochs. The results are illustrated on  
 682 Fig. 5.

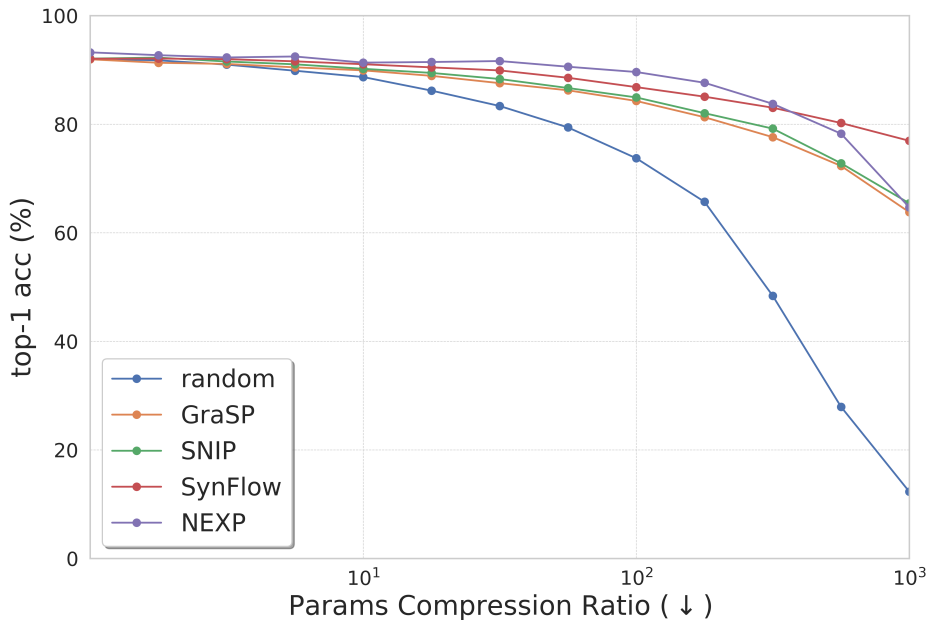


Figure 5: **Pruning VGG-16 at Initialization on CIFAR-10.** A comparative visualisation of SOTA methods across an exponential scale of params compression ratios.

683 **Observations.** Our method consistently outperforms all other algorithms, particularly in regimes of  
 684 lower compression, up to  $10^2(\downarrow)$  with an average increase of 1.21% over SynFlow, while maintaining  
 685 competitiveness at higher compression levels, above  $10^2(\downarrow)$  with an average percentage difference of  
 686 4.82%, 3.72% and -2.74%, compared to GraSP, SNIP and SynFlow respectively.

## 687 D.2 Additional Experimental Results: Tables and Figures

688 **CIFAR-10.** We present further experiments and comparisons with state-of-the-art methods,  
 689 including HRANK [30], GAL [32], ABC [31] and DCP [61], specifically for GoogLeNet and  
 690 MobileNet-v2 networks. For MobileNet-v2, our method attains an increased compression ratio of  
 691  $0.94\times$  in parameters and  $0.75\times$  in FLOPs ( $\downarrow$ ), with a minimal decrease of only -0.09% in performance  
 692 compared to DCP. In the GoogLeNet case, we demonstrate a notable enhancement in parameters  
 693 compression within the  $1.60\times$  to  $2.20\times$  FLOPs compression range, surpassing GAL and HRANK  
 694 with margins of  $1.8\times$  and  $1.52\times$  respectively, with an average improvement of 7.5% in performance  
 695 degradation.

Table 6: Analytical Comparison of Importance-based solutions and Expressiveness on CIFAR-10 using VGGNet architectures [44].

Model	Method	top-1 acc		Compression Ratio ↓	
		Base (%)	$\Delta$ (%)	#Params	#FLOPs
VGG-16	L1 [28]	93.25	+0.15	2.78×	1.52×
	GAL-0.05 [32]	93.96	-0.19	4.46×	1.65×
	GAL-0.1 [32]	93.96	-0.54	5.61×	1.82×
	HRank [30]	93.96	-0.53	5.97×	2.15×
	HRank [30]	93.96	-1.62	5.67×	2.89×
	SCP [23]	93.85	-0.06	15.38×	2.96×
	NEXP (Ours)	93.87	-0.16	5.62×	3.03×
	ABC [31]	93.02	+0.06	8.80×	3.80×
	NEXP (Ours)	93.87	-0.35	<b>13.13</b> ×	4.01×
	HRank [30]	93.96	-2.73	8.41×	4.26×
VGG-19	DCP-Adapt [61]	93.99	+0.58	15.58×	2.86×
	SCP [23]	93.84	-0.02	20.88×	3.86×
	NEXP (Ours)	94.00	-0.53	<b>22.73</b> ×	4.75×

Table 7: Analytical Comparison of Importance-based solutions and Expressiveness on CIFAR-10 using GoogLeNet [45].

Model	Method	top-1 acc		Compression Ratio ↓	
		Base (%)	$\Delta$ (%)	#Params	#FLOPs
GoogLeNet	GAL-0.5 [32]	95.05	-0.49	1.97×	1.62×
	NEXP (Ours)	94.97	-0.43	<b>3.77</b> ×	2.12×
	Hrank [30]	95.05	-0.52	2.25×	2.20×
	ABC [31]	95.05	-0.21	2.51×	2.99×
	NEXP (Ours)	94.97	-1.07	<b>7.02</b> ×	3.01×
	Hrank [30]	95.05	-0.98	3.31×	3.38×

Table 8: Analytical Comparison of Importance-based solutions and Expressiveness on CIFAR-10 using DenseNet-40 [21].

Model	Method	top-1 acc		Compression Ratio ↓	
		Base (%)	$\Delta$ (%)	#Params	#FLOPs
DenseNet-40	GAL-0.5 [32]	95.05	-0.49	1.97×	1.62×
	Hrank [30]	95.05	-0.52	2.25×	2.20×
	NEXP (Ours)	94.64	-0.89	<b>2.72</b> ×	2.25×
	NEXP (Ours)	94.64	-0.84	<b>3.12</b> ×	2.51×
	ABC [31]	95.05	-0.21	2.51×	2.99×
	Hrank [30]	95.05	-0.98	3.31×	3.38×

Table 9: Performance Outcomes for MobileNet-v2 on the CIFAR-10 Dataset.

Method	Base (%)	$\Delta$ Acc (%)	#Params ↓	#FLOPs ↓
DCP [61]	94.47	+0.22	1.31×	1.36×
NEXP (Ours)	94.32	+0.13	2.25×	2.11×

696 **YOLOv8.** Figure 6 compares Neural Expressiveness (NEXP) with Layer-Adaptive Magnitude-  
 697 Based Pruning (LAMP) [26], Network Slimming (SLIM) [35], Wang et al.'s DepGraph [11], and  
 698 Random Pruning for Object Detection on the COCO dataset, as discussed in 4.1.

699 **Motivation.** YOLOv8 [22] is the current state-of-the-art for Object Detection and Image Segmentation,  
 700 and has already been widely adopted by many for a variety of real-time applications, e.g.  
 701 Traffic Safety [1], Medical Imaging [39], Rip Currents Detection [10], and more. Such applications  
 702 could majorly benefit from model compression optimizations, achieving higher throughput ratios that  
 703 translate to increased resolution (FPS), and enabling deployment on hardware with strict resource  
 704 constraints.

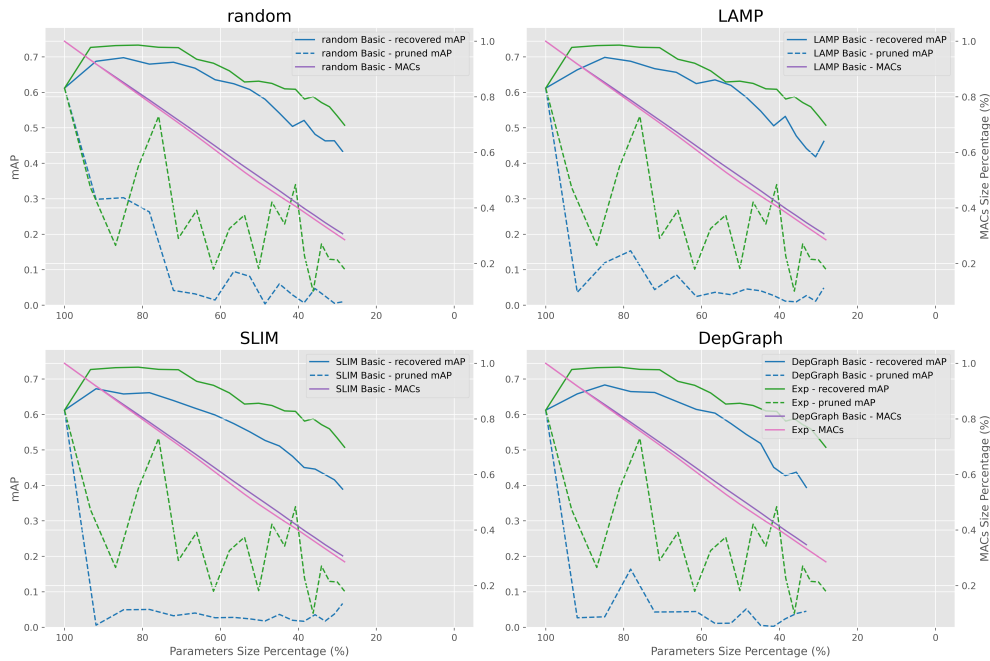


Figure 6: Pruning YOLOv8m trained on COCO for Object Detection.

705 **NeurIPS Paper Checklist**

706 **1. Claims**

707 Question: Do the main claims made in the abstract and introduction accurately reflect the  
708 paper's contributions and scope?

709 Answer: [Yes]

710 Justification: The main contributions have been reflected and discussed across the whole  
711 paper.

712 Guidelines:

- 713 • The answer NA means that the abstract and introduction do not include the claims  
714 made in the paper.
- 715 • The abstract and/or introduction should clearly state the claims made, including the  
716 contributions made in the paper and important assumptions and limitations. A No or  
717 NA answer to this question will not be perceived well by the reviewers.
- 718 • The claims made should match theoretical and experimental results, and reflect how  
719 much the results can be expected to generalize to other settings.
- 720 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
721 are not attained by the paper.

722 **2. Limitations**

723 Question: Does the paper discuss the limitations of the work performed by the authors?

724 Answer: [Yes]

725 Justification: While our work does not implicitly provide a Discussion section, we have  
726 incorporated any discussions on the limitations and the intricacies of the provided solution  
727 at its section separately.

728 Guidelines:

- 729 • The answer NA means that the paper has no limitation while the answer No means that  
730 the paper has limitations, but those are not discussed in the paper.
- 731 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 732 • The paper should point out any strong assumptions and how robust the results are to  
733 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
734 model well-specification, asymptotic approximations only holding locally). The authors  
735 should reflect on how these assumptions might be violated in practice and what the  
736 implications would be.
- 737 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
738 only tested on a few datasets or with a few runs. In general, empirical results often  
739 depend on implicit assumptions, which should be articulated.
- 740 • The authors should reflect on the factors that influence the performance of the approach.  
741 For example, a facial recognition algorithm may perform poorly when image resolution  
742 is low or images are taken in low lighting. Or a speech-to-text system might not be  
743 used reliably to provide closed captions for online lectures because it fails to handle  
744 technical jargon.
- 745 • The authors should discuss the computational efficiency of the proposed algorithms  
746 and how they scale with dataset size.
- 747 • If applicable, the authors should discuss possible limitations of their approach to  
748 address problems of privacy and fairness.
- 749 • While the authors might fear that complete honesty about limitations might be used by  
750 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
751 limitations that aren't acknowledged in the paper. The authors should use their best  
752 judgment and recognize that individual actions in favor of transparency play an impor-  
753 tant role in developing norms that preserve the integrity of the community. Reviewers  
754 will be specifically instructed to not penalize honesty concerning limitations.

755 **3. Theory Assumptions and Proofs**

756 Question: For each theoretical result, does the paper provide the full set of assumptions and  
757 a complete (and correct) proof?

758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812

Answer: [Yes]

Justification: To the best of our knowledge, all the provided set of assumptions presented in Section 3 are complete.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

**4. Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The paper specifies in great detail all the information necessary to understand the results, while also any subjectivity imposed by the experimental settings in regards to our claims and conclusions has been discussed. Detailed analysis of both the mathematical, technical and experimental intricacies have been included in our work.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: We plan to release the full code of the implementation and experiments upon acceptance.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The paper specifies all the training and test details necessary to understand the results. Each experiment is accompanied by a discussion of its experimental details and a reference to its experimental settings (Section 4). Additionally, an overview of the experiment settings can be found in Appendix C, and an in-depth analysis of the pruning procedure, including its implementation choices, is described in Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: While we do not explicitly address the statistical significance of each experiment (in a quantitative manner), we do discuss in great detail any assumptions or statistical implications of our experiments (in a qualitative manner).

Guidelines:

- The answer NA means that the paper does not include experiments.



- 865 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
866 dence intervals, or statistical significance tests, at least for the experiments that support  
867 the main claims of the paper.
- 868 • The factors of variability that the error bars are capturing should be clearly stated (for  
869 example, train/test split, initialization, random drawing of some parameter, or overall  
870 run with given experimental conditions).
- 871 • The method for calculating the error bars should be explained (closed form formula,  
872 call to a library function, bootstrap, etc.)
- 873 • The assumptions made should be given (e.g., Normally distributed errors).
- 874 • It should be clear whether the error bar is the standard deviation or the standard error  
875 of the mean.
- 876 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
877 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
878 of Normality of errors is not verified.
- 879 • For asymmetric distributions, the authors should be careful not to show in tables or  
880 figures symmetric error bars that would yield results that are out of range (e.g. negative  
881 error rates).
- 882 • If error bars are reported in tables or plots, The authors should explain in the text how  
883 they were calculated and reference the corresponding figures or tables in the text.

## 884 8. Experiments Compute Resources

885 Question: For each experiment, does the paper provide sufficient information on the com-  
886 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
887 the experiments?

888 Answer: [Yes]

889 Justification: While we do not provide the exact times of executions and memory require-  
890 ments for each experiment, we do provide an in-depth analysis of all the parameters and  
891 experimental specifications, along side with the overview of the configurations that were  
892 used for this work Appendix C and Section 4.

893 Guidelines:

- 894 • The answer NA means that the paper does not include experiments.
- 895 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
896 or cloud provider, including relevant memory and storage.
- 897 • The paper should provide the amount of compute required for each of the individual  
898 experimental runs as well as estimate the total compute.
- 899 • The paper should disclose whether the full research project required more compute  
900 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
901 didn't make it into the paper).

## 902 9. Code Of Ethics

903 Question: Does the research conducted in the paper conform, in every respect, with the  
904 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

905 Answer: [Yes]

906 Justification: We have thoroughly reviewed the research conducted in the paper and fully  
907 agree that it conforms to the NeurIPS Code of Ethics.

908 Guidelines:

- 909 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 910 • If the authors answer No, they should explain the special circumstances that require a  
911 deviation from the Code of Ethics.
- 912 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
913 eration due to laws or regulations in their jurisdiction).

## 914 10. Broader Impacts

915 Question: Does the paper discuss both potential positive societal impacts and negative  
916 societal impacts of the work performed?

917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968

Answer: [NA]

Justification: While our work does not directly discuss societal impacts, we do reference the eco-friendly implications of efficient models in the Introduction section 1. Additionally, we highlight the potential indirect societal benefits that can arise from optimizing models, such as in the case of YOLOv8 D.2. .

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

**11. Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our work emphasizes on efficiently compressing Neural Networks and does not target any specific use-case scenario, rather it addresses the greater challenge of Vision as whole.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

**12. Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

969 Justification: All the creators and original owners of the assets that were utilized for this  
970 work were properly credited through-out all parts of the paper, while also a detailed report  
971 of them can be found in Appendix C.

972 Guidelines:

- 973 • The answer NA means that the paper does not use existing assets.
- 974 • The authors should cite the original paper that produced the code package or dataset.
- 975 • The authors should state which version of the asset is used and, if possible, include a  
976 URL.
- 977 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 978 • For scraped data from a particular source (e.g., website), the copyright and terms of  
979 service of that source should be provided.
- 980 • If assets are released, the license, copyright information, and terms of use in the  
981 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
982 has curated licenses for some datasets. Their licensing guide can help determine the  
983 license of a dataset.
- 984 • For existing datasets that are re-packaged, both the original license and the license of  
985 the derived asset (if it has changed) should be provided.
- 986 • If this information is not available online, the authors are encouraged to reach out to  
987 the asset’s creators.

### 988 13. New Assets

989 Question: Are new assets introduced in the paper well documented and is the documentation  
990 provided alongside the assets?

991 Answer: [NA]

992 Justification: The paper does not release new assets besides the conceptualization and both  
993 then technical and theoretical formulation of Neural Expressiveness. However, we plan to  
994 release the full code of the implementation and experiments upon acceptance.

995 Guidelines:

- 996 • The answer NA means that the paper does not release new assets.
- 997 • Researchers should communicate the details of the dataset/code/model as part of their  
998 submissions via structured templates. This includes details about training, license,  
999 limitations, etc.
- 1000 • The paper should discuss whether and how consent was obtained from people whose  
1001 asset is used.
- 1002 • At submission time, remember to anonymize your assets (if applicable). You can either  
1003 create an anonymized URL or include an anonymized zip file.

### 1004 14. Crowdsourcing and Research with Human Subjects

1005 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1006 include the full text of instructions given to participants and screenshots, if applicable, as  
1007 well as details about compensation (if any)?

1008 Answer: [NA]

1009 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

1010 Guidelines:

- 1011 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1012 human subjects.
- 1013 • Including this information in the supplemental material is fine, but if the main contribu-  
1014 tion of the paper involves human subjects, then as much detail as possible should be  
1015 included in the main paper.
- 1016 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1017 or other labor should be paid at least the minimum wage in the country of the data  
1018 collector.

### 1019 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 1020 Subjects

1021 Question: Does the paper describe potential risks incurred by study participants, whether  
1022 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1023 approvals (or an equivalent approval/review based on the requirements of your country or  
1024 institution) were obtained?

1025 Answer: [NA]

1026 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

1027 Guidelines:

- 1028 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1029 human subjects.
- 1030 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1031 may be required for any human subjects research. If you obtained IRB approval, you  
1032 should clearly state this in the paper.
- 1033 • We recognize that the procedures for this may vary significantly between institutions  
1034 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1035 guidelines for their institution.
- 1036 • For initial submissions, do not include any information that would break anonymity (if  
1037 applicable), such as the institution conducting the review.