

Interspace Pruning: Using Adaptive Filter Representations to Improve Training of Sparse CNNs

Paul Wimmer^{*†‡§}, Jens Mehnert^{*§} and Alexandru Paul Condurache^{*†}

^{*}Automated Driving Research, Robert Bosch GmbH, 70469 Stuttgart, Germany

[†]Institute for Signal Processing, University of Lübeck, 23562 Lübeck, Germany

{paul.wimmer, jensericmarkus.mehnert, alexandrupaul.condurache}@de.bosch.com

Abstract

Unstructured pruning is well suited to reduce the memory footprint of convolutional neural networks (CNNs), both at training and inference time. CNNs contain parameters arranged in $K \times K$ filters. Standard unstructured pruning (SP) reduces the memory footprint of CNNs by setting filter elements to zero, thereby specifying a fixed subspace that constrains the filter. Especially if pruning is applied before or during training, this induces a strong bias. To overcome this, we introduce interspace pruning (IP), a general tool to improve existing pruning methods. It uses filters represented in a dynamic interspace by linear combinations of an underlying adaptive filter basis (FB). For IP, FB coefficients are set to zero while un-pruned coefficients and FBs are trained jointly. In this work, we provide mathematical evidence for IP's superior performance and demonstrate that IP outperforms SP on all tested state-of-the-art unstructured pruning methods. Especially in challenging situations, like pruning for ImageNet or pruning to high sparsity, IP greatly exceeds SP with equal runtime and parameter costs. Finally, we show that advances of IP are due to improved trainability and superior generalization ability.

1. Introduction

Deep neural networks (DNNs) have shown state-of-the-art (SOTA) performance in many artificial intelligence applications [60, 63, 83, 89, 92]. In order to solve these tasks, large models with up to billions of parameters are required. However, training, transferring, storing and evaluating such large models is costly [70, 74]. Pruning [23, 26, 28, 36, 40, 57] sets parts of the network's weights to zero. This reduces the model's complexity and memory requirements, speeds up inference [6] and may lead to an improved generalization ability [4, 30, 40]. In recent years, training sparse models became of interest, providing the benefits of reduced memory

requirements and runtime not only for inference but also for training [17, 18, 41, 53, 55, 64, 76, 82, 86].

In this work, we mainly focus on methods that prune individual parameters *before* training, while the number of zeroed coefficients is kept fixed during training. With this *unstructured pruning*, a network's memory footprint can be reduced. To lower the runtime in addition, specialized soft- and hardware is needed [15, 22, 27, 59]. For training sparse networks, we distinguish between (i) *pruning at initialization* (PaI) [12, 41, 76, 82, 86] which prunes the network at initialization and fixes zeroed parameters during training, (ii) finding the sparse architecture to be finally trained by iterative train-prune-reset cycles, a so called *lottery ticket* (LT) [18, 19], and (iii) *dynamic sparse training* (DST) [17, 45, 55] which prunes the network at initialization, but allows the pruning mask to be changed during training.

Convolutional neural networks (CNNs) are composed of layers, each having a certain number of input- and output channels. Every combination of input- and output channel is linked by a filter $h \in \mathbb{R}^{K \times K}$ with kernel size $K \times K$. A weight of h is a *spatial coefficient* $h_{i,j}$ for a spatial coordinate (i, j) . Filters h can also be modeled in an *interspace*, a linear space $\{\sum_{n=1}^{K^2} \lambda_n \cdot g^{(n)} : \lambda_n \in \mathbb{R}\}$ spanned by a *filter basis* (FB) $\mathcal{F} := \{g^{(1)}, \dots, g^{(K^2)}\} \subset \mathbb{R}^{K \times K}$ [16, 79]. One possibility for a FB is the standard basis $\mathcal{B} := \{e^{(n)} : n = 1, \dots, K^2\}$ which yields the spatial representations. General interspace representations are more flexible since bases are not fixed. We represent h in an interspace in order to learn the FB \mathcal{F} spanning this space along with the *FB coefficients* λ , and thereby obtain a better representation for h . Thus, setting coefficients of flexible, adaptive FBs to zero will improve results compared to prune spatial coefficients.

For deep networks, where the layers' purposes are usually unknown to the experts but learnt during training, we believe that filters should train their bases along with their coefficients. A FB \mathcal{F} is *dynamic*, can be shared for any number of $K \times K$ filters, and is optimized jointly with its FB coefficients λ . By fitting an interspace to sparse filters

[‡]Corresponding author. [§]Equal contribution.

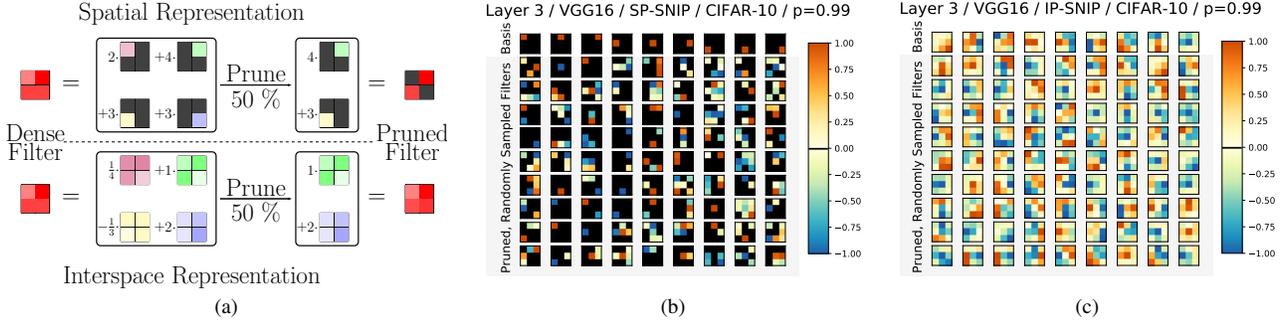


Figure 1. (a) Overview of SP and IP. Contrarily to SP (b), IP (c) produces spatially dense filters after training sparse networks. As for SP, sparsity in the interspace can be used to reduce memory requirements and, by the linearity of convolutions, also computational costs.

during training, we overcome the lack of prior knowledge for a basis that is well suited to describe filters with few non-zero coefficients. If a filter is pruned to a single FB coefficient, $h = \lambda_n \cdot g^{(n)}$, it is not restricted since $g^{(n)}$ can change. Thus, pruning interspace coefficients of dynamic FBs keeps the CNN flexible and is called *interspace pruning* (IP). A 1-sparse filter $h = h_{i_n, j_n} \cdot e^{(n)}$ directly predefines h to stay on the fixed subspace $\text{span}\{e^{(n)}\}$. Pruning spatial coefficients w.r.t. the standard basis \mathcal{B} is called *standard pruning* (SP).

During training sparse CNNs, the problem of vanishing gradients due to spatial sparsity often occurs [76, 82, 85]. In contrast, IP pruned networks are able to learn spatially dense FBs during training, even when using sparse interspace coefficients, see Fig. 1. Therefore, IP leads to an improved information flow and better trainable models.

Although IP yields dense spatial representations, the linearity of convolutions can be used to reduce the number of computations for CNNs with sparse interspace coefficients. Compared to SP, IP only increases the number of required computations by a small, constant count. However, as IP provides superior sparse models, IP generates CNNs with faster inference speed than SP while matching the dense performance. Further, the dynamic achieved by interspace representations is cheap in terms of memory. A FB \mathcal{F} has K^4 parameters as it contains K^2 filters of size $K \times K$. A single FB can be shared for all $K \times K$ filters in a CNN. Also, more than one FB can be used with just a small increase in memory requirements. For cost reasons, we do not use more FBs than the number of layers in a CNN in our experiments, resulting in *all* FBs creating an overhead of at most 0.01% of the dense network’s parameters. Despite adding only few additional costs compared to using spatial weights, interspace representations significantly improve results for sparse *and* dense training.

Our core contributions are:

- Representing and training convolutional filters in the interspace, a linear space spanned by a trainable FB. The FB is optimized jointly with the FB coefficients.

- Formulating the concept of pruning for filters with interspace representation as general method to improve performance of CNNs with sparse coefficients.
- Theoretical proof of IP’s improvements in Thm. 1.
- Experiments showing that IP exceeds SP for equal runtime and memory costs on SOTA sparse training methods and pruning methods which are applied during training or on pre-trained models. We demonstrate that IP’s superiority is achieved by improved trainability, and at lower sparsity also due to better generalization.

2. Broader Impact

Pruning can lower costs for training, storing and evaluating DNNs. We are not aware of any negative outcome directly induced by this work. Nevertheless, as tool to improve pruning, and therefore to reduce costs for CNNs, IP could be used for any CNN based application with negative ethical or societal impact. As authors, we distance ourselves from such applications and the use of our method therein.

As we show in the paper, IP improves unstructured pruning in general and is not restricted to a special scenario. We see IP as a tool which is applied in combination with SOTA SP techniques to lower costs further. Consequently, our work is to the advantage of everyone using pruning and, by the improved generalization ability obtained by training with interspace representations, deep learning in general.

3. Related work

Related work covers *general pruning* and *pruning before training and DST*. Training a sparse model allows to learn non-zero FB coefficients and FBs jointly from scratch. Such methods naturally benefit most from interspace representations and our experiments thus place a strong focus on them.

General pruning. Pruning is divided in *structured* and *unstructured* pruning. Structured pruning removes coarse structures of the network, like channels or neurons [3, 33,

44, 77, 84, 91]. This yields lean architectures and thus reduces computation time. A more fine-grained approach is unstructured pruning where single, spatial weights are zeroed [18, 21, 26, 37, 40, 41, 55]. Unstructured pruning leads to better performance than structured pruning [42, 54] but requires soft- and hardware that supports sparse tensor computations to actually reduce runtime [27, 59]. Also, storing sparse parameters in formats such as the *compressed sparse row format* [78] creates additional overhead. This can lead to non-linear dependencies between the sparsity and actual memory/runtime costs, see also Appendix Secs. C and D.4.

Pruning can be applied at any time in training. The historically first approaches [29, 36, 37, 40, 57] use trained networks and many prune and fine-tune cycles. Criteria are often based on expensive computations of the Hessian w.r.t. the loss function. Likewise, magnitudes of trained coefficients can be used as iterative pruning criterion [21, 28, 42]. By adding sparsity forcing regularizations to the loss, pruning can be integrated dynamically into training [7, 48, 88].

Closest to our work are pruning coefficients in the frequency [47] and the Winograd domain [46]. Contrarily, we do not bind representations to a fixed basis but let the network learn its FBs self-reliantly. Moreover, IP is not a pruning method by itself, but is added on top of existing ones to boost them. Also to mention is [43], a low rank approximation of CNNs. A dense, pre-trained network is approximated by learning undercomplete dictionaries for 3D filters. We, on the contrary, represent 2D filters $h \in \mathbb{R}^{K \times K}$, prune the network instead of using low rank approximations and learn FBs jointly with the coefficients in one training.

Pruning before training and dynamic sparse training.

In [18], an iterative procedure is proposed which consists of training un-pruned weights to convergence, applying magnitude pruning with a small pruning rate to the trained weights and resetting the non-zero weights to their initial value. Finally, this leads to sparse, randomly initialized networks which are well trainable – so called lottery tickets. For SOTA CNNs, resetting un-pruned weights not to the initialization but a value from an early iteration improves performance significantly [19, 68]. By applying other criteria, like *information flow* in the sparse network [62, 76, 82, 86] or influence of non-zero weights on changing the loss [12, 41, 81, 86], pruning can be successfully applied at initialization without pre-training the network. GraSP [82], SNIP [41] and SynFlow [76] are SOTA for PaI [20]. Dynamic sparse training [13, 17, 45, 55] adjusts pruning masks during training to ensure sparse networks while adapting the architecture to different conditions. SET [55] frequently prunes the network based on magnitudes and activates as many un-trained parameters randomly. RigL [17] improves this by recovering those weights with the biggest gradient magnitude.

4. Filter bases and interspace pruning

Inspired by *sparse dictionary learning* (SDL) (Sec. 4.1) we introduce interspace representations of convolutional filters and propose computations of resulting FB convolutions in Sec. 4.2. Further, Sec. 4.3 discusses FB sharing and the initialization of FBs and their coefficients. Finally, interspace pruning is formally defined in Sec. 4.4.

4.1. Inspiration from sparse dictionary learning

Sparse dictionary learning [2, 16, 52] optimizes a dictionary $\mathbf{F} \in \mathbb{R}^{m \times m}$ jointly with coefficients $R \in \mathbb{R}^{m \times n}$ to approximate a target $U \in \mathbb{R}^{m \times n}$ by using only s non-zero coefficients. Setting the *pruning mask* $\text{supp } R := \{(i, j) : R_{i,j} \neq 0\}$, this defines a non-convex optimization problem

$$\inf_{\mathbf{F}, R} \|U - \mathbf{F} \cdot R\|_F \quad \text{s.t.} \quad \|R\|_0 := \#\text{supp } R \leq s. \quad (1)$$

Usually, SDL allows $\mathbf{F} \in \mathbb{R}^{m \times M}$ with arbitrary M . Since FBs are bases, we restrict \mathbf{F} to be quadratic. In our context, U corresponds to all flattened filters of a convolutional layer, the dictionary \mathbf{F} to the layer’s flattened FB \mathcal{F} and R to the FB coefficients. For a layer $h \in \mathbb{R}^{c_{out} \times c_{in} \times K \times K}$ with an associated FB $\mathcal{F} = \{g^{(1)}, \dots, g^{(K^2)}\} \subset \mathbb{R}^{K \times K}$, we have $m = K^2$ and $n = c_{out} \cdot c_{in}$. Standard magnitude pruning is a special case of SDL where \mathbf{F} is fixed to form the standard basis $\mathbf{F} = \text{id}_{\mathbb{R}^m}$. Accordingly,

$$\min_{\bar{R}} \|U - \bar{R}\|_F \quad \text{s.t.} \quad \|\bar{R}\|_0 \leq s \quad (2)$$

is minimized for magnitude pruning. Since we train sparse, randomly initialized CNNs, our overall goal is not to *mimic* a given dense CNN, but to train the sparse network to *generalize* well. We consequently use Eqs. (1) and (2) only to find a decent subset of coefficients to be pruned. In contrast to SDL, deep learning methods are used to further optimize the un-pruned coefficients and additionally the FBs in the case of IP. In our experimental evaluation, we also test other methods than magnitude pruning, *i.e.* Eqs. (1) and (2). Still, Eqs. (1) and (2) measure the ability of a sparse layer to function as well as a dense layer and thus are good indicators for the general performance of IP and SP, respectively.

Most SDL algorithms [2, 16, 52] optimize \mathbf{F} and R alternately. Whereas, SP-PaI fixes the basis as $\text{id}_{\mathbb{R}^m}$ and the pruning mask $\text{supp } \bar{R}$ too. This simplifies the task, but reduces the solution space. IP overcomes the small, fixed solution space by adapting the basis during training. For IP-PaI, the pruning mask $\text{supp } R$ is determined heuristically and also fixed which still leads to sub-optimal architectures. As shown in this work, using expensive pre-training to find a better pruning mask via LTs or adapting $\text{supp } R$ during training via DST further improves IP’s performance.

Theorem 1 shows that a dynamic \mathbf{F} leads to better approximations than using the standard basis. Consequently,

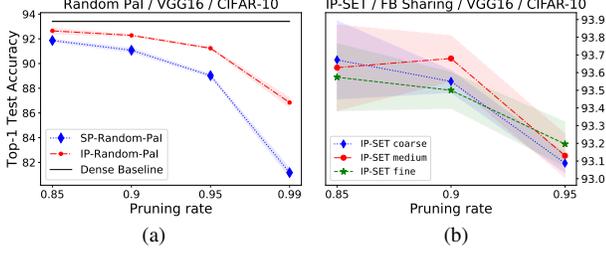


Figure 2. VGG16 on CIFAR-10: (a) Random SP-PaI and random IP-PaI. (b) Coarse, medium and fine FB sharing for IP-SET.

the FBs’ adaptivity improves performance after pruning and Thm. 1 is a theoretical motivation for IP.

Assume a convolutional layer with c_{out} output and c_{in} input channels, kernel size $K \times K$ and $s = (1 - p) \cdot c_{out} \cdot c_{in} \cdot K^2$ un-pruned coefficients. For $m = K^2 \geq 9$, the δ in Eq. (3) is numerically equal to zero if $n = c_{out} \cdot c_{in} \geq 100$ and $0 < s < c_{out} \cdot c_{in} \cdot K^2$. Thus, for each non-trivial sparsity, the adaptivity of the FB improves results. This even holds if the pruning mask for Eq. (1) is fixed to be the one of the minimizer of Eq. (2), *i.e.* starting with an arbitrary pruned network and adding an adaptive FB always improves results. The proof of Thm. 1 is shown in Appendix Sec. J. It uses the fact that Eq. (1) is smaller or equal to Eq. (2). Equality is only possible if Eq. (2) has a solution such that each $K \times K$ filter is either fully pruned or dense. This is almost impossible for big layers and a non trivial sparsity. If $\text{supp } R$ is further not fixed for Eq. (1), (F, R) can be chosen such that Eq. (1) is always strictly smaller than Eq. (2).

Theorem 1. Let $0 < s < m \cdot n$, $m > 1$ and $U_{i,j} \sim \mathcal{N}(0, 1)$ i.i.d. Let $\varepsilon_{(1)}$ be the infimum of Eq. (1), and $\varepsilon_{(2)}$ the minimum of Eq. (2) solved by \bar{R}^* . Then, $\varepsilon_{(1)} < \varepsilon_{(2)}$ with $\mathbb{P} = 1$. If $\text{supp } R$ is fixed for Eq. (1) to be $\text{supp } \bar{R}^*$, $\varepsilon_{(1)} \leq \varepsilon_{(2)}$ is true and strict inequality holds with $\mathbb{P} \geq 1 - \delta$, where

$$\delta = \begin{cases} \left(\frac{\binom{n}{s}}{\binom{m}{s}}\right) / \binom{m-n}{s} & \text{if } s \equiv 0 \pmod{m} \\ 0 & \text{else} \end{cases}. \quad (3)$$

Figure 2(a) compares SP and IP for random PaI for a VGG16 [72] trained on CIFAR-10 [38]. IP improves results tremendously compared to SP. This experimentally shows that sparse training performs better when coefficients of adaptive FBs are pruned than if spatial weights are pruned. This holds even though fixed pruning masks are used.

4.2. Interspace representation and convolutions

For a convolutional layer, let c_{out} denote its number of output channels, c_{in} its number of input channels and $K \times K$ its kernel size. To simplify formulas, we restrict the formulation to 2D convolutions with quadratic kernel, no padding, 1×1 stride and dilation. Generalizing the FB formulations to arbitrary convolutions is straightforward. A

2D convolution h describing this layer consists of $c_{out} \cdot c_{in} K \times K$ filters $h^{(\alpha,\beta)} \in \mathbb{R}^{K \times K}$, *i.e.* $h = (h^{(\alpha,\beta)})_{\alpha,\beta} \in \mathbb{R}^{c_{out} \times c_{in} \times K \times K}$. Inspired by the discussion in Sec. 4.1, we now represent all $h^{(\alpha,\beta)}$ in the interspace spanned by the layer’s FB $\mathcal{F} = \{g^{(1)}, \dots, g^{(K^2)}\} \subset \mathbb{R}^{K \times K}$. The FB coefficients $\lambda = (\lambda_n^{(\alpha,\beta)})_{\alpha,\beta,n} \in \mathbb{R}^{c_{out} \times c_{in} \times K^2}$ define the interspace representation of $h^{(\alpha,\beta)}$, given by

$$h^{(\alpha,\beta)} = \sum_{n=1}^{K^2} \lambda_n^{(\alpha,\beta)} \cdot g^{(n)}. \quad (4)$$

This is a basis transformation of the spatial representation

$$h^{(\alpha,\beta)} = \sum_{n=1}^{K^2} h_{i_n, j_n}^{(\alpha,\beta)} \cdot e^{(n)}, \quad e_{i,j}^{(n)} = \delta_{i,i_n} \cdot \delta_{j,j_n}. \quad (5)$$

Normally, $h^{(\alpha,\beta)}$ is defined in spatial representation. Thus, spatial coefficients are stored in $h^{(\alpha,\beta)} \in \mathbb{R}^{K \times K}$. Whereas, FB coefficients are specified by vectors $\lambda^{(\alpha,\beta)} \in \mathbb{R}^{K^2}$. By linearity, a 2D FB convolution $(Y^{(\alpha)})_{\alpha} = Y = h \star X$ with input feature map $X = (X^{(\beta)})_{\beta} \in \mathbb{R}^{c_{in} \times h \times w}$ can be computed for each output channel $\alpha \in \{1, \dots, c_{out}\}$ as

$$Y^{(\alpha)} = \sum_{\beta=1}^{c_{in}} h^{(\alpha,\beta)} \star X^{(\beta)} \stackrel{(4)}{=} \sum_{\beta=1}^{c_{in}} \sum_{n=1}^{K^2} \lambda_n^{(\alpha,\beta)} (g^{(n)} \star X^{(\beta)}). \quad (6)$$

Gradients of the loss \mathcal{L} are needed to train the FB coefficients λ and the FB \mathcal{F} . Backpropagation formulas for them are derived in Appendix Sec. D.2. It holds for all n, α, β

$$\frac{\partial \mathcal{L}}{\partial \lambda_n^{(\alpha,\beta)}} = \left\langle g^{(n)}, \frac{\partial \mathcal{L}}{\partial h^{(\alpha,\beta)}} \right\rangle, \quad \frac{\partial \mathcal{L}}{\partial g^{(n)}} = \sum_{\alpha,\beta} \lambda_n^{(\alpha,\beta)} \frac{\partial \mathcal{L}}{\partial h^{(\alpha,\beta)}}. \quad (7)$$

4.3. Filter basis sharing and initialization

For kernel size 1×1 , the FB formulation is, up to a rescaling, equivalent to the spatial representation. Thus, we assume a CNN with L_c convolutional layers with $K > 1$ to be given and do not apply the FB formulation to 1×1 convolutions. In this work, we test three versions of FB sharing. Our FB sharing schemes differ in their *granularity*. The *coarse* scheme shares one global FB \mathcal{F} for all layers $l = 1, \dots, L_c$. Whereas, the *fine* scheme shares a FB $\mathcal{F}^{(l)}$ for each layer l , thus it uses L_c FBs. In between lies the *medium* scheme with 5 FBs in total. For ResNets [32], one FB is shared for each of the 5 convolutional blocks. For VGG16 [72], convolutional layers $\{1, 2\}$, $\{3, 4\}$, $\{5, 6, 7\}$, $\{8, 9, 10\}$ and $\{11, 12, 13\}$ share one FB each. The number of FBs increases from *fine* to *coarse*. The total number of FBs in the network, J , satisfies $J \leq L_c$. Consequently, the number of parameters in *all* FBs in the network is bounded from above by $L_c \cdot K^4$. Note for the CNNs used

in this work, $L_c \cdot K^4$ is at most 0.01% of all parameters in the model. Thus, the additional parameter costs for IP with our proposed sharing schemes are neglectable.

The dimension of the space spanned by each layer does not change for different FB sharing schemes and is equal to using spatial representations. However, `coarse` sharing correlates all layers in the network by using and updating the same interspace. For `fine` sharing, each layer has its own interspace which is adapted more fine-grained. For spatial representations, the basis \mathcal{B} is fixed, not updated and does not induce correlations between weights. We found different sharing schemes to work best for varying training/model/dataset combinations. Figure 2(b) shows our FB sharing schemes for different pruning rates. `Coarse` sharing works best for higher numbers of trained parameters. By correlating all layers through a global FB, we assume it to have a regularizing effect on training, see also Sec. 5.4. `Fine` sharing makes the network more flexible. Thus, results are the best ones for high pruning rates where the network is not able to overfit on the training data anymore. In between, `medium` sharing reaches the best results by combining the best of both worlds.

In this work, we use a simple initialization for FBs and FB coefficients. We initialize each FB as \mathcal{B} and the FB coefficients with a `kaiming normal` initialization [31]. This scheme is equivalent to the `kaiming normal` initialization for standard CNNs – which is also used for dense baselines and SP experiments. In Appendix Sec. G, we propose further initialization schemes for the interspace.

4.4. Interspace pruning and cost comparison

SP is modeled by superimposing *pruning masks* $\bar{\mu}^{(\alpha,\beta)} \in \{0, 1\}^{K \times K}$ over filters $h^{(\alpha,\beta)} \in \mathbb{R}^{K \times K}$. This results in sparse filters $h^{(\alpha,\beta)} \odot \bar{\mu}^{(\alpha,\beta)}$, with the Hadamard product \odot . Filters represented in the interspace have coefficients $\lambda^{(\alpha,\beta)} \in \mathbb{R}^{K^2}$ w.r.t. a FB \mathcal{F} . Thus, interspace pruning is defined by masking FB coefficients with pruning masks $\mu^{(\alpha,\beta)} \in \{0, 1\}^{K^2}$ via $\lambda^{(\alpha,\beta)} \odot \mu^{(\alpha,\beta)}$. Combined with Eq. (6), IP yields sparse computations of convolutions:

$$Y^{(\alpha)} = \sum_{\beta=1}^{c_{in}} \sum_{n \in \text{supp } \mu^{(\alpha,\beta)}} \lambda_n^{(\alpha,\beta)} \cdot \left(g^{(n)} \star X^{(\beta)} \right). \quad (8)$$

The *pruning rate* p for SP (p_{SP}) and IP (p_{IP}) is defined as

$$p_{SP} = 1 - \frac{\|\Lambda\|_0}{D}, \quad p_{IP} = 1 - \frac{\|\Lambda\|_0 + \sum_{j=1}^J \|\mathcal{F}^{(j)}\|_0}{D}. \quad (9)$$

For SP, $\Lambda \in \mathbb{R}^D$ denotes the network’s parameters, whereas $\Lambda \in \mathbb{R}^D$ contains all parameters except the FBs themselves in the IP setting. Thus, Λ has exactly the same number of elements for IP and SP. The pruning rates Eq. (9) are the fractions of parameters being equal to zero. To have a fair

Algorithm 1 FB 2D Convolution with IP

```

1: instance variables ▷ of IP_FB_2DConv
2:   filter_basis:  $\{g^{(1)}, \dots, g^{(K^2)}\} \subset \mathbb{R}^{K \times K}$ 
3:   fb_coefficients:  $(\lambda_n^{(\alpha,\beta)})_{\alpha,\beta,n} \in \mathbb{R}^{c_{out} \times c_{in} \times K^2}$ 
4:   pruning_mask:  $(\mu_n^{(\alpha,\beta)})_{\alpha,\beta,n} \in \{0, 1\}^{c_{out} \times c_{in} \times K^2}$ 
5:   conv_args ▷ e.g. stride, padding, groups, ...
6: def FORWARD_PASS( $X$ ) ▷ input  $X \in \mathbb{R}^{c_{in} \times h \times w}$ 
7:   for all  $\beta \in \{1, \dots, c_{in}\}, n \in \{1, \dots, K^2\}$  do
8:      $Z_n^{(\beta)} = \text{Conv2D}(g^{(n)}, X^{(\beta)}, \text{conv\_args})$ 
9:   for all  $\alpha \in \{1, \dots, c_{out}\}$  do
10:     $Y^{(\alpha)} = \sum_{\{(\beta,n): \mu_n^{(\alpha,\beta)}=1\}} \lambda_n^{(\alpha,\beta)} \cdot Z_n^{(\beta)}$ 
11:   return  $Y = (Y^{(\alpha)})_{\alpha=1}^{c_{out}}$ 

```

comparison between IP and SP, we normalize the number of non-zero parameters with the total count of coefficients in the standard dense network, *i.e.* the dense network without FBs. The number of bias and batch normalization parameters is tiny compared to convolutional and fully connected layers. Also, all parameters of FBs together are at most 0.01% of D in our experiments. Consequently, we only prune weights of fully connected layers as well as spatial- and FB coefficients of convolutional layers. FBs, bias and batch normalization parameters are all trained.

Computational cost comparison. As discussed, parameter costs for IP with our FB sharing schemes are only negligibly bigger than for SP. By the linearity of convolutions, the sparsity of filters in the interspace can be used to reduce computational costs, see Eq. (8). In Appendix Sec. D, computational costs are calculated and compared for IP and SP. Costs are measured by the number of theoretically required *floating point operations* (FLOPs) for a convolutional layer and are independent of the used FB sharing scheme. IP’s overhead is composed of additional costs in the forward and backward pass. For inference, only the additional cost of the forward pass counts. Both, SP and IP, need specialized soft- and hardware that supports sparse computations to actually reduce runtime.

Assume a layer with kernel size $K \times K$, c_{in} input and c_{out} output channels. In the forward pass, SP has $1 - p$ times the FLOPs cost of the dense layer. Due to l. 7-8 in Alg. 1, IP has a constant overhead K^2/c_{out} . In total, IP has $1 - p + K^2/c_{out}$ times the FLOPs cost of the dense layer.

In the backward pass, the number of FLOPs for IP is in $\mathcal{O}(\text{cost}(\frac{\partial \mathcal{L}}{\partial h}))$, *i.e.* comparable to the cost of computing the dense gradient of layer h in spatial representation.

As discussed, IP needs more computations for inference than SP for equal sparsity. However, since IP finds superior sparse models, IP actually achieves a *higher speed up* in real time measurements than SP while reaching similar or even better performance, as will be shown Fig. 5(a).

Pruning methods. Algorithm 1 describes sparse FB 2D convolutions with IP in pseudo code. Since automatic differentiation is standard in modern deep learning frameworks, backpropagation formulas for FB convolutions are computed automatically and are not included in Alg. 1. The FB in Alg. 1 might be shared over several layers, see Sec. 4.3. Our experiments in Sec. 5 compare SP and IP on various sparse training and other pruning methods, namely:

DST randomly prunes the model at initialization. During training, unimportant coefficients are pruned based on their magnitude. In each layer, the same number of parameters is *regrown* by activating their gradients. SET regrows coefficients randomly whereas RigL regrows those with high gradient magnitude. The pruning mask is updated each 1,500 iterations for SET and 4,000 for RigL. A cosine schedule is used to reduce the number of pruned/regrown coefficients.

LT pre-trains the network for $t_0 = 500$ steps. Then, the network is trained to convergence. Now, 20% of the non-zero coefficients are pruned based on their magnitude. The un-pruned part of the CNN is reset to its value at t_0 . The whole procedure is applied k times in total until the desired pruning rate $p = 1 - 0.8^k$ is reached. Ultimately, the final sparse network is trained, starting at t_0 .

PaI prunes the model at initialization without pre-training or changing the pruning mask during training. Random PaI prunes weights i.i.d. with probability p . SNIP trains coefficients which have high influence on changing the loss \mathcal{L} when training starts. GraSP finds coefficients which improve the gradient flow at the beginning of training most. SynFlow keeps coefficients with high information throughput which is measured by their influence on the total *path norm* of the sparse network.

Gradual Magnitude Pruning (GMP) [21] starts training with dense coefficients. During training, the CNN is gradually sparsified based on the coefficients’ magnitudes. Pruned parameters are fixed at zero, thus never regrow.

Fine-Tuning (FT) [68] uses a pre-trained network. The $p \cdot \bar{D}$ coefficients with smallest magnitude are pruned. The pre-trained coefficients of the sparse CNN are fine-tuned with the learning rate schedule of the dense training.

All these methods were developed for SP. Yet, in our experiments they are applied unchanged to the interspace setting. For more details see Appendix Secs. F and G.

5. Experiments and discussion

Section 5.1 covers the experimental setup. Next, Sec. 5.2 compares the three SOTA PaI methods [41, 76, 82] for IP and SP. In Sec. 5.3, we discuss IP and SP for more sophisticated sparse training methods, namely LTs [19] and the DST methods SET [55] and RigL [17]. Furthermore, we show that IP also improves SP on classical pruning methods applied during training, GMP [21], and on pre-trained models, FT [68]. Improved trainability and generalization ability of

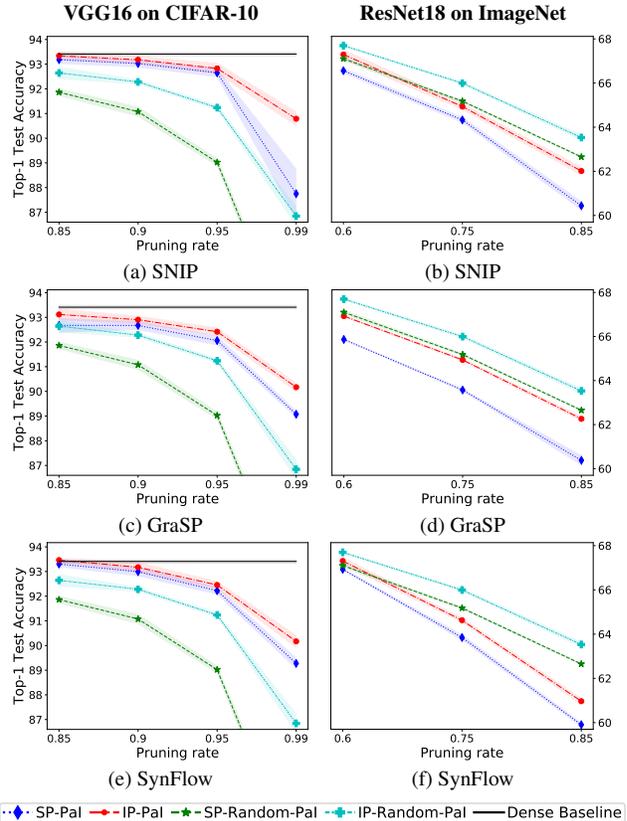


Figure 3. Comparing SP and IP for PaI methods, SNIP, GraSP and SynFlow together with random PaI for CIFAR-10 and ImageNet.

IP compared to SP is shown and discussed in Sec. 5.4.

5.1. Experimental setup

We compare IP and SP for a VGG16 [72] on CIFAR-10 [38] and ResNets 18 and 50 [32] on ImageNet ILSVRC2012 [69]. Models are trained with cross entropy loss. We report mean and std of five runs for CIFAR-10 and three for ImageNet. Weight decay is applied on coefficients but not on FBs. Coefficients of 3×3 filters and their FBs are trained jointly, whereas fixed FBs $\mathcal{F} = \mathcal{B}$ are used for 1×1 filters. For ResNet18 we fix the FB $\mathcal{F} = \mathcal{B}$ for the 7×7 convolution whereas the 7×7 FB is trained for ResNet50. We use medium FB sharing for CIFAR-10 experiments, fine for ResNet50 and coarse sharing for all 3×3 convolutions for the ResNet18 on ImageNet. For SP and dense baselines, standard CNNs are used. As common in the literature, we report ImageNet results on the validation set. Note, we use training schedules intended for the corresponding SP method for both, SP and IP. In particular, FBs are trained without optimized hyperparameters. Thus, they use the same learning rate as all parameters. More details on hyperparameters, evaluation and used CNN architectures are given in Secs. H and I in the Appendix.

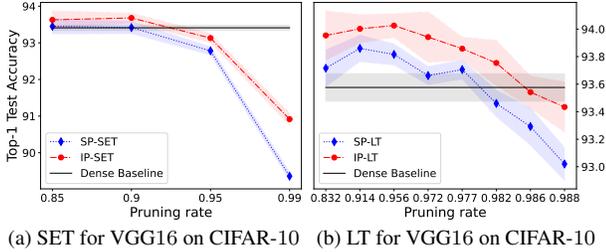


Figure 4. Comparison between SP and IP for (a) the DST method SET and (b) LT on a VGG16 trained on CIFAR-10.

5.2. Pruning at initialization methods

Figure 3 compares SP and IP for PaI methods SNIP [41], GraSP [82] and SynFlow [76] together with random PaI for a VGG16 on CIFAR-10 and a ResNet18 on ImageNet.

The experiments show that pruning FB coefficients instead of spatial parameters leads to significant improvements in top-1 test accuracy while having the same memory costs. This holds true for all PaI methods, pruning rates and for high p in particular. In comparison to CIFAR-10, IP improves results on ImageNet even more. However, the three methods SNIP, GraSP and SynFlow are all outperformed by random PaI for ResNet18 on ImageNet. This demonstrates that these methods perform well for smaller datasets but show inferior results for small networks on big scale datasets like ImageNet. Still, as discussed earlier, the use of IP significantly improves all PaI methods, including random PaI. Section 5.3 shows that IP benefits from a stronger underlying pruning method to improve results further.

Despite optimizing FBs in addition to FB coefficients, IP does not induce instability compared to SP, see Fig. 5(b) and standard deviations in Fig. 3. In Appendix Sec. D.3, we show that the upper bounds for the gradient norms of FBs $\frac{\partial \mathcal{L}}{\partial \mathcal{F}}$ and FB coefficients $\frac{\partial \mathcal{L}}{\partial \lambda}$ are both determined by $\|\frac{\partial \mathcal{L}}{\partial h}\|$. This boundedness of the gradients leads to stable convergence for both, \mathcal{F} and λ , while the convergence behavior of λ and the standard coefficients h is similar, see Fig. 5(b).

5.3. DST, LTs and classical pruning methods

For SP, more expensive or sophisticated methods like LT and DST improve sparse training results compared to PaI. We want to analyze whether this also applies to the IP setting. Furthermore, we want to check if IP boosts the SOTA methods LT and RigL as well. Finally, we benchmark IP and SP on various SOTA unstructured pruning methods for a ResNet50 on ImageNet.

DST and LT on CIFAR-10. IP improves DST and LTs significantly, see Figs. 4(a) and (b). For all p , IP-LT surpasses SP-LT. IP needs to train 3.7 times less parameters ($p = 0.977$) than SP to reach SP’s best result for $p = 0.914$. IP-LT matches the dense baseline while training only 1.4%

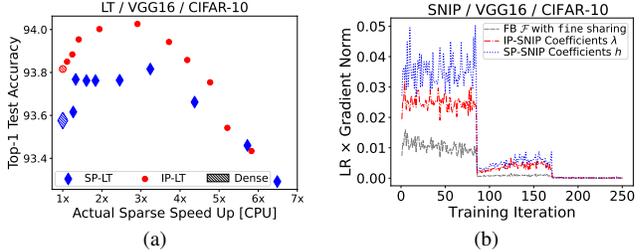


Figure 5. VGG16 on CIFAR-10: (a) Top-1 test accuracy over real time acceleration for IP- and SP-LT. (b) Gradient L_2 norm \times learning rate (LR) for SP- and IP-SNIP for layer 1 and $p = 0.85$.

Top-1 Accuracy for ResNet50 on ImageNet			
Method	$p = 0.0$	$p = 0.8$	$p = 0.9$
SP-FT	77.15 \pm 0.04	77.02 \pm 0.03	75.67 \pm 0.09
IP-FT	77.30 \pm 0.04	77.18 \pm 0.01	75.89 \pm 0.09
SP-GMP	76.64 \pm 0.06	75.37 \pm 0.01	73.57 \pm 0.06
IP-GMP	<u>77.16 \pm 0.04</u>	<u>75.71 \pm 0.05</u>	<u>74.20 \pm 0.07</u>
SP-RigL	77.15 \pm 0.04	75.75 \pm 0.10	73.88 \pm 0.06
IP-RigL	77.30 \pm 0.04	<u>76.03 \pm 0.08</u>	<u>74.32 \pm 0.11</u>

Table 1. ResNet50 trained on ImageNet for 100 epochs.

of its parameters and outperforms it for all $p \leq 0.98$. Comparable results hold for SET. IP-SET improves the dense baseline for $p \leq 0.9$, whereas SP-SET only matches it. Similar to PaI, IP-SET greatly exceeds SP-SET for high p . Comparing Figs. 3 and 4 shows that spending more effort in finding the sparse architecture (LT) or adapting it during training (SET) improves performance compared to PaI for both, SP and IP.

ResNet50 on ImageNet. Table 1 compares IP and SP on the SOTA pruning methods RigL [17], GMP [21] and FT [68]. As shown, IP outperforms all underlying SP methods for a ResNet50 on ImageNet. Results are significantly improved with interspace representations even though more than 50% of the coefficients of a ResNet50 are 1×1 convolutions which are equivalent for IP and SP. For example, IP-FT has similar performance as a standard dense model while training only 20% of its parameters. Note, using FBs does not only boost training sparse CNNs but dense training too, which will be discussed in more detail in Sec. 5.4.

Computational costs. Up to now, IP and SP were compared for equal memory costs. As analyzed in Sec. 4.4, IP has a small computational overhead compared to SP for equal sparsity. In applications, the actual runtime is more important than the theoretically required FLOPs. Thus, we compare the performance of IP and SP w.r.t. the actual acceleration on a CPU achieved by using sparse representations. Details on the implementation are provided in the Appendix Sec. D.4. IP indeed has a longer runtime for equal spar-

VGG16 on CIFAR-10				
Method	$p = 0.85$		$p = 0.99$	
	Train	Test	Train	Test
SP-SET	99.85	93.45	94.20	89.36
IP-SET	99.89	93.63	96.89	90.92
SP-SNIP	99.94	93.18	93.96	87.75
IP-SNIP	99.96	93.34	98.38	90.79

ResNet50 on ImageNet				
Method	$p = 0.8$		$p = 0.9$	
	Train	Test	Train	Test
SP-RigL	74.64	75.75	71.30	73.88
IP-RigL	75.39	76.03	72.08	74.32

Table 2. Generalization gaps for various pruning methods.

sity due to the mentioned extra computations. However, by boosting performance of sparse models, IP reaches similar results than dense training with 5.2 times speed up and better results than SP for equal runtime, see Fig. 5(a).

5.4. Generalization and trainability

We consider generalization as the ability to correctly classify unseen data [50]. In this context a major aspect is the relationship between performance on the train and test set. Ideally, the performance on the train set should be optimal and a strong indicator for the performance on the test set. The *generalization gap* is the difference between train and test accuracy. Generalization can be improved by regularizations [8, 34, 39, 73, 90], enabling the model to use geometrical prior knowledge about the scene [10, 11, 35, 65, 66], shifting the model back to an area where it generalizes well [49, 51, 67, 71] but also by pruning the network [4, 30, 40].

Table 2 shows training and test accuracy for the IP- and SP versions of SET and SNIP for a VGG16 on CIFAR-10 as well as RigL for a ResNet50 on ImageNet. IP pruned networks train better than SP pruned ones for all p . Note, the used ImageNet training is highly regularized. Thus, the test accuracy is *higher* than the train accuracy. For ImageNet and $p = 0.99$ on CIFAR-10, IP has a bigger generalization gap than SP. This is due to a much better training accuracy for IP, which in the end leads to an improved test accuracy. However, IP has a smaller generalization gap than SP for $p = 0.85$ on CIFAR-10 where the model overfits.

Table 3 further shows that IP can generally improve results for pruning rates where training overfits. Note, $p = 0$ is dense training and SP for $p = 0$ is standard dense training. Improved performance in the dense setting can not be explained by IP’s superior expressiveness (Thm. 1) since IP and SP can represent the same if all parameters are un-pruned. We hypothesize that correlating filters in a CNN via FB sharing regularizes training, thereby improving generalization. One indicator of this is the fact that correlating *all* filters via *coarse* sharing shows the best results while *fine* sharing

Method	Pruning rate p			
	0.0	0.35	0.6	0.85
SNIP				
SP	93.4 ± 0.1	93.4 ± 0.1	93.3 ± 0.2	93.2 ± 0.2
IP-coarse	93.9 ± 0.2	93.7 ± 0.2	93.8 ± 0.1	93.5 ± 0.0
IP-medium	93.9 ± 0.2	93.6 ± 0.2	93.7 ± 0.2	93.3 ± 0.2
IP-fine	93.7 ± 0.1	93.3 ± 0.2	93.3 ± 0.2	93.2 ± 0.1
SET				
SP	93.4 ± 0.1	93.5 ± 0.2	93.3 ± 0.2	93.5 ± 0.2
IP-coarse	93.9 ± 0.1	93.9 ± 0.2	93.8 ± 0.1	93.7 ± 0.2
IP-medium	93.9 ± 0.1	93.7 ± 0.2	93.8 ± 0.1	93.6 ± 0.2
IP-fine	93.7 ± 0.1	93.6 ± 0.2	93.6 ± 0.2	93.6 ± 0.2

Table 3. Varying FB sharing schemes for lower pruning rates p .

has comparable results to SP. Consequently, interspace representations can also be used to regularize dense training even for ResNet50 on ImageNet, see Tab. 1. After training, dense interspace representations can be converted to standard ones to reduce computational costs for inference. By optimizing weight decay and initialization schemes, IP’s performance can be increased even further, as shown in Appendix Sec. B.

6. Conclusions and directions for future work

IP significantly improves results compared to pruning spatial coefficients. We demonstrate this by achieving SOTA results with the application of IP to SOTA standard PaI, LT, DST as well as classical pruning methods.

Theorem 1 proves that IP leads to better sparse approximations than SP. Especially, IP generates models with *higher sparsity and equal performance* than SP. Also, FB representations combined with FB sharing *improve generalization* of overfitting CNNs, even for dense training. This comes with the prize of a small computational overhead for inference and additional gradient computations during training. Nevertheless, we show that sparse interspace representations accelerate dense baselines more than SP while keeping or even improving the baseline’s performance.

We believe that IP can be enhanced by adapting more advanced strategies of SDL to the joint training of \mathcal{F} and λ . Adapting IP to structured pruning is an option to maintain the network’s accuracy while reducing inference time for arbitrary soft- and hardware. Combining IP with low rank tensor approximations lowers computational costs as well and is discussed in Appendix Secs. B and D. The interspace representation is an adaptive basis transformation of a finite dimensional vector space. Therefore, FBs \mathcal{F} are not limited to represent convolutional filters but can express arbitrary vectors, like columns or small blocks of a matrix. This makes the concept of IP available for MLPs or self-attention modules.

Acknowledgements

The authors would like to thank their colleagues Julia Lust, Matthias Rath and Rinor Cakaj for their valuable contributions and fruitful discussions.

References

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Gregory S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian J. Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Józefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Gordon Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul A. Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda B. Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *CoRR*, abs/1603.04467, 2016. [iv](#)
- [2] M. Aharon, M. Elad, and A. Bruckstein. K-svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on Signal Processing*, 54(11):4311–4322, 2006. [3](#)
- [3] Sajid Anwar, Kyuyeon Hwang, and Wonyong Sung. Structured pruning of deep convolutional neural networks. *ACM Journal on Emerging Technologies in Computing Systems*, 13(3):1–18, 2017. [2](#)
- [4] Brian Bartoldson, Ari Morcos, Adrian Barbu, and Gordon Erlebacher. The generalization-stability tradeoff in neural network pruning. In *Advances in Neural Information Processing Systems 33*, 2020. [1, 8](#)
- [5] Guillaume Bellec, David Kappel, Wolfgang Maass, and Robert Legenstein. Deep rewiring: Training very sparse deep networks. In *International Conference on Learning Representations*, 2018. [xii](#)
- [6] Davis W. Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John V. Guttag. What is the state of neural network pruning? In *Proceedings of Machine Learning and Systems 2*, 2020. [1](#)
- [7] Miguel A. Carreira-Perpinan and Yerlan Idelbayev. "Learning-compression" algorithms for neural net pruning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018. [3](#)
- [8] Rich Caruana, Steve Lawrence, and Lee Giles. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *Advances in Neural Information Processing Systems 13*, 2000. [8](#)
- [9] Albert Cohen, Wolfgang Dahmen, and Ronald Devore. Compressed sensing and best k -term approximation. *Journal of the American Mathematical Society*, 22(1):211–231, 2009. [iii](#)
- [10] Taco S. Cohen and Max Welling. Group equivariant convolutional networks. In *Proceedings of the 33rd International Conference on Machine Learning*, 2016. [8](#)
- [11] Benjamin Coors, Alexandru Paul Condurache, and Andreas Geiger. Spherenet: Learning spherical representations for detection and classification in omnidirectional images. In *Proceedings of the European Conference on Computer Vision*, 2018. [8](#)
- [12] Pau de Jorge, Amartya Sanyal, Harkirat Behl, Philip Torr, Grégory Rogez, and Puneet K. Dokania. Progressive skeletonization: Trimming more fat from a network at initialization. In *International Conference on Learning Representations*, 2021. [1, 3](#)
- [13] Tim Dettmers and Luke Zettlemoyer. Sparse networks from scratch: Faster training without losing performance. *CoRR*, abs/1907.04840, 2019. [3, xii](#)
- [14] D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, 2006. [iii](#)
- [15] Erich Elsen, Marat Dukhan, Trevor Gale, and Karen Simonyan. Fast sparse convnets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. [1](#)
- [16] K. Engan, S. O. Aase, and J. H. Husøy. Method of optimal directions for frame design. *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing 5*, 1999. [1, 3](#)
- [17] Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *Proceedings of the 37th International Conference on Machine Learning*, 2020. [1, 3, 6, 7, ix, xii, xiii](#)
- [18] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2018. [1, 3](#)
- [19] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis. In *Proceedings of the 37th International Conference on Machine Learning*, 2020. [1, 3, 6, ix, xi, xii, xiii](#)
- [20] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Pruning neural networks at initialization: Why are we missing the mark? In *International Conference on Learning Representations*, 2021. [3](#)
- [21] Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. *CoRR*, abs/1902.09574, 2019. [3, 6, 7, ix, xii, xiii](#)
- [22] Trevor Gale, Matei Zaharia, Cliff Young, and Erich Elsen. Sparse gpu kernels for deep learning. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, 2020. [1](#)
- [23] Shangqian Gao, Feihu Huang, Weidong Cai, and Heng Huang. Network pruning via performance maximization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. [1](#)
- [24] Stuart Geman, Elie Bienenstock, and Rene Doursat. Neural networks and the bias/variance dilemma. *Neural Computation*, 4(1):1–58, 1992. [i](#)
- [25] Priya Goyal, Piotr Dollár, Ross B. Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: training imagenet in 1 hour. *CoRR*, abs/1706.02677, 2017. [xiii](#)

- [26] Yiwen Guo, Anbang Yao, and Yurong Chen. Dynamic network surgery for efficient dnns. In *Advances in Neural Information Processing Systems* 29. 2016. [1](#), [3](#)
- [27] Song Han, Xingyu Liu, Huizi Mao, Jing Pu, Ardavan Pedram, Mark A. Horowitz, and William J. Dally. Eie: Efficient inference engine on compressed deep neural network. *ACM SIGARCH Computer Architecture News*, 44(3):243–254, 2016. [1](#), [3](#)
- [28] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In *Advances in Neural Information Processing Systems* 28. 2015. [1](#), [3](#)
- [29] Stephen Jose Hanson and Lorien Y. Pratt. Comparing biases for minimal network construction with back-propagation. In *Advances in Neural Information Processing Systems 1*. 1989. [3](#)
- [30] Babak Hassibi and David Stork. Second order derivatives for network pruning: Optimal brain surgeon. In *Advances in Neural Information Processing Systems*, 1992. [1](#), [8](#)
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *IEEE International Conference on Computer Vision*, 2015. [5](#), [x](#)
- [32] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. [4](#), [6](#), [xiv](#)
- [33] Zehao Huang and Naiyan Wang. Data-driven sparse structure selection for deep neural networks. *Proceedings of the European conference on computer vision*, 2018. [2](#)
- [34] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning*, 2015. [8](#), [xii](#)
- [35] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. In *Advances in Neural Information Processing Systems*, 2015. [8](#)
- [36] Steven A. Janowsky. Pruning versus clipping in neural networks. *Physical Review A*, 39:6600–6603, 1989. [1](#), [3](#)
- [37] Ehud D. Karnin. A simple procedure for pruning back-propagation trained neural networks. *IEEE Transactions on Neural Networks*, 1(2):239–242, 1990. [3](#)
- [38] Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 2012. <http://www.cs.toronto.edu/~kriz/cifar.html>. [4](#), [6](#), [xiii](#)
- [39] Anders Krogh and John A. Hertz. A simple weight decay can improve generalization. In *Advances in Neural Information Processing Systems 4*. 1992. [8](#), [i](#)
- [40] Yann LeCun, John S. Denker, and Sara A. Solla. Optimal brain damage. In *Advances in Neural Information Processing Systems 2*. 1990. [1](#), [3](#), [8](#)
- [41] Namhoon Lee, Thalaiyasingam Ajanthan, and Philip H.S. Torr. SNIP: Single-shot network pruning based on connection sensitivity. In *International Conference on Learning Representations*, 2019. [1](#), [3](#), [6](#), [7](#), [iii](#), [vii](#), [ix](#), [xii](#), [xiii](#)
- [42] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *International Conference on Learning Representations*, 2017. [3](#)
- [43] Yawei Li, Shuhang Gu, Luc Van Gool, and Radu Timofte. Learning filter basis for convolutional neural network compression. In *IEEE International Conference on Computer Vision*, 2019. [3](#)
- [44] Zhengang Li, Geng Yuan, Wei Niu, Pu Zhao, Yanyu Li, Yuxuan Cai, Xuan Shen, Zheng Zhan, Zhenglun Kong, Qing Jin, Zhiyu Chen, Sijia Liu, Kaiyuan Yang, Bin Ren, Yanzhi Wang, and Xue Lin. Npas: A compiler-aware framework of unified network pruning and architecture search for beyond real-time mobile acceleration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. [2](#)
- [45] Shiwei Liu, Lu Yin, Decebal Constantin Mocanu, and Mykola Pechenizkiy. Do we actually need dense over-parameterization? In-time over-parameterization in sparse training. In *Proceedings of the 38th International Conference on Machine Learning*, 2021. [1](#), [3](#), [xii](#), [xiii](#)
- [46] Xingyu Liu, Jeff Pool, Song Han, and William J. Dally. Efficient sparse-winograd convolutional neural networks. In *International Conference on Learning Representations*, 2018. [3](#)
- [47] Zhenhua Liu, Jizheng Xu, Xiulian Peng, and Ruiqin Xiong. Frequency-domain dynamic pruning for convolutional neural networks. In *Advances in Neural Information Processing Systems 31*, 2018. [3](#)
- [48] Christos Louizos, Max Welling, and Diederik P. Kingma. Learning sparse neural networks through l0 regularization. In *International Conference on Learning Representations*, 2018. [3](#)
- [49] Julia Lust and Alexandru Paul Condurache. Gran: An efficient gradient-norm based detector for adversarial and misclassified examples. In *28th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, 2020. [8](#)
- [50] Julia Lust and Alexandru Paul Condurache. A survey on assessing the generalization envelope of deep neural networks at inference time for image classification. *CoRR*, abs/2008.09381, 2020. [8](#)
- [51] Julia Lust and Alexandru Paul Condurache. Efficient detection of adversarial, out-of-distribution and other misclassified samples. *Neurocomputing*, 470:335–343, 2022. [8](#)
- [52] Julien Mairal, Francis Bach, Jean Ponce, and Guillermo Sapiro. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11:19–60, 2010. [3](#)
- [53] Eran Malach, Gilad Yehudai, Shai Shalev-Schwartz, and Ohad Shamir. Proving the lottery ticket hypothesis: Pruning is all you need. In *Proceedings of the 37th International Conference on Machine Learning*, 2020. [1](#)
- [54] Huizi Mao, Song Han, Jeff Pool, Wenshuo Li, Xingyu Liu, Yu Wang, and William J. Dally. Exploring the granularity of sparsity in convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017. [3](#)
- [55] Decebal Mocanu, Elena Mocanu, Peter Stone, Phuong Nguyen, Madeleine Gibescu, and Antonio Liotta. Scalable

- training of artificial neural networks with adaptive sparse connectivity inspired by network science. *Nature Communications*, 9, 2018. 1, 3, 6, ix, xii, xiii
- [56] Hesham Mostafa and Xin Wang. Parameter efficient training of deep convolutional neural networks by dynamic sparse reparameterization. In *Proceedings of the 36th International Conference on Machine Learning*, 2019. xii
- [57] Michael C. Mozer and Paul Smolensky. Skeletonization: A technique for trimming the fat from a network via relevance assessment. In *Advances in Neural Information Processing Systems 1*. 1989. 1, 3
- [58] NVIDIA, Péter Vingelmann, and Frank H.P. Fitzek. Cuda, release: 10.2.89, 2020. xiii
- [59] Angshuman Parashar, Minsoo Rhu, Anurag Mukkara, Antonio Puglielli, Rangharajan Venkatesan, Bruce Khailany, Joel Emer, Stephen W. Keckler, and William J. Dally. Scnn. *Proceedings of the 44th Annual International Symposium on Computer Architecture*, 2017. 1, 3
- [60] Daniel S. Park, Yu Zhang, Chung-Cheng Chiu, Youzheng Chen, Bo Li, William Chan, Quoc V. Le, and Yonghui Wu. Spcaugment on large scale datasets. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2020. 1
- [61] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*. 2019. iv, xiii
- [62] Shreyas Malakarjun Patil and Constantine Dovrolis. PHEW: Constructing sparse networks that learn fast and generalize well without training data. In *Proceedings of the 38th International Conference on Machine Learning*, 2021. 3
- [63] Hieu Pham, Zihang Dai, Qizhe Xie, Minh-Thang Luong, and Quoc V. Le. Meta pseudo labels. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 1
- [64] Vivek Ramanujan, Mitchell Wortsman, Aniruddha Kembhavi, Ali Farhadi, and Mohammad Rastegari. What’s hidden in a randomly weighted neural network? In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 1
- [65] Matthias Rath and Alexandru Paul Condurache. Invariant integration in deep convolutional feature space. In *28th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, 2020. 8
- [66] Matthias Rath and Alexandru Paul Condurache. Improving the sample-complexity of deep classification networks with invariant integration. In *Proceedings of the 17th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, 2022. 8
- [67] Jie Ren, Peter J. Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. In *Advances in Neural Information Processing Systems*, 2019. 8
- [68] Alex Renda, Jonathan Frankle, and Michael Carbin. Comparing rewinding and fine-tuning in neural network pruning. In *International Conference on Learning Representations*, 2020. 3, 6, 7, ix, xiii
- [69] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 6, xiii
- [70] Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. Green AI. *Communications of the ACM*, 63(12):54–63, 2020. 1
- [71] Joan Serra, David Álvarez, Vicenç Gómez, Olga Slizovskaia, José F. Núñez, and Jordi Luque. Input complexity and out-of-distribution detection with likelihood-based generative models. In *International Conference on Learning Representations*, 2020. 8
- [72] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 4, 6, iii, xiv
- [73] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014. 8
- [74] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for modern deep learning research. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020. 1
- [75] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. xiii
- [76] Hidenori Tanaka, Daniel Kunin, Daniel L Yamins, and Surya Ganguli. Pruning neural networks without any data by iteratively conserving synaptic flow. In *Advances in Neural Information Processing Systems 33*, 2020. 1, 2, 3, 6, 7, vii, ix, xii, xiii
- [77] Yehui Tang, Yunhe Wang, Yixing Xu, Yiping Deng, Chao Xu, Dacheng Tao, and Chang Xu. Manifold regularized dynamic network pruning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 2
- [78] W.F. Tinney and J.W. Walker. Direct solutions of sparse network equations by optimally ordered triangular factorization. *Proceedings of the IEEE*, 55(11):1801–1809, 1967. 3, iii, vii
- [79] Karen Ullrich, Edward Meeds, and Max Welling. Soft weight-sharing for neural network compression. In *International Conference on Learning Representations*, 2017. 1
- [80] M. Unser and T. Blu. Mathematical properties of the jpeg2000 wavelet filters. *IEEE Transactions on Image Processing*, pages 1080–1090, 2003. iii
- [81] Stijn Verdenius, Maarten Stol, and Patrick Forré. Pruning via iterative ranking of sensitivity statistics. *CoRR*, abs/2006.00896, 2020. 3
- [82] Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by preserving gradient flow. In

- International Conference on Learning Representations*, 2020. [1](#), [2](#), [3](#), [6](#), [7](#), [vii](#), [ix](#), [x](#), [xii](#), [xiii](#)
- [83] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Scaled-yolov4: Scaling cross stage partial network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. [1](#)
- [84] Zi Wang, Chengcheng Li, and Xiangyang Wang. Convolutional neural network pruning with structural redundancy reduction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. [2](#)
- [85] Paul Wimmer, Jens Mehnert, and Alexandru Condurache. FreezeNet: Full performance by reduced storage costs. In *Proceedings of the Asian Conference on Computer Vision*, 2020. [2](#), [iii](#)
- [86] Paul Wimmer, Jens Mehnert, and Alexandru Condurache. COPS: Controlled pruning before training starts. In *International Joint Conference on Neural Networks*, 2021. [1](#), [3](#)
- [87] Ian H. Witten, Radford M. Neal, and John G. Cleary. Arithmetic coding for data compression. *Commun. ACM*, 30(6):520–540, 1987. [iv](#)
- [88] Huanrui Yang, Wei Wen, and Hai Li. DeepHoyer: Learning sparser neural network with differentiable scale-invariant sparsity measures. In *International Conference on Learning Representations*, 2020. [3](#)
- [89] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic segmentation. In *Proceedings of the European conference on computer vision*, 2020. [1](#)
- [90] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *5th International Conference on Learning Representations*, 2017. [8](#)
- [91] Tao Zhuang, Zhixuan Zhang, Yuheng Huang, Xiaoyi Zeng, Kai Shuang, and Xiang Li. Neuron-level structured pruning using polarization regularizer. In *Advances in Neural Information Processing Systems 33*, 2020. [2](#)
- [92] Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin Dogus Cubuk, and Quoc Le. Rethinking pre-training and self-training. In *Advances in Neural Information Processing Systems 33*, 2020. [1](#)

Appendix

A. Structure of the Appendix

The Appendix is divided into the following Sections:

- A** Describes the structural organization of the Appendix.
- B** Contains ablation studies for IP methods which are not shown in the main body of the text.
- C** Discussion about storing unstructured sparse networks.
- D** Gives detailed information about the computations of FB-CNNs including backpropagation formulas. Especially, a comparison between the number of FLOPs required to evaluate and train a convolutional layer in the spatial and interspace representation is drawn. Finally, details on the real time measurements of sparse speed ups are given.
- E** Computes transformation rules between the spatial and interspace coefficients, their gradients and Hessian matrices.
- F** Here, the computations of the pruning scores used for experiments in the main body of the text are proposed.
- G** Shows three different Algorithms to initialize FB-CNNs, including the standard initialization scheme used in the main body of the text. Further, details of implementations of the pruning methods are proposed.
- H** Describes used training setups, hyperparameters, datasets and evaluation procedure for experiments in the main body of the text.
- I** Presents network architectures, used in the experimental evaluation.
- J** Concludes the Appendix with a mathematical proof of Thm. 1.

B. Additional ablations

B.1. Using different initializations for the interspace

For SP-SNIP, the problem of vanishing gradients occurs, see Fig. A6(a). Filters which are spatially too sparse induce a vanishing gradient for high pruning rates. As shown in Fig. 1, IP leads to less zeros in the spatial representation of filters than SP *after* training. But, a pruned CNN has a spatially sparse topology *before* training if a standard initialization is used. This seems not to be the optimal initial situation for training FBs jointly with their coefficients. To analyze different starting conditions for IP, we initialized the interspace with standard, random ONB and random initializations. For details on these different initialization schemes, see Sec. G.

Experimental results can be seen in Fig. A1(a) for a VGG16 trained on CIFAR-10. For lower pruning rates, starting with \mathcal{B} and a random ONB behaves similar. For high pruning rates, random ONBs are even better suited to be used. With them, the forward and backward dynamics of a pruned network are not impaired by spatially sparse filters

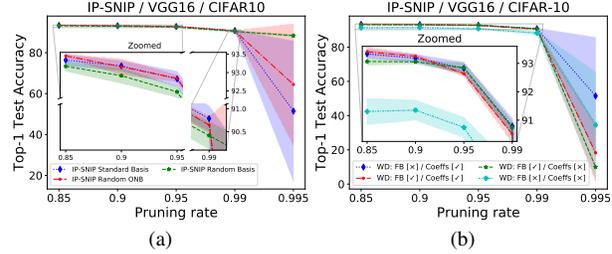


Figure A1. Both: IP-SNIP for a VGG16 trained on CIFAR-10. (a) Standard, random ONB and random initialization are compared. (b) Weight decay applied [✓] / not applied [×] on FBs and FB coefficients.

at the beginning of training. Using non-orthonormal FBs leads to worse results than ONBs for lower pruning rates. Elements of a random basis are likely to be more similar to each other than those of ONBs. This redundancy worsens performance for lower pruning rates, but significantly improves results for higher sparsity.

B.2. Top-5 accuracy for PaI on ImageNet

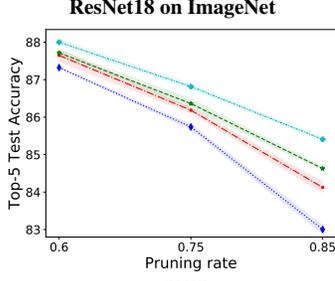
Figures A2(a), (b) and (c) show the top-5 test accuracies for the PaI ImageNet experiment with a ResNet18. Using IP instead of SP again improves results significantly as already shown and discussed for top-1 test accuracies in Figs 3(b), (d) and (f) and Sec. 5.2, respectively. Similar to the top-1 accuracy, random PaI reaches better top-5 results than SNIP, GraSP and SynFlow.

B.3. Impact of weight decay

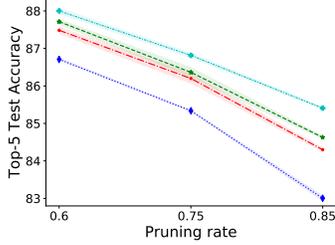
Weight decay (WD) [39] reduces the network’s capacity by shrinking parameters smoothly during training. Due to the bias-variance trade-off [24], WD can help to increase the network’s generalization ability. To find the best way to combine WD and IP, we tested all combinations of WD turned on/off for FBs and their coefficients. For this purpose, we used IP-SNIP on VGG16 and CIFAR-10, see Fig. A1(b). For lower p , not using WD at all yields the worst performance whereas the best results are obtained by applying WD on both, FBs and FB coefficients. For higher p , applying WD on the FBs reduces the network’s capacity too much. On average, using WD on the FB coefficients but not on the FBs themselves leads to the best results.

B.4. Similarity of filter bases

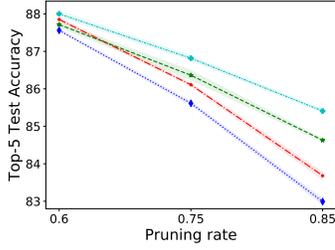
In Fig. A3, the development of the \cos similarity for the coarse FB \mathcal{F} is tracked at training time for different pruning rates for IP-SNIP with a VGG16 trained on CIFAR-10. For \mathcal{F} , random initialization is used. The \cos similarity of \mathcal{F} is the sum of all absolute values of \cos similarities of



(a) SNIP



(b) GraSP



(c) SynFlow

◆ SP-Pal ■ IP-Pal ★ SP-Random-Pal ◆ IP-Random-Pal

Figure A2. Comparison between top-5 test accuracies for SP on SNIP (a), GraSP (b) and SynFlow (c), and their adapted IP methods for a ResNet18 on ImageNet.

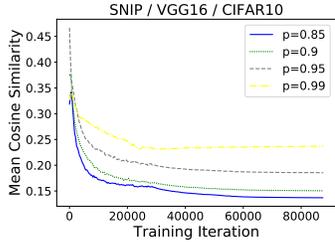


Figure A3. The mean cosine similarity of all elements of the coarse FB for a VGG16 pruned with IP-SNIP and trained on CIFAR-10 is tracked over training for varying pruning rates.

distinct elements in \mathcal{F} , i.e.

$$\frac{2}{K^2 \cdot (K^2 - 1)} \sum_{j=1}^{K^2} \sum_{k=j+1}^{K^2} \frac{|(g^{(j)}, g^{(k)})|}{\|g^{(j)}\|_2 \cdot \|g^{(k)}\|_2}. \quad (\text{A.1})$$

It therefore measures how similar two elements in \mathcal{F} are on average. Figure A3 shows that the bases have approximately the same similarity at the beginning of training for

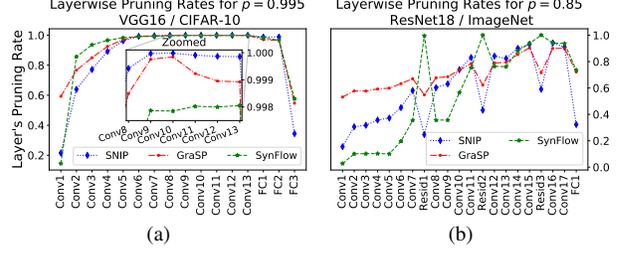


Figure A4. Layerwise pruning rates for PaI methods on (a) VGG16 on CIFAR-10 and (b) ResNet18 on ImageNet.

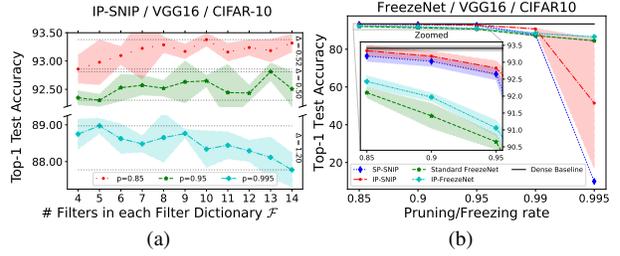


Figure A5. (a): Different sizes of FDs for varying pruning rates p . (b): IP and SP versions of freezing parameters compared to pruning them. Frozen/pruned parameters are selected before training by the SNIP criterion.

all pruning rates. For lower pruning rates, the final similarity is much smaller than for higher ones. Therefore, we assume that increasing the FB to more than 9 filters for lower pruning rates might reduce the number of needed FB coefficients, as there is “enough space” left between the 9 basis filters. On the other hand, for high pruning rates we should be able to reduce the elements in the FB, since the basis elements tend to assimilate, *i.e.* “do not need the whole space”. Experimental justifications of these assumptions are shown in Sec. B.6.

B.5. Layerwise pruning rates for PaI

As shown in Fig. A4(a), SNIP has the problem of pruning big layers too much. For the VGG16, convolutional layers 9 and 10 are pruned almost completely. This will lead to a vanishing gradient, see Fig. A6(a). With IP, the gradient flow can be increased, but if a layer is pruned completely, even an adaptive basis can not repair the damage.

SynFlow tends to fully prune 1×1 convolutional residual connections in ResNets. As shown in Fig. A4(b), all three residual connections are pruned completely. Consequently, IP- and SP-SynFlow show worse results than IP-/SP-SNIP and GraSP for ResNets, see for example Fig. 3(f). Both, SNIP and GraSP prune residual connections even less than surrounding layers.

B.6. Generalizing filter bases

Up to now, we discussed experiments where FBs \mathcal{F} formed bases. But, the spanning system $\mathcal{F} \subset \mathbb{R}^{K \times K}$ does not need to form a basis. The interspace can also be spanned by an overcomplete \mathcal{F} , *i.e.* $\#\mathcal{F} > K^2$ or an undercomplete \mathcal{F} with $\#\mathcal{F} < K^2$. This leads to the more generalized formulation of *filter dictionaries* (FDs) which include all sizes of $\#\mathcal{F}$. Of course, a FB defines a FD with $\#\mathcal{F} = K^2$ elements which are additionally assumed to be linearly independent.

As discussed in Sec. D, undercomplete FDs can be used to reduce the number of computations needed for a 2D FB convolution. However, overcomplete FDs might lead to representations of filters needing less coefficients, see [9, 14, 80]. It is not clear which elements of a basis \mathcal{B} should be removed to obtain an undercomplete FD, or added for overcomplete ones. Thus, we initialized all elements of the FDs randomly in this experiment.

A VGG16 contains 3×3 filters, thus a FB has 9 filters. Figure A5(a) shows IP-SNIP for a VGG16 trained on CIFAR-10. Reported results are those with the best validation accuracy from coarse, medium and fine FD sharing. Here, p measures the pruning rate for IP with a FB, *i.e.* $\mathcal{F} = 9$. For $\#\mathcal{F} \neq 9$, the number of non-zero FD coefficients is equal to $\#\mathcal{F} = 9$. Thus, the representation of a filter in the interspace spanned by its dictionary is more sparse if $\#\mathcal{F} > 9$ and less sparse if $\#\mathcal{F} < 9$ compared to $\#\mathcal{F} = 9$.

More than 9 elements in a FD improve results if coefficients are not too sparse, *e.g.* $\#\mathcal{F} = 10$ for $p = 0.85$ or $\#\mathcal{F} = 13$ for $p = 0.95$. Using more sophisticated methods to determine initial FDs might help to exploit overcomplete FDs better. Since $\#\mathcal{F} > 9$ increases the sparsity of FD coefficients, the performance for high pruning rates drops drastically for overcomplete FDs compared to bases.

If undercomplete FDs are used, performance worsens for lower pruning rates. Here, the capacity of the network is too low as the interspace is only $\#\mathcal{F}$ dimensional. Due to only few non-zero FB coefficients, this is not a limiting factor for high pruning rates anymore. The reduced dimensionality of the interspace even increases performance compared to $\#\mathcal{F} \geq 9$. A reason for this might be the increased information flow induced by a denser structure of the interspace. The best result for $p = 0.995$, with test accuracy 88.9%, is achieved with $\#\mathcal{F} = 5$. In comparison, SP-SNIP has 10.0% test accuracy for the same number of non-zero parameters.

B.7. Freezing coefficients

FreezeNet [85] is closely related to pruning before training via SNIP [41]. FreezeNet trains the same parameters as SNIP but *freezes* the un-trained coefficients during training instead of pruning them. By using pseudo random initializations for the network, the frozen coefficients do not have

Pruning rate	Size in kB	$\frac{\text{sparse size}}{\text{dense size}} \cdot 100\%$	Sparse & mask
Dense training	53, 256	—	—
0.2	60, 765	114.1	82.3
0.35	49, 345	92.7	67.9
0.6	30, 446	57.2	43.0
0.85	11, 461	21.5	16.9
0.995	410	0.8	0.6

Table A4. Compression for the PaI method IP-SNIP after training a VGG16 on CIFAR-10. All stored network parameters are in full precision, *i.e.* 32bit floating points. Sparse networks are stored in the CSR format whereas the dense one is stored raw. Moreover, dense and sparse networks are compressed by using `numpy.savez_compressed`. *Sparse and mask* denotes the percentage of the theoretically needed memory if only sparse parameters are stored together with the entropy encoded pruning mask.

to be stored after training but can be recovered with the used random seed. By always guaranteeing a strong gradient signal, FreezeNet outperforms SNIP significantly for low numbers of trained parameters as shown in Fig. A5(b). The opposite is true if more parameters are trained.

We further compare freezing of spatial coefficients, *standard FreezeNet*, and freezing interspace coefficients, *IP-FreezeNet*. Using adaptive FBs instead of freezing the spatial coefficients again significantly improves performance. Thus, improvements induced by interspace representations are not limited to pruning but also hold for other dimensionality reductions like freezing parts of a CNN during training.

C. Storing unstructured sparse networks

Storing sparse parameters in formats such as the *compressed sparse row format* (CSR) [78] creates additional overhead. The CSR format stores all non-zero elements of a matrix together with an array that contains the column indices and an additional array with the number of elements in each row. Therefore, additional parameters have to be stored for each non-zero element to determine the corresponding column- and row index. However, the two additional arrays do not need to be stored in 32bit full precision, but only as integers. The CSR format can be used for efficiently computing sparse matrix vector products which we also used for determining the sparse speed up for IP and SP, see Sec. D.4.

We empirically tested the overhead for real memory costs of sparse networks stored in the CSR format, see Tab. A4. Note, IP or SP pruned networks have, up to some insignificant differences, equal memory costs in practice and theory. Therefore, we report IP pruned networks in Tab. A4. Training 0.5% of all parameters compressed the network to 0.8% of the dense network’s size for IP-SNIP with a VGG16 [72] trained on CIFAR-10. Of course, for such a small number of non-zero elements, the overhead of the CSR format is

also quite small. For pruning 85% of the parameters, 21.8% of the dense memory is needed. As can be seen, additional index memory for sparse row formats increases with a decreasing pruning rate. Thus, for $p \leq 0.5$ CSR will not lead to good compression results and finally even lead to a higher memory requirement than storing the network in a dense format. As shown in Fig. A6(b), using the CSR format for such low pruning rates does not significantly speed up the network inference.

Therefore, other formats for storing the sparse network can be used for lower pruning rates. By storing the pruning mask via entropy encoding, *e.g.* [87], at most 1bit is needed for each mask parameter. To be exact, storing the network’s pruning mask for a pruning rate $p \in (0, 1)$ ideally needs

$$1 \geq S = -p \cdot \log_2 p - (1-p) \cdot \log_2(1-p) \text{ bits} \quad (\text{A.2})$$

for each element in the mask. If the mask is known, only the non-zero parameters have to be stored in the right order and in full precision. Thus, storing the sparse network of total size d with pruning rate p needs, in the ideal case, $d \cdot (S + (1-p) \cdot 32)$ bits, compared to $d \cdot 32$ bits for the dense network. In total, using entropy encoding for the pruning mask compresses the sparse network to $S/32 + (1-p)$ of its original size. As shown in Tab. A4, storing the pruning mask together with the non-zero coefficients is cheaper than CSR for all pruning rates.

D. Comparing computational costs for convolutions with spatial and interspace representations

For simplicity we will do the analysis with a FB \mathcal{F} consisting of K^2 elements in the following. But it is straight forward to do similar computations with an arbitrary FD \mathcal{F} of size N .^{A.1} As a results, all computational costs for the interspace setting are multiplied by a factor N/K^2 to get the costs for the arbitrary FD case.^{A.2} This shows that computations for IP are more expensive if an overcomplete FD with $\#\mathcal{F} > K^2$ is used. On the other hand, by reducing the size of a FD, the computations can be sped up.

In this Section, we determine the number of FLOPs needed to evaluate a standard 2D convolutional layer and a FB 2D convolutional layer. We use FLOPs as a measure since they are easy to determine and replicable in a mathematical framework but can also be measured in real time applications. A FLOP corresponds to either a multiplication or a summation.

For the forward pass, we show that the number of required FLOPs is increased by a small, constant amount

^{A.1}Summing from 1 to N instead of K^2 or doing needed computations N times instead of K^2 times.

^{A.2}Except the costs for computing $\frac{\partial \mathcal{L}}{\partial h}$ needed to update \mathcal{F} which are equal for all sizes of \mathcal{F} .

for FB-CNNs compared to standard CNNs for all pruning rates. Since the FB formulation can easily be converted to a standard representation, dense FB-CNNs therefore could be transformed into standard CNNs after training. If a CNN is pruned, this transformation is not advisable since it usually destroys the sparsity of the network.

In the backward pass, similar results hold. Moreover, we need to compute the gradient of the FB which of course requires additional resources in the IP setting.

In the following, we will assume the convolutions to have quadratic $K \times K$ kernels as well as no zero padding, stride 1×1 and dilation 1×1 .

D.1. Computations in the forward pass

D.1.1 Standard convolution.

Let $h = (h^{(\alpha,\beta)})_{\alpha,\beta} \in \mathbb{R}^{c_{out} \times c_{in} \times K \times K}$ denote a convolutional layer of a CNN. Furthermore, let $X = (X^{(\beta)})_{\beta} \in \mathbb{R}^{c_{in} \times h \times w}$ be the input feature map of the corresponding layer. In the following, we determine the number of FLOPs needed to evaluate this layer. In order to do so, we first analyze the costs for one cross-correlation \star , used in practice to compute 2D convolutional layer [1, 61], *i.e.*

$$\begin{aligned} h^{(\alpha,\beta)} \star X^{(\beta)} &= \left((h^{(\alpha,\beta)} \star X^{(\beta)})_{i,j} \right)_{i,j} \quad (\text{A.3}) \\ &= \left(\sum_{m,n=1}^K h_{m,n}^{(\alpha,\beta)} \cdot X_{m+i,n+j}^{(\beta)} \right)_{i,j} \in \mathbb{R}^{d_1 \times d_2}, \quad (\text{A.4}) \end{aligned}$$

with $d_1 := h + 1 - K$ and $d_2 := w + 1 - K$, the dimensions of the output. Equation (A.4) shows that the cost for one cross-correlation is given by $2 \cdot K^2 \cdot d_1 \cdot d_2$ FLOPs. The output of a 2D convolutional layer is given by

$$Y = \left(Y^{(\alpha)} \right)_{\alpha} = \left(\sum_{\beta=1}^{c_{in}} h^{(\alpha,\beta)} \star X^{(\beta)} \right)_{\alpha} \in \mathbb{R}^{c_{out} \times d_1 \times d_2}, \quad (\text{A.5})$$

which finally leads to $c_{out} \cdot c_{in}$ times the costs to compute a single cross-correlation Eq. (A.4). Therefore, $2 \cdot c_{out} \cdot c_{in} \cdot K^2 \cdot d_1 \cdot d_2$ FLOPs are needed in total to compute a standard 2D convolutional layer.

D.1.2 FB convolution.

Let $h = (h^{(\alpha,\beta)})_{\alpha,\beta} = \left(\sum_n \lambda_n^{(\alpha,\beta)} \cdot g^{(n)} \right)_{\alpha,\beta} \in \mathbb{R}^{c_{out} \times c_{in} \times K \times K}$ be the interspace representation of h , where the FB is given by $\mathcal{F} = \{g^{(1)}, \dots, g^{(K^2)}\} \subset \mathbb{R}^{K \times K}$. The 2D convolution of this layer with input $X \in \mathbb{R}^{c_{in} \times h \times w}$ can

be computed via

$$Y = \left(Y^{(\alpha)} \right)_\alpha = \left(\sum_{\beta=1}^{c_{in}} \sum_{n=1}^{K^2} \lambda_n^{(\alpha,\beta)} \cdot \left(g^{(n)} \star X^{(\beta)} \right) \right)_{\alpha} \quad (\text{A.6})$$

Using the last equation in Eq. (A.6), we see that $g^{(n)} \star X^{(\beta)}$ has to be computed once for each combination of β and n , i.e. $c_{in} \cdot K^2$ many times. The costs for computing all $g^{(n)} \star X^{(\beta)}$ is therefore given by $2 \cdot c_{in} \cdot K^4 \cdot d_1 \cdot d_2$ FLOPs. For each combination of α, β and n , $g^{(n)} \star X^{(\beta)}$ has to be multiplied by the scalar $\lambda_n^{(\alpha,\beta)}$. These are $d_1 \cdot d_2$ many FLOPs for each α, β and n . Summing over β and n yields another $c_{out} \cdot c_{in} \cdot K^2 \cdot d_1 \cdot d_2$ FLOPs in total. Thus, the total costs for computing a FB 2D convolutional layer is given by $2 \cdot c_{out} \cdot c_{in} \cdot K^2 \cdot d_1 \cdot d_2 + 2 \cdot c_{in} \cdot K^4 \cdot d_1 \cdot d_2$ FLOPs.

By using FB convolutions, the numbers of needed FLOPs is therefore slightly increased by $2 \cdot c_{in} \cdot K^4 \cdot d_1 \cdot d_2$. Which is a relative increase of $K^2/c_{out} \cdot 100\%$ compared to the standard case.

D.1.3 Pruned networks.

In the following we assume the convolutional layer $h \in \mathbb{R}^{c_{out} \times c_{in} \times K \times K}$ to be pruned with a pruning rate of $p \in [0, 1]$.^{A.3} We suppose all zero coefficients to be known. Thus, the corresponding multiplications do not have to be computed in Eqs. (A.4) and (A.6).

The required number of computations for a standard 2D convolutional layer with pruning rate p is therefore given by

$$2 \cdot c_{out} \cdot c_{in} \cdot K^2 \cdot d_1 \cdot d_2 \cdot (1 - p) \text{ FLOPs} . \quad (\text{A.7})$$

For a pruned FB 2D convolutional layer,

$$2 \cdot c_{out} \cdot c_{in} \cdot K^2 \cdot d_1 \cdot d_2 \cdot (1 - p) + 2 \cdot c_{in} \cdot K^4 \cdot d_1 \cdot d_2 \text{ FLOPs} \quad (\text{A.8})$$

are needed for evaluation.

The number of FLOPs for IP is increased for all pruning rates by $2 \cdot c_{in} \cdot K^4 \cdot d_1 \cdot d_2$ compared to SP. These are exactly the costs for computing all combinations of $g^{(n)} \star X^{(\beta)}$, needed for the forward pass for FB 2D Convolutions. These costs are independent of the pruning rate and therefore a constant overhead of IP compared to SP. Thus the additional costs for IP in the forward pass compared to SP are K^2/c_{out} times the costs of the dense forward pass.

D.2. Backward Pass

Up to now, we have computed additional FLOP costs for IP compared to SP in the forward pass. Now we want to

^{A.3}For simplicity, we assume the number of non-zero coefficients for IP and SP to be equal here. Due to extra FB parameters, the number of non-zero interspace coefficients is always slightly smaller than for the standard case in our experiments.

have a closer look at the backward pass. We note that $\frac{\partial \mathcal{L}}{\partial Y}$ always has the same cost for the standard- and the FB 2D convolution layer. This holds since $\hat{X} = \sigma(Y)$ for some activation function σ and consequently $\frac{\partial \mathcal{L}}{\partial Y} = \frac{\partial \mathcal{L}}{\partial \hat{X}} \odot \sigma'(Y)$.

D.2.1 Computing the gradient for X .

Furthermore, it is known that

$$\frac{\partial \mathcal{L}}{\partial X^{(\beta)}} = \sum_{\alpha=1}^{c_{out}} h^{(\alpha,\beta)} \hat{\star} \frac{\partial \mathcal{L}}{Y^{(\alpha)}} \quad (\text{A.9})$$

with a strided convolution $\hat{\star}$ that corresponds to the forward pass and which needs $2 \cdot c_{out} \cdot c_{in} \cdot K^2 \cdot h \cdot w$ FLOPs. By representing $h^{(\alpha,\beta)} = \sum_{n=1}^{K^2} \lambda_n^{(\alpha,\beta)} \cdot g^{(n)}$ and using the linearity of $\hat{\star}$, we now get the computational overhead of $2 \cdot c_{out} \cdot K^4 \cdot h \cdot w$ FLOPs which are the costs for computing

$$g^{(n)} \hat{\star} \frac{\partial \mathcal{L}}{Y^{(\alpha)}} \quad (\text{A.10})$$

for all $n \in \{1, \dots, K^2\}$ and $\alpha \in \{1, \dots, c_{out}\}$. This results in an overhead of K^2/c_{in} compared to the costs of the standard, dense network.

In the sparse case, again the FLOP costs for SP are decreased by a factor $(1 - p)$. Furthermore, the overhead K^2/c_{in} is constant since the $g^{(n)}$ are not pruned. Similar formulas to Eqs. (A.7) and (A.8) hold also in the backpropagation case which results in a constant overhead of IP compared to SP for computing $\frac{\partial \mathcal{L}}{\partial X}$ equal to K^2/c_{in} times the costs of the dense computation of $\frac{\partial \mathcal{L}}{\partial X}$.

D.2.2 Gradients for coefficients.

The backpropagation formulas for the spatial and FB coefficients are given by

$$\frac{\partial \mathcal{L}}{\partial h_{i,j}^{(\alpha,\beta)}} = \left(\frac{\partial \mathcal{L}}{\partial Y^{(\alpha)}} \star X^{(\beta)} \right)_{i,j} \quad (\text{A.11})$$

and

$$\frac{\partial \mathcal{L}}{\partial \lambda_n^{(\alpha,\beta)}} = \left\langle \frac{\partial \mathcal{L}}{\partial Y^{(\alpha)}}, g^{(n)} \star X^{(\beta)} \right\rangle , \quad (\text{A.12})$$

respectively. Since $g^{(n)} \star X^{(\beta)}$ is already computed in the forward pass, both computations for the standard case and the FB representation have equal FLOP costs. In total, this equals to $2 \cdot c_{out} \cdot c_{in} \cdot K^2 \cdot d_1 \cdot d_2$ FLOPs for computing $\frac{\partial \mathcal{L}}{\partial h}$ or $\frac{\partial \mathcal{L}}{\partial \lambda}$. If pruning is applied, this reduces to $2 \cdot c_{out} \cdot c_{in} \cdot K^2 \cdot d_1 \cdot d_2 \cdot (1 - p)$ FLOPs for IP and SP, since gradients for pruned coefficients do not need to be computed.

Note, if the size $d_1 \cdot d_2$ of $Y \in \mathbb{R}^{c_{out} \times d_1 \times d_2}$ is bigger than the kernel size K^2 it is even cheaper to compute the

gradient of $\lambda_n^{(\alpha,\beta)}$ via

$$\frac{\partial \mathcal{L}}{\partial \lambda_n^{(\alpha,\beta)}} = \left\langle g^{(n)}, \frac{\partial \mathcal{L}}{\partial Y^{(\alpha)}} \star X^{(\beta)} \right\rangle. \quad (\text{A.13})$$

In Eq. (A.12) there are $2 \cdot d_1 \cdot d_2$ FLOPs needed (if $g^{(n)} \star X^{(\beta)}$ is known which we can assume due to the forward pass) whereas Eq. (A.13) needs $2 \cdot K^2$ FLOPs if $\frac{\partial \mathcal{L}}{\partial Y^{(\alpha)}} \star X^{(\beta)}$ is known. As we will see in the following, $g^{(n)}$ needs $\frac{\partial \mathcal{L}}{\partial Y^{(\alpha)}} \star X^{(\beta)}$ to be computed for all α, β and consequently we can assume them to be known. In summary we can say that the computation of the interspace coefficients λ requires the same number of FLOPs, or even less, compared to the spatial coefficients.

D.2.3 Gradient for the filter base.

The computation of the gradients $\frac{\partial \mathcal{L}}{\partial g^{(n)}}$ also generates extra costs for the backward pass of training interspace representations. It holds

$$\frac{\partial \mathcal{L}}{\partial g^{(n)}} = \sum_{\alpha=1}^{c_{out}} \sum_{\beta=1}^{c_{in}} \lambda_n^{(\alpha,\beta)} \cdot \left(\frac{\partial \mathcal{L}}{\partial Y^{(\alpha)}} \star X^{(\beta)} \right). \quad (\text{A.14})$$

As shown in Eq. (A.14), $\frac{\partial \mathcal{L}}{\partial g^{(n)}}$ first needs to compute all $\frac{\partial \mathcal{L}}{\partial Y^{(\alpha)}} \star X^{(\beta)} = \frac{\partial \mathcal{L}}{\partial h^{(\alpha,\beta)}}$. This is exactly the cost for computing the dense gradient $\frac{\partial \mathcal{L}}{\partial h}$ which needs $2 \cdot c_{in} \cdot c_{out} \cdot K^2 \cdot d_1 \cdot d_2$ FLOPs. The scaling and summation in the sum Eq. (A.14) requires $2 \cdot c_{in} \cdot c_{out} \cdot K^2$ FLOPs in total. If pruning is applied, this reduces to $2 \cdot c_{in} \cdot c_{out} \cdot K^2 \cdot (1-p)$. Altogether, computing the gradients of $g^{(1)}, \dots, g^{(K^2)}$ needs $2 \cdot c_{out} \cdot c_{in} \cdot K^2 \cdot (d_1 \cdot d_2 + (1-p) \cdot K^2)$ FLOPs. In simple terms, the total computation of $\frac{\partial \mathcal{L}}{\partial g}$ lies in $\mathcal{O}(\text{costs}(\frac{\partial \mathcal{L}}{\partial h}))$.

D.2.4 Summary for the backward pass.

In summary, the computation of $\frac{\partial \mathcal{L}}{\partial X}$ of IP induces a constant overhead compared to IP. This corresponds to K^2/c_{in} times the costs of computing the dense gradient of $\frac{\partial \mathcal{L}}{\partial X}$ by using spatial coefficients. On top of that, IP also needs to compute the gradient for the FB \mathcal{F} which is in $\mathcal{O}(\text{costs}(\frac{\partial \mathcal{L}}{\partial h}))$.

D.3. Upper bounds for gradients

As Eq. (A.13) and Eq. (A.14) show, jointly optimizing \mathcal{F} and λ leads to non trivial correlations between them. With a slight abuse of notation we assume for the following discussion \mathcal{F} to be the $K^2 \times K^2$ matrix containing all flattened $g^{(n)}$. Further, let $h, \lambda \in \mathbb{R}^{K^2 \times c_{out} c_{in}}$ contain all spatial and interspace coefficients of the layer, respectively. Therefore, it holds $h = \mathcal{F} \cdot \lambda$. By using $\frac{\partial \mathcal{L}}{\partial h^{(\alpha,\beta)}} = \frac{\partial \mathcal{L}}{\partial Y^{(\alpha)}} \star X^{(\beta)}$ and the Cauchy-Schwartz inequality, the gradients for \mathcal{F} and λ

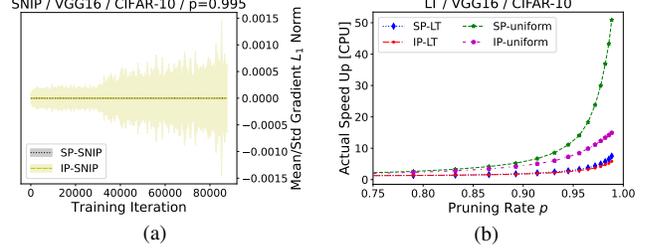


Figure A6. (a) Mean and std of gradient L_1 norm for IP- and SP-SNIP for $p = 0.995$. (b) Speed up on CPU for varying p for SP and IP for uniform sparsity and LT's sparsity.

are bounded by

$$\left\| \frac{\partial \mathcal{L}}{\partial \lambda} \right\|_F \leq \|\mathcal{F}\|_F \left\| \frac{\partial \mathcal{L}}{\partial h} \right\|_F \quad \text{and} \quad \left\| \frac{\partial \mathcal{L}}{\partial \mathcal{F}} \right\|_F \leq \|\lambda\|_F \left\| \frac{\partial \mathcal{L}}{\partial h} \right\|_F. \quad (\text{A.15})$$

This shows that upper bounds for $\frac{\partial \mathcal{L}}{\partial \mathcal{F}}$ and $\frac{\partial \mathcal{L}}{\partial \lambda}$ are determined by the spatial gradient $\frac{\partial \mathcal{L}}{\partial h}$. This boundedness of the gradients leads to stable convergence for both, \mathcal{F} and λ , while the convergence behavior of λ is similar to the standard coefficients h , see Fig. 5(b). Moreover, Fig. A6(a) even shows that adaptive FBs help to overcome vanishing gradients for SNIP by becoming spatially dense. IP-SNIP can use that to recover during training from a complete, PaI induced information loss, while SP-SNIP is stuck with zero gradient flow.

D.4. Real runtime measurements

To measure and compare the real runtime accelerations of IP and SP for inference, we used `scipy`'s sparse package. To be precise, we used `scipy.sparse.csr_matrix`, see [online documentation](#). As discussed in Sec. D.1, sparse FB convolutions can be computed by first convolving all $X^{(\beta)}$ with all $g^{(n)}$. Afterwards, sparse matrix multiplications can be used to compute the actual output Y . To rule out runtime differences induced by mismatches between sparse implementations of matrix multiplications and convolutions, we simulated sparse convolutions with sparse matrix multiplications of matching dimensions. Therefore, a sparse convolution $h \star X$ with $h \in \mathbb{R}^{c_{out} \times c_{in} \times K \times K}$ and $X \in \mathbb{R}^{c_{in} \times H \times W}$ corresponds to a matrix-vector multiplication $\hat{h} \cdot \hat{X}$ with $\hat{h} \in \mathbb{R}^{c_{out} \cdot H \cdot W \times c_{in} \cdot K^2}$ and $\hat{X} \in \mathbb{R}^{c_{in} \cdot K^2}$.

We measured the runtime of a VGG16 on input images $X \in \mathbb{R}^{3 \times 32 \times 32}$ (i.e. CIFAR-10) with two different sparsity configurations, the sparsity distribution found by pruning with LT, see Sec. 5.3, and uniform sparsity for each layer. For simplicity, we omit the batch normalization layers and non-linearities. Runtime is measured on one core of an Intel XEON E5-2680 v4 2.4 GHz CPU where we used batch size 1 and the mean runtime of 25 runs.

Figure Fig. A6(b) shows the comparison between model sparsity and the actual runtime speed up on a CPU. Since the

used CSR [78] format for sparse coefficients adds additional overhead to the actually executed computations, runtime is sped up significantly only for $p \geq 0.75$. Note, different sparsity distributions can lead to varying accelerations for a similar global pruning rate p . IP indeed has a longer runtime due to the mentioned extra computations. But by boosting performance of sparse models, IP reaches similar results than dense training with 5.2 times speed up and better results than SP for equal runtime, as shown in Fig. 5(a).

E. Transformation rules in the interspace

Since the interspace representation is obtained by a linear transformation of the standard, spatial representation, we will derive formulas for this transformation. By knowing them, it will be straight forward to also determine transformation rules for the corresponding gradients and higher derivatives. Those transformation rules might be useful if pruning methods that need first or second order information are used. We test three methods in our work that need the information of the gradient, SNIP [41], GraSP [82] and SynFlow [76]. Moreover, GraSP needs second order information as well. Since all these methods are applied at initialization and we use an initialization equivalent to the standard network in the main part of this work, the gradient and Hessian are equivalent at that time. Still, if such pruning methods are applied with different initializations, the knowledge of these transformation rules might be helpful to overcome scaling problems. Furthermore, we believe the transformation rules to be fruitful for analyzing the information flow in FB-CNNs which we think is an interesting direction for future work.

Again we will assume the special case considered in the paper, *i.e.* \mathcal{F} forming a basis.

E.1. Transformation rules for filters

For a given layer in a CNN, let α denote the output channel, β the corresponding input channel and $K \times K$ be the kernel size of a filter $h^{(\alpha,\beta)}$. For the layer's FB $\mathcal{F} = \{g^{(1)}, \dots, g^{(K^2)}\}$, the filter's interspace representation is given by

$$h^{(\alpha,\beta)} = \sum_{n=1}^{K^2} \lambda_n^{(\alpha,\beta)} \cdot g^{(n)}. \quad (\text{A.16})$$

Here, the FB coefficients of $h^{(\alpha,\beta)}$ are given by $\lambda^{(\alpha,\beta)} = (\lambda_n^{(\alpha,\beta)})_n \in \mathbb{R}^{K^2}$. Let \mathcal{B} be the standard basis for $\mathbb{R}^{K \times K}$. Then, the spatial representation of filter $h^{(\alpha,\beta)}$ is given by

$$h^{(\alpha,\beta)} = \sum_{n=1}^{K^2} h_{i_n, j_n}^{(\alpha,\beta)} \cdot e^{(n)} = \sum_{n=1}^{K^2} \varphi_n^{(\alpha,\beta)} \cdot e^{(n)}, \quad (\text{A.17})$$

with spatial coefficients $\varphi^{(\alpha,\beta)} := (\varphi_n^{(\alpha,\beta)})_n = ((h^{(\alpha,\beta)}, e^{(n)}))_n \in \mathbb{R}^{K^2}$ and standard basis \mathcal{B} given by

$$\mathcal{B} = \{e^{(1)}, \dots, e^{(K^2)}\} \text{ and}$$

$$e_{i,j}^{(n)} = \delta_{i,i_n} \cdot \delta_{j,j_n}, (i_n, j_n) \in \{1, \dots, K\}^2, \\ (i_n, j_n) \neq (i_m, j_m) \text{ for } n \neq m. \quad (\text{A.18})$$

Consequently,

$$\varphi^{(\alpha,\beta)} = \Psi \cdot \lambda^{(\alpha,\beta)}, \Psi = \left(\langle g^{(m)}, e^{(n)} \rangle \right)_{n,m} \in \mathbb{R}^{K^2 \times K^2} \quad (\text{A.19})$$

holds. Note, since FBs are shared for at least one layer, the basis transformation matrix Ψ is not labeled with the input- and output channels α and β , respectively. But of course, formulas can be adapted to the case of more than one FB per layer. Since we assume \mathcal{F} to form a basis, the reverse is given by

$$\lambda^{(\alpha,\beta)} = \Psi^{-1} \cdot \varphi^{(\alpha,\beta)}. \quad (\text{A.20})$$

Note, if we use \mathcal{F} as a general dictionary, and not a basis anymore, a reverse can still be computed by the Moore-Penrose pseudo inverse $\Psi^\dagger = \Psi^T \cdot (\Psi \cdot \Psi^T)^{-1}$ if \mathcal{F} forms a *generating system* for $\mathbb{R}^{K \times K}$. If \mathcal{F} forms a linear independent, undercomplete dictionary, we can express $\lambda^{(\alpha,\beta)} = \hat{\Psi} \varphi^{(\alpha,\beta)}$ for a suitable $\hat{\Psi} \in \mathbb{R}^{\#\mathcal{F} \times K^2}$. A reverse of this is given by $\varphi^{(\alpha,\beta)} = \hat{\Psi}^\dagger \lambda^{(\alpha,\beta)}$, where again $\hat{\Psi}^\dagger = \hat{\Psi}^T \cdot (\hat{\Psi} \cdot \hat{\Psi}^T)^{-1}$ forms the Moore-Penrose pseudo inverse.

E.2. Transformation rules for gradients

Let \mathcal{L} denote the loss function used to train the CNN. Assuming $h^{(\alpha,\beta)}$ to be given in the interspace representation Eq. (A.16),

$$\frac{\partial \mathcal{L}}{\partial \lambda^{(\alpha,\beta)}} = \left(\frac{\partial \varphi^{(\alpha,\beta)}}{\partial \lambda^{(\alpha,\beta)}} \right)^T \cdot \frac{\partial \mathcal{L}}{\partial \varphi^{(\alpha,\beta)}} = \Psi^T \cdot \frac{\partial \mathcal{L}}{\partial \varphi^{(\alpha,\beta)}} \quad (\text{A.21})$$

holds by the chain rule and Eq. (A.19). Consequently,

$$\frac{\partial \mathcal{L}}{\partial \varphi^{(\alpha,\beta)}} = (\Psi^{-1})^T \cdot \frac{\partial \mathcal{L}}{\partial \lambda^{(\alpha,\beta)}} \quad (\text{A.22})$$

is true for the gradient as well.

By comparing Eq. (A.20) with Eq. (A.21) (or Eq. (A.19) with Eq. (A.22)) we see that the coefficients and their corresponding gradients transform complementary to each other if $\Psi^T \neq \Psi^{-1}$. Note, $\Psi^{-1} = \Psi^T$ if and only if Ψ is orthonormal which is equivalent to \mathcal{F} forming an orthonormal basis (ONB).

E.3. Transformation rules for Hessian

In order to compute the Hessian of the loss function, we have to index all possible filters in a CNN. Let $\lambda^{(\alpha,\beta;l)}$ and $\varphi^{(\alpha,\beta;l)}$ denote the interspace and spatial coefficients of a filter in layer l corresponding to input channel β and output channel α . Here, the basis transformation in layer l is given by $\Psi^{(l)}$. If a FB is shared for layers l and $l+1$, then

$\Psi^{(l)} = \Psi^{(l+1)}$ would hold. Furthermore, $K_{(l)}$ is the filter size in layer l and $c_{out}^{(l)}$ and $c_{in}^{(l)}$ denote the number of output and input channels, respectively. We assume the CNN to have L_c convolutional layers in total. Let $H^{\mathcal{B}}$ be the Hessian matrix of \mathcal{L} w.r.t. to coefficients of \mathcal{B} . The corresponding values of the Hessian are given by

$$H^{\mathcal{B}}((\alpha, \beta, n; l), (\alpha', \beta', n'; l')) = \frac{\partial^2 \mathcal{L}}{\partial \varphi_n^{(\alpha, \beta; l)} \partial \varphi_{n'}^{(\alpha', \beta'; l')}}. \quad (\text{A.23})$$

Equivalently, the Hessian w.r.t. \mathcal{F} is given by $H^{\mathcal{F}}$ with values

$$H^{\mathcal{F}}((\alpha, \beta, n; l), (\alpha', \beta', n'; l')) = \frac{\partial^2 \mathcal{L}}{\partial \lambda_n^{(\alpha, \beta; l)} \partial \lambda_{n'}^{(\alpha', \beta'; l')}}. \quad (\text{A.24})$$

Using multi-index notation, we can describe the transformation of the Hessian matrix compactly. For $\mathbf{m} := (\alpha, \beta, n; l)$ and $\mathbf{m}' := (\alpha', \beta', n'; l')$, define

$$\varphi_{\mathbf{m}} := \varphi_n^{(\alpha, \beta; l)}, \quad \varphi := (\varphi_{\mathbf{m}})_{\mathbf{m} \in \mathcal{M}} \quad (\text{A.25})$$

$$\lambda_{\mathbf{m}} := \lambda_n^{(\alpha, \beta; l)}, \quad \lambda := (\lambda_{\mathbf{m}})_{\mathbf{m} \in \mathcal{M}} \quad (\text{A.26})$$

$$\Psi_{\mathbf{m}, \mathbf{m}'} := \delta_{\alpha, \alpha'} \cdot \delta_{\beta, \beta'} \cdot \delta_{l, l'} \cdot \Psi_{n, n'}^{(l)}, \quad (\text{A.27})$$

$$\Psi := (\Psi_{\mathbf{m}, \mathbf{m}'})_{\mathbf{m}, \mathbf{m}' \in \mathcal{M}}, \quad (\text{A.28})$$

where all possible multi-indices are given by

$$\mathcal{M} := \bigcup_{l=1}^{L_c} \bigcup_{\alpha=1}^{c_{out}^{(l)}} \bigcup_{\beta=1}^{c_{in}^{(l)}} \bigcup_{n=1}^{K_{(l)}} \{(\alpha, \beta, n; l)\}. \quad (\text{A.29})$$

Let $d := \#\mathcal{M}$ be the dimension of the CNN. For all matrices/vectors $A \in \mathbb{R}^{d_1 \times d}$ and $B \in \mathbb{R}^{d \times d_2}$ with $d_1, d_2 \in \{1, d\}$, indexed with multi-indices, we define multi-index multiplication via

$$(A \cdot B)_{\mathbf{m}, \mathbf{m}'} := \sum_{\mathbf{m}'' \in \mathcal{M}} A_{\mathbf{m}, \mathbf{m}''} \cdot B_{\mathbf{m}'', \mathbf{m}'}. \quad (\text{A.30})$$

Using the multi-index notation, together with Eqs. (A.25), (A.26) and (A.28), leads to simple transformations of coefficients and their gradients for the whole CNN, given by

$$\varphi = \Psi \cdot \lambda, \quad \frac{\partial \mathcal{L}}{\partial \varphi} = (\Psi^{-1})^T \cdot \frac{\partial \mathcal{L}}{\partial \lambda} \quad (\text{A.31})$$

with the transpose of a multi-index matrix A defined via

$$A_{\mathbf{m}, \mathbf{m}'}^T := A_{\mathbf{m}', \mathbf{m}} \quad (\text{A.32})$$

and

$$\Psi^{-1} := (\Psi_{\mathbf{m}, \mathbf{m}'}^{-1})_{\mathbf{m}, \mathbf{m}' \in \mathcal{M}} \quad (\text{A.33})$$

$$\Psi_{\mathbf{m}, \mathbf{m}'}^{-1} := \delta_{\alpha, \alpha'} \cdot \delta_{\beta, \beta'} \cdot \delta_{l, l'} \cdot (\Psi_{n, n'}^{(l)})^{-1}. \quad (\text{A.34})$$

It holds

$$H_{\mathbf{m}, \mathbf{m}'}^{\mathcal{B}} = \frac{\partial^2 \mathcal{L}}{\partial \varphi_{\mathbf{m}} \partial \varphi_{\mathbf{m}'}} \quad (\text{A.35})$$

$$= \frac{\partial}{\partial \varphi_{\mathbf{m}}} \left(\sum_{\mathbf{m}'' \in \mathcal{M}} \Psi_{\mathbf{m}'', \mathbf{m}'}^{-1} \cdot \frac{\partial \mathcal{L}}{\partial \lambda_{\mathbf{m}''}} \right) \quad (\text{A.36})$$

$$= \sum_{\mathbf{m}'', \mathbf{m}''' \in \mathcal{M}} \Psi_{\mathbf{m}''', \mathbf{m}}^{-1} \cdot H_{\mathbf{m}''', \mathbf{m}''}^{\mathcal{F}} \cdot \Psi_{\mathbf{m}'', \mathbf{m}'}^{-1} \quad (\text{A.37})$$

$$= ((\Psi^{-1})^T \cdot H^{\mathcal{F}} \cdot \Psi^{-1})_{\mathbf{m}, \mathbf{m}'}. \quad (\text{A.38})$$

F. Computation of pruning scores

In the following, we will derive the computations of the pruning scores used in the experimental evaluation in Sec. 5 in the main paper. We will present the original scores for SP and their corresponding IP version. In this Section, we assume all FBs \mathcal{F} to form bases and fully connected layers to be described by 1×1 convolutions.

F.1. Pruning scores in general

In Secs. F.2 - F.6, five different methods for computing a pruning score vector $S \in \mathbb{R}^d$ are presented. These are the pruning scores used in our experimental evaluation.

A global pruning score means that the whole network is pruned altogether based on this score vector. Here, d denotes the number of all prunable parameters – for simplicity pooled in a big vector $\lambda \in \mathbb{R}^d$. For each parameter λ_j , there exists exactly one corresponding pruning score S_j . The higher a pruning score, the more important the corresponding parameter is. Thus, for pruning a network with prunable parameters λ to pruning rate $p \in [0, 1]$, only the k biggest entries in S are not pruned, where

$$k := \lfloor (1 - p) \cdot d \rfloor. \quad (\text{A.39})$$

Consequently, the global pruning mask $\mu \in \{0, 1\}^d$ is defined via

$$\mu_j = \begin{cases} 1, & S_j \text{ belongs to the } k \text{ biggest entries in } S \\ 0, & \text{else} \end{cases}. \quad (\text{A.40})$$

After the pruning mask is computed, the network's prunable parameters are masked with the pruning mask via $\lambda \odot \mu$.

Note, for DST methods, we use *layerwise* pruning. The pruning score is computed equivalent to the global case but each layer is pruned with an individual pruning rate. Consequently, the number of pruned parameters is computed for each layer individually.

F.2. Random pruning score

Random pruning scores are easy to obtain. For each coefficient $\varphi_n^{(\alpha, \beta; l)}$ in the SP case or $\lambda_n^{(\alpha, \beta; l)}$ in the IP case, a

corresponding random number is drawn i.i.d. from a $\mathcal{N}(0, 1)$ distribution.

F.3. Magnitude pruning score

Magnitudes are used as pruning criterion for LTs [19], the DST methods SET [55] and RigL [17], FT [68] and GMP [21]. Using magnitudes as pruning criterion assumes that big coefficients are more likely to significantly influence the network’s output than smaller ones. The corresponding formula is straight forward and given by

$$S_{Mag}(\alpha, \beta, n; l, \mathcal{B}) := \left| \varphi_n^{(\alpha, \beta; l)} \right|. \quad (\text{A.41})$$

The corresponding formula for the FB representation is given by

$$S_{Mag}(\alpha, \beta, n; l, \mathcal{F}) := \left| \lambda_n^{(\alpha, \beta; l)} \right|. \quad (\text{A.42})$$

By the transformation rules for the coefficients Eq. (A.19), it holds

$$\left\| \varphi_n^{(\alpha, \beta; l)} \right\| = \left\| \Psi^{(l)} \lambda_n^{(\alpha, \beta; l)} \right\|. \quad (\text{A.43})$$

Consequently, $K \times K$ filters $h^{(\alpha, \beta; l)}$ do not need to have the same total pruning score in the spatial and FB representation, but might be scaled differently. We used magnitude pruning during or after training, where usually $\varphi_n^{(\alpha, \beta; l)} \neq \lambda_n^{(\alpha, \beta; l)}$. Even though magnitude pruning is normally used for the spatial representation, we did not have any scaling issues in the IP setting. This again indicates that jointly optimizing the FBs and their coefficients is stable.

F.4. SynFlow pruning score

SynFlow [76] is a pruning score, calculating the contribution of a parameter to the CNN’s overall information flow. This is done, by differentiating the so called L_1 path norm of the network. The formula is given by

$$S_{SynFlow}(\alpha, \beta, n; l, \mathcal{B}) := \frac{\partial \mathcal{R}}{\partial \varphi_n^{(\alpha, \beta; l)}} \cdot \varphi_n^{(\alpha, \beta; l)}, \quad (\text{A.44})$$

with

$$\mathcal{R} := \sum_{p \in \mathcal{P}} \prod_{\varphi_n^{(\alpha, \beta; l)} \in p} |\varphi_n^{(\alpha, \beta; l)}|. \quad (\text{A.45})$$

Here,

$$\mathcal{P} = \left\{ \left\{ \varphi_{n_1}^{(\alpha_1, \beta; 1)}, \varphi_{n_2}^{(\alpha_2, \alpha_1; 2)}, \dots, \varphi_{n_{L_c}}^{(\alpha_{L_c}, \alpha_{L_c-1}; L_c)} \right\} \right\} \quad (\text{A.46})$$

describes all existing paths in a CNN which start in the input layer and end in the output layer.

Consequently, the corresponding SynFlow score w.r.t. coefficients for the FBs is given by

$$S_{SynFlow}(\alpha, \beta, n; l, \mathcal{F}) := \frac{\partial \mathcal{R}}{\partial \lambda_n^{(\alpha, \beta; l)}} \cdot \lambda_n^{(\alpha, \beta; l)}. \quad (\text{A.47})$$

The basis transformation between \mathcal{B} and \mathcal{F} does not change the total pruning score of a filter $h^{(\alpha, \beta; l)}$. This can be seen by

$$\sum_{n=1}^{K^2} S_{SynFlow}(\alpha, \beta, n; l, \mathcal{B}) \quad (\text{A.48})$$

$$= \left\langle \frac{\partial \mathcal{R}}{\partial \varphi^{(\alpha, \beta; l)}}, \varphi^{(\alpha, \beta; l)} \right\rangle \quad (\text{A.49})$$

$$= \left\langle \left(\Psi^{(l)-1} \right)^T \frac{\partial \mathcal{R}}{\partial \lambda^{(\alpha, \beta; l)}}, \Psi^{(l)} \lambda^{(\alpha, \beta; l)} \right\rangle \quad (\text{A.50})$$

$$= \left\langle \frac{\partial \mathcal{R}}{\partial \lambda^{(\alpha, \beta; l)}}, \Psi^{(l)-1} \Psi^{(l)} \lambda^{(\alpha, \beta; l)} \right\rangle \quad (\text{A.51})$$

$$= \left\langle \frac{\partial \mathcal{R}}{\partial \lambda^{(\alpha, \beta; l)}}, \lambda^{(\alpha, \beta; l)} \right\rangle \quad (\text{A.52})$$

$$= \sum_{n=1}^{K^2} S_{SynFlow}(\alpha, \beta, n; l, \mathcal{F}). \quad (\text{A.53})$$

The second equality is induced by the transformation formulas Eqs. (A.19) and (A.22). By having the same total pruning score for a filter for coefficients w.r.t. \mathcal{F} and \mathcal{B} , we do not need to worry about possible scaling issues for the SynFlow score.

F.5. SNIP pruning score

SNIP [41] computes a so called *saliency score* for each parameter of a CNN before training. The idea is to measure the effect of changing the activation of a coefficient on the loss function. If this effect is big, the corresponding coefficient is trained, otherwise it is pruned. Let $\varphi^{(\alpha, \beta; l)} \in \mathbb{R}^{K^2_{(l)}}$ be the vector consisting of all spatial coefficients in the l -th layer of a CNN with input channel β and output channel α . Its saliency score is then computed as

$$S_{SNIP}(\alpha, \beta, n; l, \mathcal{B}) := \left| \frac{\partial \mathcal{L}(m \cdot \varphi_n^{(\alpha, \beta; l)})}{\partial m} \Bigg|_{m=1} \right| \quad (\text{A.54})$$

$$= \left| \frac{\partial \mathcal{L}}{\partial \varphi_n^{(\alpha, \beta; l)}} \cdot \varphi_n^{(\alpha, \beta; l)} \right|, \quad (\text{A.55})$$

where $m \in \mathbb{R}$ models the activation of the filter value and \mathcal{L} is the used loss function. The second equality is induced by using the chain rule [82].

The corresponding SNIP score w.r.t. \mathcal{F} is given by

$$S_{SNIP}(\alpha, \beta, n; l, \mathcal{F}) := \left| \frac{\partial \mathcal{L}(m \cdot \lambda_n^{(\alpha, \beta; l)})}{\partial m} \Bigg|_{m=1} \right| \quad (\text{A.56})$$

$$= \left| \frac{\partial \mathcal{L}}{\partial \lambda_n^{(\alpha, \beta; l)}} \cdot \lambda_n^{(\alpha, \beta; l)} \right|. \quad (\text{A.57})$$

By inserting the transformation formulas (A.19) and (A.22) into Eq. (A.55), we get the relationship for the SNIP score

of a filter $h^{(\alpha,\beta;l)}$ as

$$S_{SNIP}(\alpha, \beta; l, \mathcal{B}) \quad (\text{A.58})$$

$$:= (S_{SNIP}(\alpha, \beta, n; l, \mathcal{B}))_{n=1}^{K^2} \quad (\text{A.59})$$

$$= \left| (\Psi^{(l)})^T \frac{\partial \mathcal{L}}{\partial \lambda^{(\alpha,\beta;l)}} \right| \odot \left| \Psi^{(l)} \lambda^{(\alpha,\beta;l)} \right|. \quad (\text{A.60})$$

By comparing Eqs. (A.57) and (A.60), we see that changing the basis from \mathcal{B} to \mathcal{F} leads to different transformations of the gradient and the basis coefficient for a non-orthonormal FB \mathcal{F} .

In our experiments in the main body of the work, we computed the SNIP score with $\mathcal{F} = \mathcal{B}$, thus spatial and FB SNIP scores are equivalent. But, if arbitrary FBs are used, the scaling Eq. (A.60) might cause problems and should be taken into account.

F.6. GraSP pruning score

The GraSP score [82] approximates the influence of the removal of a spatial coefficient onto the network's gradient flow before training starts, the so called *importance score*. Using multi-index notation, it is computed as

$$S_{GraSP}(\mathbf{m}; \mathcal{B}) := - \left(H^{\mathcal{B}} \cdot \frac{\partial \mathcal{L}}{\partial \varphi} \right)_{\mathbf{m}} \cdot \varphi_{\mathbf{m}}. \quad (\text{A.61})$$

The corresponding score w.r.t. to basis coefficients \mathcal{F} is given by

$$S_{GraSP}(\mathbf{m}; \mathcal{F}) := - \left(H^{\mathcal{F}} \cdot \frac{\partial \mathcal{L}}{\partial \lambda} \right)_{\mathbf{m}} \cdot \lambda_{\mathbf{m}}. \quad (\text{A.62})$$

By inserting the transformation rules for the coefficient, gradient and Hessian matrix, we derive

$$- \left(H^{\mathcal{B}} \cdot \frac{\partial \mathcal{L}}{\partial \varphi} \right)_{\mathbf{m}} \cdot \varphi_{\mathbf{m}} \quad (\text{A.63})$$

$$= - \left((\Psi^{-1})^T \cdot H^{\mathcal{F}} \cdot \Psi^{-1} \cdot (\Psi^{-1})^T \cdot \frac{\partial \mathcal{L}}{\partial \lambda} \right)_{\mathbf{m}} \cdot (\Psi \cdot \lambda)_{\mathbf{m}}. \quad (\text{A.64})$$

Therefore, the GraSP score is scaled differently for varying layers. Similar to the SNIP score, scaling issues might need to be handled if FBs do not form ONBs.

On the other hand, if all FBs form ONBs, Eq. (A.64) reduces to

$$S_{GraSP}(\mathbf{m}; \mathcal{B}) = - \left(\Psi \cdot H^{\mathcal{F}} \cdot \frac{\partial \mathcal{L}}{\partial \lambda} \right)_{\mathbf{m}} \cdot (\Psi \cdot \lambda)_{\mathbf{m}}. \quad (\text{A.65})$$

Similar to SynFlow, it therefore holds

$$\sum_{\mathbf{m} \in \mathcal{M}_{\alpha,\beta;l}} S_{GraSP}(\mathbf{m}; \mathcal{B}) = \sum_{\mathbf{m} \in \mathcal{M}_{\alpha,\beta;l}} S_{GraSP}(\mathbf{m}; \mathcal{F}) \quad (\text{A.66})$$

Algorithm A1 Standard interspace initialization for a FB 2D convolutional layer

Require: Filter size $K \times K$, number of output channels c_{out} , number of input channels c_{in}

- 1: $\mu_h \leftarrow 0$ mean of spatial coefficients
 - 2: $\sigma_h \leftarrow \sqrt{\frac{2}{c_{in} \cdot K^2}}$ variance of spatial coefficients
 - 3: Initialize $g^{(1)}, \dots, g^{(K^2)}$ via $g^{(n)} = e^{(n)}$ for all $n = 1, \dots, K^2$
 - 4: Initialize $\lambda_n^{(\alpha,\beta)} \sim \mathcal{N}(\mu_h, \sigma_h^2)$ i.i.d. for all $\alpha \in \{1, \dots, c_{out}\}, \beta \in \{1, \dots, c_{in}\}, n \in \{1, \dots, K^2\}$
 - 5: **return** FB $\mathcal{F} = \{g^{(1)}, \dots, g^{(K^2)}\}$, FB coefficients $\lambda = (\lambda_n^{(\alpha,\beta)})_{\alpha,\beta,n}$
-

Algorithm A2 Random ONB interspace initialization for a FB 2D convolutional layer

Require: Filter size $K \times K$, number of output channels c_{out} , number of input channels c_{in}

- 1: $\mu_h \leftarrow 0$ mean of spatial coefficients
 - 2: $\sigma_h \leftarrow \sqrt{\frac{2}{c_{in} \cdot K^2}}$ variance of spatial coefficients
 - 3: Initialize $\tilde{g}^{(1)}, \dots, \tilde{g}^{(K^2)} \in \mathbb{R}^{K \times K}$ with $\tilde{g}_{i,j}^{(n)} \sim \mathcal{N}(0, 1)$ i.i.d. $\triangleright \{\tilde{g}^{(1)}, \dots, \tilde{g}^{(K^2)}\}$ with $\mathbb{P} = 1$ lin. independent
 - 4: Apply Gram-Schmidt on $\{\tilde{g}^{(1)}, \dots, \tilde{g}^{(K^2)}\}$ to obtain ONB $\mathcal{F} = \{g^{(1)}, \dots, g^{(K^2)}\}$
 - 5: Initialize spatial coefficients $\varphi_n^{(\alpha,\beta)} \sim \mathcal{N}(\mu_h, \sigma_h^2)$ i.i.d. for all $\alpha \in \{1, \dots, c_{out}\}, \beta \in \{1, \dots, c_{in}\}, n \in \{1, \dots, K^2\}$
 - 6: Compute basis transformation matrix Ψ according to Eq. (A.19)
 - 7: $\lambda^{(\alpha,\beta)} \leftarrow \Psi^T \cdot \varphi^{(\alpha,\beta)}$ FB coefficients
 - 8: **return** FB $\mathcal{F} = \{g^{(1)}, \dots, g^{(K^2)}\}$, FB coefficients $\lambda = (\lambda_n^{(\alpha,\beta)})_{\alpha,\beta,n}$
-

for $\mathcal{M}_{\alpha,\beta;l} := \{(\alpha, \beta, n; l) : n = 1, \dots, K_{(l)}^2\}$, the multi-indices corresponding to an arbitrary filter $h^{(\alpha,\beta;l)}$. In this case, the total pruning score of a filter $h^{(\alpha,\beta;l)}$ does not depend on the representation.

G. Pruning methods and initialization of the interspace

In our experiments, we used the so called kaiming normal initialization [31] for the standard CNNs. Meaning that $h_{i,j}^{(\alpha,\beta)} \sim \mathcal{N}(\mu_h, \sigma_h^2)$ i.i.d. with

$$\mu_h = 0 \text{ and } \sigma_h = \sqrt{\frac{2}{c_{in} \cdot K^2}}. \quad (\text{A.67})$$

Algorithm A3 Random interspace initialization for a FD 2D convolutional layer with $\#\mathcal{F} = N$ arbitrary

Require: Size of filter dictionary N , filter size $K \times K$, number of output channels c_{out} , number of input channels c_{in}

- 1: $\mu_h \leftarrow 0$ mean of spatial coefficients
 - 2: $\sigma_h \leftarrow \sqrt{\frac{2}{c_{in} \cdot K^2}}$ variance of spatial coefficients
 - 3: Initialize $\tilde{g}^{(1)}, \dots, \tilde{g}^{(N)} \in \mathbb{R}^{K \times K}$ with $\tilde{g}_{i,j}^{(n)} \sim \mathcal{N}(0, 1)$ i.i.d. $\triangleright \{\tilde{g}^{(i_1)}, \dots, \tilde{g}^{(i_m)}\}$ with $i_1 \neq \dots \neq i_m$ and $m \leq K^2$ with $\mathbb{P} = 1$ lin. independent
 - 4: Compute pixelwise sample mean $\tilde{\mu}_{i,j}$ and sample variance $\tilde{\sigma}_{i,j}$ according to Eq. (A.68)
 - 5: $g_{i,j}^{(n)} \leftarrow \sqrt{\frac{1}{N} - \frac{1}{N^2}} \cdot \frac{\tilde{g}_{i,j}^{(n)} - \tilde{\mu}_{i,j}}{\tilde{\sigma}_{i,j}} + \frac{1}{N}$ \triangleright rescale FD
 - 6: Initialize $\lambda_n^{(\alpha, \beta)} \sim \mathcal{N}(\mu_h, \sigma_h^2)$ i.i.d. for all $\alpha \in \{1, \dots, c_{out}\}, \beta \in \{1, \dots, c_{in}\}, n \in \{1, \dots, N\}$
 - 7: **return** FD $\mathcal{F} = \{g^{(1)}, \dots, g^{(N)}\}$, FD coefficients $\lambda = (\lambda_n^{(\alpha, \beta)})_{\alpha, \beta, n}$
-

We initialized all FB-CNNs such that their *spatial representations* follow a `kaiming normal` initialization, see Algs. A1 - A3. For simplicity, we propose the initialization of FB coefficients together with the FB. Of course, if a FB is shared for more than one layer, it has to be initialized just once.

Derivation of rescaling in Algorithm A3. In Alg. A3, \mathcal{F} may contain $N \neq K^2$ elements. Thus, obtaining an equivalent initialization to the spatial `kaiming normal` initialization can not always be obtained by a simple basis transformation. Consequently, we rescale \mathcal{F} in order to mimic a `kaiming normal` initialization of spatial coefficients if FB coefficients are initialized with a `kaiming normal` initialization as well.

Let $\mu_{i,j}$ and $\sigma_{i,j}^2$ be the pixel wise sample mean and sample variance of the FD \mathcal{F} with arbitrary size $N \geq 1$, i.e.

$$\mu_{i,j} := \frac{1}{N} \sum_{n=1}^N g_{i,j}^{(n)} \quad \text{and} \quad \sigma_{i,j}^2 := \frac{1}{N} \sum_{n=1}^N (g_{i,j}^{(n)} - \mu_{i,j})^2. \quad (\text{A.68})$$

By using Eq. (A.68) it holds for an arbitrary i.i.d. initialization of λ with mean μ_λ and variance σ_λ^2

$$\mathbb{E}[h_{i,j}] = \mathbb{E}\left[\sum_{n=1}^N \lambda_n g_{i,j}^{(n)}\right] = \sum_{n=1}^N \mathbb{E}[\lambda_n] g_{i,j}^{(n)} = \mu_\lambda N \mu_{i,j} \quad (\text{A.69})$$

and

$$\mathbb{E}[h_{i,j}^2] = \sum_{n,m} \mathbb{E}[\lambda_n \lambda_m] g_{i,j}^{(n)} g_{i,j}^{(m)} \quad (\text{A.70})$$

$$= \sum_n \mathbb{E}[\lambda_n^2] g_{i,j}^{(n)2} + \sum_n \sum_{n \neq m} \mathbb{E}[\lambda_n^2] g_{i,j}^{(n)} g_{i,j}^{(m)} \quad (\text{A.71})$$

$$= (\sigma_\lambda^2 + \mu_\lambda^2) \sum_n g_{i,j}^{(n)2} + \mu_\lambda^2 \sum_n \sum_{n \neq m} g_{i,j}^{(n)} g_{i,j}^{(m)} \quad (\text{A.72})$$

$$= \sigma_\lambda^2 \sum_n g_{i,j}^{(n)2} + \mu_\lambda^2 \sum_{n,m} g_{i,j}^{(n)} g_{i,j}^{(m)} \quad (\text{A.73})$$

$$= N \sigma_\lambda^2 (\sigma_{i,j}^2 + \mu_{i,j}^2) + \mu_\lambda^2 N^2 \mu_{i,j}^2. \quad (\text{A.74})$$

We now want to determine $\mu_{i,j}$ and $\sigma_{i,j}$ such that $\mathbb{E}[h_{i,j}] = \mu_\lambda$ and $\mathbb{E}[h_{i,j}^2] = \sigma_\lambda^2 + \mu_\lambda^2$ holds, i.e. the distribution of λ and h have the same mean and variance. By Eq. (A.69), setting $\mu_{i,j} = 1/N$ guarantees $\mathbb{E}[h_{i,j}] = \mu_\lambda$. By inserting $\mu_{i,j} = 1/N$ into Eq. (A.74), we see that $\sigma_{i,j}^2 = 1/N - 1/N^2$ implies $\mathbb{E}[h_{i,j}^2] = \sigma_\lambda^2 + \mu_\lambda^2$. Consequently, FDs are rescaled in Alg. A3 to have pixelwise sample mean $\mu_{i,j} = 1/N$ and sample variance $\sigma_{i,j}^2 = 1/N - 1/N^2$.

G.1. General setup for all pruning methods

For SP and the dense baselines, we initialized the standard CNN with the `kaiming normal` initialization. For IP, we initialize networks according to Alg. A1 for all experiments in the main body of the text. In Sec. B, also Algs. A2 and A3 are used as initializations for the interspace.

Pruning masks are computed according to the formulas described in Sec. F for SP and IP. In the following we will describe the used pruning methods in detail.

G.2. Lottery tickets with resetting coefficients

LTs with resetting coefficients to an early training iteration [19] are obtained as follows. We first train the network to epoch $t_0 = 500$ and store the corresponding model, optimizer, etc. Afterwards, the network is trained to convergence. Then, 20% of the coefficients are pruned, based on their magnitudes. All non-zero parameters are reset to their values at training time $t_0 = 500$ and pruned ones are fixed at zero from now on. Thus, contrarily to DST, a pruned coefficient will never be able to recover. The training schedule parameters, like learning rate or moving averages for batch normalization and SGD with momentum, are reset to their corresponding value at t_0 as well. For IP, also the FBs are reset to step t_0 . Then, the network is again trained to convergence, 20% of the non-zero coefficients are pruned and the remaining non-zero parameters are reset again. This is done, until the desired pruning rate is reached. Note, for the last pruning step also $< 20\%$ of the non-zero parameters

might be pruned to exactly match the desired pruning rate. If the final pruning rates is reached, the network with desired sparsity is trained for a last, final time. All in all,

$$k = 1 + \left\lceil \frac{\log(1-p)}{\log 4 - \log 5} \right\rceil \quad (\text{A.75})$$

trainings are needed to obtain and train a network with sparsity p using this iterative approach.

Following [19], we do not prune the fully connected layer for LTs but keep it dense.

G.3. Dynamic sparse training

Dynamic sparse training methods adapt the network’s pruning mask during training [5, 13, 17, 45, 55, 56]. In this work we use SET [55] which is based on estimating the importance of coefficients via magnitude pruning and re-growing coefficients due to a random selection. RigL [17] improves this approach by regrowing coefficients which have the biggest gradient magnitudes.

Before training, the networks are pruned randomly. It was shown in [17] that using layerwise sparsity corresponding to an Erdős-Rényi-kernel leads to good results. This means that each layer has sparsity depending on its size, *i.e.*

$$1 - \varepsilon \cdot \frac{c_{out}^{(l)} + c_{in}^{(l)} + 2 \cdot K^{(l)}}{c_{out}^{(l)} \cdot c_{in}^{(l)} \cdot K^{(l)}}, \quad (\text{A.76})$$

and ε is a global parameter, tuned such that a global sparsity of p is obtained.

During training, the pruning mask is frequently updated. For this, parameters are pruned for each layer with rate p_t . To be precise, all *non-zero* parameters in a layer are pruned with the rate p_t . The pruning rate p_t depends on the training step t and decays with a cosine schedule in order to improve convergence [17]. It holds

$$p_t = p_{min} + \frac{1}{2} (p_{init} - p_{min}) \cdot \left(1 + \cos \left(\frac{t\pi}{T} \right) \right) \quad (\text{A.77})$$

where T is the number of total training steps, $p_{min} = 0.005$ the minimal pruning rate and $p_{init} = 0.5$ the initial rate used for updating the pruning mask. Of course, in each layer an equal number of non trained coefficients are regrown after pruning. Following [55] and [17], regrown coefficients are initialized with value 0 but are updated via SGD from this moment on.

SET and RigL work optimal for different pruning mask update frequencies [45]. According to [45], we update pruning masks each 1, 500 training steps for SET and each 4, 000 steps for RigL.

G.4. Pruning at initialization

We test PaI methods, SNIP [41], GraSP [82] and SynFlow [76] together with random PaI. In contrast to LTs and

DST, PaI is quite simple. For SNIP and GraSP we compute the pruning scores described in Secs. F.5 and F.6 with the help of 100 batches of training data for the CIFAR-10 experiments and 15 for ImageNet. As proposed by [76], we compute the pruning scores for SNIP and GraSP with all batch normalization layers [34] set to PyTorch’s `train` mode. Afterwards, all coefficients are pruned one-shot.

The GraSP scores for SP, Eq. (A.61), and IP, Eq. (A.62), require the computation of a Hessian vector product $H \cdot g$. Fortunately, not the whole Hessian H needs to be computed to evaluate such products. For this, we use the linearity of the derivative. For $H \in \mathbb{R}^{d \times d}$ and an arbitrary vector $v \in \mathbb{R}^d$, it holds

$$(H \cdot v)_i = \sum_j \frac{\partial^2 \mathcal{L}}{\partial \lambda_i \partial \lambda_j} \cdot v_j \quad (\text{A.78})$$

$$= \frac{\partial}{\partial \lambda_i} \left(\sum_j \frac{\partial \mathcal{L}}{\partial \lambda_j} v_j \right) \quad (\text{A.79})$$

$$= \frac{\partial}{\partial \lambda_i} \left\langle \frac{\partial \mathcal{L}}{\partial \lambda}, v \right\rangle, \quad (\text{A.80})$$

and consequently $H \cdot v = \frac{\partial \langle \frac{\partial \mathcal{L}}{\partial \lambda}, v \rangle}{\partial \lambda}$. By using v as the fixed gradient $g_v = \frac{\partial \mathcal{L}}{\partial \lambda}$, we can compute $H \cdot g_v$ with only three backward passes. The first one is needed to compute the fixed gradient g_v . The second and third are required to compute $H \cdot g_v = \frac{\partial \langle \frac{\partial \mathcal{L}}{\partial \lambda}, g_v \rangle}{\partial \lambda}$. An implementation of this can be found in the official code base for GraSP, [see this link](#) (MIT license).

In contrast to these one-shot methods, pruning masks for SynFlow are computed in an iterative fashion. For this purpose, we compute the pruning score proposed in Sec. F.4 with one forward- and backward pass, prune a small fraction of elements and repeat it 100 times. The pruning rate grows with an exponential schedule [76] which gives the pruning rate

$$p_k = 1 - (1 - p)^{k/100} \quad (\text{A.81})$$

after the k^{th} pruning iteration. For computing the SynFlow score, we set all batch normalization layers to `eval` mode, as suggested by [76].

G.5. Gradual magnitude pruning

Following [21] we gradually sparsify the model based on the coefficients’ magnitudes. After each N training iterations, the pruning rate is increased and new weights are pruned. The pruning rate at iteration $k \cdot N$ is given by

$$p(k \cdot N) = \begin{cases} 0, & k \cdot N < t_0 \\ p \cdot \left(1 - \left(1 - \frac{k \cdot N - t_0}{t_1 - t_0} \right)^3 \right), & k \cdot N \in [t_0, t_1] \\ 1, & k \cdot N > t_1 \end{cases} \quad (\text{A.82})$$

Here, p denotes the final pruning rate and t_0, t_1 are the training iterations where the gradual pruning begins and ends, respectively. For each iteration $k \cdot N$ with $k \in \mathbb{N}$, the $p(k \cdot N) \cdot d$ weights with smallest magnitude are pruned. Since pruned weights are frozen at 0, those weights will be pruned again and they will therefore never recover.

We follow the suggestions in [21] for choosing t_0, t_1 and N for the final pruning rates $p \in \{0.8, 0.9\}$ for the ResNet50 on ImageNet (summarized in [this table](#)). Since we use a different batch size than [21] in our experiments (256 compared to 1,024), we adapt their choices for t_0, t_1 and N to our batch size by multiplying them by 4. However, these choices are not optimized for batch size 256 and we therefore report slightly worse results than [21]. For our experiments, we use $N = 8,000, t_0 = 160,000$ and $t_1 = 400,000$ for $p = 0.8$. For $p = 0.9$ we set $N = 8,000, t_0 = 160,000$ and $t_1 = 304,000$.

G.6. Fine tuning

We also compare IP and SP for magnitude pruning applied on a pre-trained, dense network while the sparse network is fine-tuned afterwards. As suggested by [68], we do not use a classical fine-tuning setup beginning with a small learning rate, but use the setup of the dense pre-training also for fine-tuning. Thus, the sparse fine-tuning starts with a high learning rate. Naturally, for IP we use the pre-trained FBs as starting point for the fine-tuning. For FT, we use the same training setup as for RigL.

H. Experimental setup

In this Section, we describe the setup for the experiments discussed in the main body of the paper and in Sec. B. Our experiments were conducted on an internal cluster with CentOS Linux release 7.9.2009 (Core). We used Python 3.8 with the deep learning framework PyTorch1.9 [61] (BSD license) together with cudatoolkit 10.2 [58]. As hardware we had an Intel XEON E5-2680 v4 à 2.4 GHz CPU and n NVIDIA GeForce 1080ti GPUs, where $n = 1$ for the CIFAR-10 experiments, $n \in \{2, 4\}$ for ResNet50 on ImageNet and $n = 8$ for ResNet18 on ImageNet. Further details on the training time for one epoch, the number of used CPU cores, used RAM and GPU memory are given in Tab. A5.

Since we compare and adapt different pruning methods, we used the publicly available codes for reproducing the results and modifying them for the IP setting. These are:

- [Link to code](#) for SynFlow [76] (unknown license),
- [Link to code](#) for GraSP [82] (MIT license),
- PyTorch adaption of [link to code](#) for SNIP [41] (MIT license),
- [Link to code](#) as the official PyTorch version for Lottery Tickets [19] (MIT license),
- [Link to code](#) as a base for DST methods which is the official code for [45] (unknown license). In [45], train-

ing schedules for SET [55] and RigL [17] are improved.

This code base is also used for FT [68] and GMP [21].

For the IP versions of these pruning methods, we additionally updated the code to use Alg. 1 as a 2D convolution.

The initializations of the used CNNs and FB-CNNs are described in Sec. G. Used hyperparameters are summarized in Tab. A6. We also trained the dense standard networks with the same training schedules as the pruned ones. Results for the dense models can be found in Tab. A5.

Training schedule. As common in the literature of sparse training, see for example [41, 76, 82], no hyperparameter tuning is done in this work. We want to highlight, that all used training setups are chosen from one of the adapted pruning methods. Especially, FBs and FB coefficients are trained with the standard learning rates, optimized for training spatial coefficients.

For PaI on CIFAR-10, we used the setup from [41], whereas SET uses the training schedule from [45]. The CIFAR-10 experiment for LTs is equal to the one used in [19]. For ImageNet on ResNet18, we used the standard PyTorch ImageNet training ([see this link](#)) (BSD 3-clause license) as baseline for our training – the same as used in [82]. Results can be further improved by adding learning rate warm-up for 5 epochs [25] and label smoothing [75] with smoothing parameter $\varepsilon = 0.1$. This improvement is inherited from [45] and applied to train RigL on ResNet50. For FT, we use the same training hyperparameters as for RigL whereas we use the suggested one from the original paper [21] for GMP.

CIFAR-10. CIFAR-10 [38] consists of 60,000 32×32 RGB images. CIFAR-10 is a publicly available dataset with, best to our knowledge, no existing licenses. Authors are allowed to use the datasets for publications if the tech report [38] is referenced. CIFAR-10 has 10 classes with 6,000 images per class. The data is split into 50,000 training and 10,000 test images. For each training, we randomly split the training images into two parts, 45,000 images for training and 5,000 images for validation. All images for training, validation and testing are normalized by their channel wise mean and standard deviation. Furthermore, we additionally use the standard data augmentation for CIFAR on the training images. This is given by cropping and random horizontal flipping of the images. Test results are reported for the early stopping epoch, the epoch with the highest validation accuracy. Used hyperparameters are summarized in Tab. A6.

ImageNet. The ImageNet ILSVRC2012 [69] dataset is an image classification dataset, containing approximately 1.2 million RGB images for training and 150,000 RGB images for validation, divided into 1,000 classes. ImageNet has a custom license allowing non-commercial research. To be allowed to use ImageNet for non-commercial research, access to the image database has to be requested – which we did. Full terms for the usage of the ImageNet database can be found in [this link](#).

Dataset/Model	Mean \pm Std	# Params	# FB Params fine sharing	Time 1 epoch [s]	#GPUs	GPU memory	#CPU cores	RAM
CIFAR-10/VGG16	93.41 \pm 0.07%	15.3mio	1,053	23.7	1	11 GB	1	12 GB
CIFAR-10/VGG16-LT	93.58 \pm 0.10%	14.7mio	1,053	21.9	1	11 GB	1	12 GB
ImageNet/ResNet18	69.77 \pm 0.06%	11.7mio	1,296	2,770.5	8	8 \times 11 GB	8	8 \times 12 GB
ImageNet/ResNet50 (RigL & FT)	77.15 \pm 0.04%	25.6mio	3,697	4,354.2	2	2 \times 11 GB	4	4 \times 12 GB
ImageNet/ResNet50 (GMP)	76.64 \pm 0.06%	25.6mio	3,697	2,622.8	4	4 \times 11 GB	8	8 \times 12 GB

Table A5. Top-1 test accuracies for densely trained models with additional information, including the hardware setup. Note, # FB parameters for fine sharing is the *most* number of extra parameters induced by the interspace representation. We suggest to use at least 2 times the number of CPU cores than GPUs, otherwise data loading becomes a bottleneck (compare ResNet18 and ResNet50 for GMP which have approximately the same runtime despite different network sizes and # GPUs).

Experiment	CIFAR-10 (no LTs)	CIFAR-10 (LTs)	ImageNet (PaI)	ImageNet (RigL & FT)	ImageNet (GMP)
Network	VGG16	VGG16-LT	ResNet18	ResNet50	ResNet50
# Trainings	5	5	3	3	3
# Epochs	250	160	90	100	100
Batch Size	128	128	512	128	256
# GPUs	1	1	8	2	4
Optimizer: SGD- Momentum	Momentum	Momentum	Momentum	Momentum	Momentum
	0.9	0.9	0.9	0.9	0.9
Learning Rate	0.1	0.1	0.1	0.1	0.1
LR Decay	$\times 0.1$	$\times 0.1$	$\times 0.1$	$\times 0.1$	$\times 0.1$
	every 30k iterations	epochs 80/120	epochs 30/60	epochs 30/60/90	epochs 30/60/80
LR Warm-up	\times	\times	\times	5 epochs, linear	5 epochs, linear
Weight Decay	$5 \cdot 10^{-4}$	10^{-4}	10^{-4}	10^{-4}	10^{-4}
Label smoothing	\times	\times	\times	$\varepsilon = 0.1$	$\varepsilon = 0.1$

Table A6. Setups for experiments in the main body of the paper and Sec. B. Each run of the # trainings was executed with a different random seed.

Again, all images are normalized for each channel. Training images are randomly cropped to size 224×224 and randomly flipped in the horizontal direction. The validation images are resized to size 256×256 and their central 224×224 pixels are used for validation. For the ResNet50 experiment, we also add label smoothing on the training loss with a smoothing factor 0.1. As common in the literature, we report results on the validation set, since labels for the test set are publicly not available. Hyperparameters are provided in Tab. A6.

I. Network architectures

In this Section we provide the used network architectures VGG16 [72] and the adapted version VGG16-LT for CIFAR-10, as well as ResNet18 and ResNet50 [32] for ImageNet.

The architectures are shown in Tabs. A7 to A9. Note, we use two different versions of a VGG16, a small one for the LT experiments and a bigger one for the remaining experiments. A graphical description of the residual block and the bottleneck block used for ResNets, is shown in Figs. A7 and A8, respectively. Furthermore, all architectures are used in their standard form or with FB convolutions. Therefore,

(FB) always indicates, that a FB version of the filter is used for the FB-CNN. Not all convolutional layers are marked with a (FB), since we only apply the FB formulation on convolutional layers with kernel size $K > 1$. Additionally, we indicate the layers which share one FB for coarse, medium and fine FB sharing in Tabs. A7 to A9.

All tensor dimensions of standard 2D convolutional filters are given as $c_{out} \times c_{in} \times K \times K$, where c_{in} equals the number of input channels, c_{out} the number of output channels and $K \times K$ the size of each convolutional kernel. FB convolutions have coefficients represented by a tensor of size $c_{out} \times c_{in} \times K^2$. For pooling layers, $K \times K$ denotes the tiling size. In the case of linear layers, the tensor size is given as $c_{out} \times c_{in}$, where c_{in} is the number of incoming neurons and c_{out} the number of outgoing neurons. For training the networks, we used the cross entropy loss function.

J. Proof of Theorem 1

SDL optimizes a dictionary $\mathbf{F} \in \mathbb{R}^{m \times M}$ jointly with its coefficients $R \in \mathbb{R}^{M \times n}$ w.r.t. the non-convex problem

$$\inf_{\mathbf{F}, R} \|U - \mathbf{F} \cdot R\|_F \text{ s.t. } \|R\|_0 \leq s, \quad (\text{A.83})$$

Module	Output Size	c_{in}	c_{out}	K	Repeat	Stride	Padding	Bias	BatchNorm	ReLU	Coarse	Medium	Fine
(FB) Conv2D	32×32	3	64	3×3	$\times 1$	1×1	1×1	✓	✓	✓			■
(FB) Conv2D	32×32	64	64	3×3	$\times 1$	1×1	1×1	✓	✓	✓			
MaxPool2D	16×16	64	64	2×2	$\times 1$	2×2	0×0	×	×	×			■
(FB) Conv2D	16×16	64	128	3×3	$\times 1$	1×1	1×1	✓	✓	✓			
(FB) Conv2D	16×16	128	128	3×3	$\times 1$	1×1	1×1	✓	✓	✓			■
MaxPool2D	8×8	128	128	2×2	$\times 1$	2×2	0×0	×	×	×			
(FB) Conv2D	8×8	128	256	3×3	$\times 1$	1×1	1×1	✓	✓	✓			■
(FB) Conv2D	8×8	256	256	3×3	$\times 2$	1×1	1×1	✓	✓	✓			
MaxPool2D	4×4	256	256	2×2	$\times 1$	2×2	0×0	×	×	×			■ $\times 2$
(FB) Conv2D	4×4	256	512	3×3	$\times 1$	1×1	1×1	✓	✓	✓			
(FB) Conv2D	4×4	512	512	3×3	$\times 2$	1×1	1×1	✓	✓	✓			■ $\times 2$
MaxPool2D	2×2	512	512	2×2	$\times 1$	2×2	0×0	×	×	×			
(FB) Conv2D	2×2	512	512	3×3	$\times 3$	1×1	1×1	✓	✓	✓			■ $\times 3$
MaxPool2D	1×1	512	512	2×2	$\times 1$	2×2	0×0	×	×	×			
Linear	512	512	512	—	$\times 2$	—	—	✓	✓	✓	Removed for VGG16-LT		
Linear	10	512	10	—	$\times 1$	—	—	✓	×	×			

Table A7. VGG16 and VGG16-LT for CIFAR-10 with *coarse*, *medium* and *fine* FB sharing schemes specified in the last three columns. Each ■, connected by either another ■ or |, corresponds to exactly one FB \mathcal{F} shared for all filters in the corresponding layers. The thin connection | corresponds to `MaxPool2D` layers which do not use the FBs themselves. Note, for VGG16-LT, the first and second `Linear` layers are removed.

Module	Output Size	c_{in}	c_{out}	K	Stride	Padding	Bias	BatchNorm	ReLU	Coarse	Medium	Fine
Conv2D	112×112	3	64	7×7	2×2	3×3	×	✓	✓			■ $\times 4$
MaxPool2D	56×56	64	64	3×3	2×2	1×1	×	×	×			
(FB) ResBlock $\times 2$	56×56	64	64	3×3	1×1	1×1	×	✓	✓			
(FB) ResBlock $\times 2$	28×28	64	128	3×3	2×2	1×1	×	✓	✓			
(FB) ResBlock $\times 2$	14×14	128	256	3×3	2×2	1×1	×	✓	✓			
(FB) ResBlock $\times 2$	7×7	256	512	3×3	2×2	1×1	×	✓	✓			■ $\times 4$
AvgPool2D	1×1	512	512	7×7	0×0	0×0	×	×	×			
Linear	1,000	512	1,000	—	—	—	✓	×	×			

Table A8. ResNet18 for ImageNet with (FB) `ResBlock $\times 2$` , shown in Fig. A7. Last three columns declare *coarse*, *medium* and *fine* FB sharing schemes. Each ■ connected by another ■ corresponds to exactly one FB \mathcal{F} shared for all filters in the corresponding layers. Note, the first convolutional layer has kernel size 7×7 and we do not use a FB formulation for this layer.

for a target $U \in \mathbb{R}^{m \times n}$ and sparsity constraint s . In our context U corresponds to a convolutional layer, the dictionary \mathbf{F} to the layer’s FB (FD) \mathcal{F} and R to the FB (FD) coefficients. Standard magnitude pruning can be seen as a special case of SDL where the dictionary \mathbf{F} is fixed to form the standard basis, i.e. $\mathbf{F} = \text{id}_{\mathbb{R}^m}$. Accordingly,

$$\inf_{\Phi} \|U - \Phi\|_F \text{ s.t. } \|R\|_0 \leq s \quad (\text{A.84})$$

is minimized.

Theorem A.1. *Let $1 < m \leq M$, $0 < s < m \cdot n$ and $U_{i,j} \sim \mathcal{N}(0, 1)$ i.i.d. Further assume that $\varepsilon_{(1)}$ is the infimum of Eq. (A.83) and $\varepsilon_{(2)}$ the minimum of Eq. (A.84). Assume Φ^* to be the minimizer for Eq. (A.84). Then $\varepsilon_{(1)} < \varepsilon_{(2)}$ holds with probability $\mathbb{P} = 1$.*

If furthermore $\text{supp } R$ for Eq. (A.83) is fixed to be equal to $\text{supp } \Phi^$, then $\varepsilon_{(1)} \leq \varepsilon_{(2)}$ and strict inequality holds with*

$\mathbb{P} \geq 1 - \delta$, where

$$\delta = \begin{cases} 0 & , \text{ if } s \not\equiv 0(\text{mod } m) \\ \frac{\binom{n}{s}}{\binom{m \cdot n}{s}} & , \text{ if } s \equiv 0(\text{mod } m) \end{cases}. \quad (\text{A.85})$$

Figure A9 shows the probability of the solution to Eq. (A.84) being strictly bigger than the solution of Eq. (A.83) if $\text{supp } R$ is restricted to be $\text{supp } \Phi^*$. Precisely, it shows $1 - \delta$ for varying pruning rates p . It can be seen that, except for the trivial case of a network being completely pruned or being not pruned at all, δ is numerically equal to zero, even for a network with only 100 filters.^{A.4} Thus,

^{A.4}The minimum δ we computed for non-trivial pruning rates $p \in (0, 1)$ for a network with $n = 10$ filters of size 3×3 was given by $\delta \leq 10^{-10}$. For $n \geq 100$, numerically $\delta = 0$ for all $p \in (0, 1)$.

Module	Output Size	c_{in}	c_{mid}	c_{out}	n_b	K	Stride	Padding	Bias	BN	ReLU	Coarse adapt	Medium	Fine
(FB) Conv2D	112×112	3	—	64	—	7×7	2×2	3×3	×	✓	✓	■	■	■ $\times 1$
MaxPool2D	56×56	64	—	64	—	3×3	2×2	1×1	×	×	×			
(FB) BottleneckBlock	56×56	64	64	256	3	3×3	1×1	1×1	×	✓	✓	■	■	■ $\times 3$
(FB) BottleneckBlock	28×28	256	128	512	4	3×3	2×2	1×1	×	✓	✓	■	■	■ $\times 4$
(FB) BottleneckBlock	14×14	512	256	1,024	6	3×3	2×2	1×1	×	✓	✓	■	■	■ $\times 6$
(FB) BottleneckBlock	7×7	1,024	512	2,048	3	3×3	2×2	1×1	×	✓	✓	■	■	■ $\times 3$
AvgPool2D	1×1	512	—	512	—	7×7	0×0	0×0	×	×	×			
Linear	1,000	512	—	1,000	—	—	—	—	✓	×	×			

Table A9. ResNet50 for ImageNet with (FB) BottleneckBlock, shown in Fig. A8. Last three columns declare the adapted coarse, medium and fine FB sharing schemes. Each ■ connected by another ■ corresponds to exactly one FB \mathcal{F} shared for all filters in the corresponding layers.

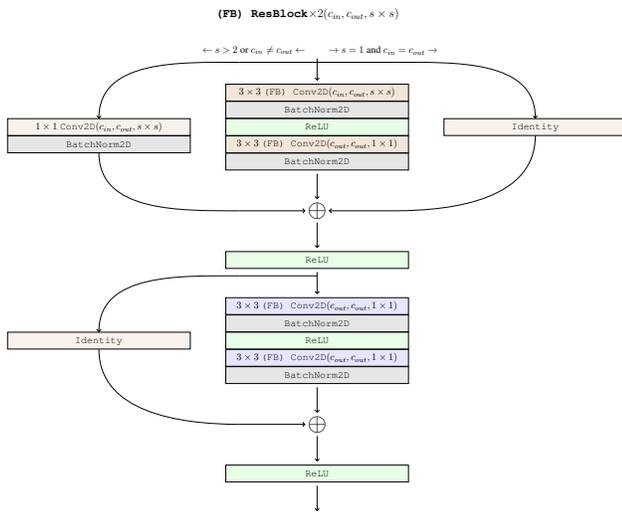


Figure A7. Architecture of a (FB) ResBlock $\times 2$ with c_{in} input channels, c_{out} output channels and stride $s \times s$ for the first (FB) convolution. The first residual connection is a standard 1×1 2D Convolution followed by a BatchNorm2D layer if $s > 1$ or $c_{in} \neq c_{out}$. However, if $c_{in} = c_{out}$ and $s = 1$, the first residual connection is simply an identity mapping. The second residual connection is always an identity mapping. All (FB) Conv2D layers do not have biases.

despite $\delta > 0$, numerically the chance of $\varepsilon_{(1)} = \varepsilon_{(2)}$ is equal to zero.

The proof of Thm. A.1 is split in several parts.

- Lemma A.2 shows that Eq. (A.84) always has a minimum and constructs the minimizing Φ^* .
- Lemma A.3 will show that for each feasible point Φ_0 for Eq. (A.84), there exists an equivalent feasible point (\mathbf{F}_0, R_0) for Eq. (A.83) with the same sparsity $\|R_0\|_0 = \|\Phi_0\|_0$ and distance $\|U - \mathbf{F}_0 \cdot R_0\|_F = \|U - \Phi_0\|_F$.
- Consequently, Corollary A.4 shows that the solution to Eq. (A.83) is always smaller or equal to the solution of

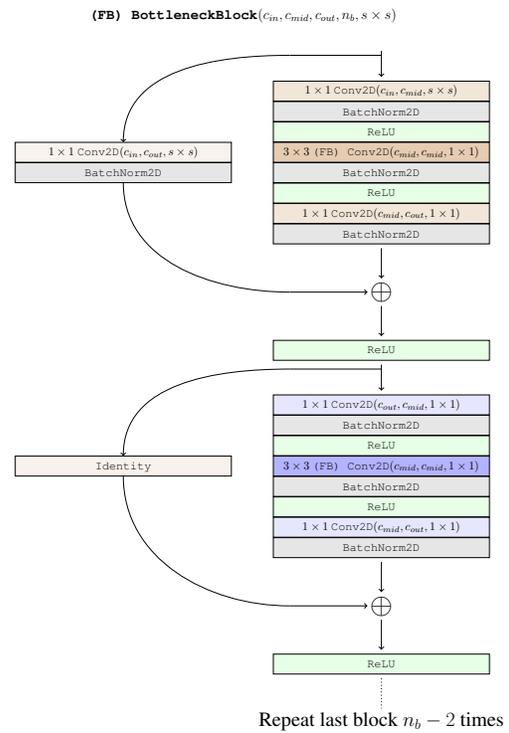


Figure A8. Architecture of a (FB) BottleneckBlock with c_{in} input channels, c_{mid} middle channels and c_{out} output channels and stride $s \times s$ for the first 1×1 (FB) convolution. The number of small blocks forming the (FB) BottleneckBlock is given by n_b . The first residual connection is a standard 1×1 2D Convolution followed by a BatchNorm2D layer. The following residual connections are always identity mappings. All (FB) Conv2D layers do not have biases.

Eq. (A.84).

- The first part of the proof of Thm. A.1 shows that the solution obtained by Eq. (A.84) can only with a small chance $\delta_0 \leq \delta$ be the optimum of Eq. (A.83). This is based on two facts, first we construct the equivalent

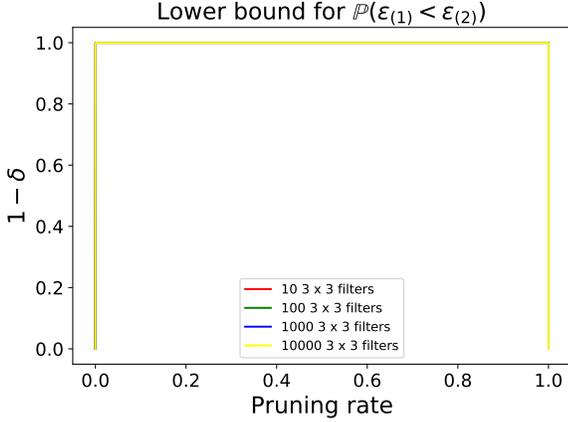


Figure A9. Lower bound $1 - \delta$ for $\mathbb{P}(\varepsilon_{(1)} < \varepsilon_{(2)})$, i.e. the optimization problem (A.83) having a *smaller* solution than the problem (A.84). Computed for varying numbers n of 3×3 filters in a layer and varying pruning rates $p \in [0, 1]$.

point (\mathbf{F}_0, R_0) to Φ^* according to Lemma A.3. Then, we show that with a probability of at most δ , (\mathbf{F}_0, R_0) fulfills a necessary condition for solving Eq. (A.83). This condition is given by \mathbf{F}_0 yielding a local optimum^{A.5} of the smooth, convex function

$$f : \mathbb{R}^{m \times M} \rightarrow \mathbb{R}, \mathbf{F} \mapsto \|U - \mathbf{F} \cdot R_0\|_F^2, \quad (\text{A.86})$$

which is evaluated by looking at the probability of \mathbf{F}_0 being a root of $\frac{\partial f}{\partial \mathbf{F}}$.

- The second part of Thm. A.1 adapts (\mathbf{F}_0, R_0) if $\text{supp } R_0$ is not fixed to be equal to $\text{supp } \Phi^*$. By setting one column of Φ^* as a new basis element, the number of coefficients needed to match $\Phi^* = \mathbf{F}^* \cdot R^*$ with the adapted (\mathbf{F}^*, R^*) is reduced. This of course provides new unused coefficients which are used to *better approximate* the target U .

Lemma A.2. *The optimization problem (A.84) always has a solution $\Phi^* \in \mathbb{R}^{m \times n}$ obtained by*

$$\Phi^* = (\Phi_{i,j}^*)_{i,j} \text{ with} \quad (\text{A.87})$$

$$\Phi_{i,j}^* = \begin{cases} U_{i,j}, & \text{if } (i,j) \in \text{TOP}_s(U) \\ 0, & \text{else} \end{cases}$$

Here,

$$\text{TOP}_s(U) := \{(i_0, j_0) \in \{1, \dots, m\} \times \{1, \dots, n\} : U_{i_0, j_0} \text{ belongs to the top } s \text{ magnitudes of } \{U_{i,j}\}\} \quad (\text{A.88})$$

defines the indices corresponding to the s highest magnitudes of U .

^{A.5}By convexity of f it is therefore a global minimum.

Proof of Lemma A.2. To solve Eq. (A.84), we rewrite the optimization problem into its equivalent, squared form

$$\inf_{\Phi \in \mathbb{R}^{m \times n}} \|U - \Phi\|_F^2 \text{ s.t. } \|\Phi\|_0 \leq s. \quad (\text{A.89})$$

The problem (A.89) is equivalent to

$$\inf_{\Phi \in \bigcup_{k=1}^r \mathcal{S}_k} \|U - \Phi\|_F^2 = \min_{k \in \{1, \dots, r\}} \inf_{\Phi \in \mathcal{S}_k} \|U - \Phi\|_F^2 \quad (\text{A.90})$$

with $r = \binom{n+m}{s}$,

$$\mathcal{S}_k = \{A \in \mathbb{R}^{m \times n} : \text{supp } A \subset S_k\}, \quad (\text{A.91})$$

$$S_k \subset \{1, \dots, m\} \times \{1, \dots, n\}, \#S_k = s$$

satisfying $S_k \neq S_j$ for $k \neq j$ and $\bigcup_{k=1}^r \mathcal{S}_k = \{A \in \mathbb{R}^{m \times n} : \|A\|_0 \leq s\}$.

In order to solve Eq. (A.90), we minimize for each k

$$\inf_{\Phi \in \mathcal{S}_k} \|U - \Phi\|_F^2 = \inf_{\Phi \in \mathcal{S}_k} \sum_{i,j} (U_{i,j} - \Phi_{i,j})^2 \quad (\text{A.92})$$

individually. The problem (A.92) is minimized by $\Phi^{(k)*}$ with

$$\Phi_{i,j}^{(k)*} = \begin{cases} U_{i,j}, & \text{if } (i,j) \in S_k \\ 0, & \text{else} \end{cases}. \quad (\text{A.93})$$

Thus, the minimum of Eq. (A.92) for a $k \in \{1, \dots, r\}$ is given by

$$\|U - \Phi^{(k)*}\|_F^2 = \sum_{(i,j) \notin S_k} U_{i,j}^2 \quad (\text{A.94})$$

Equation (A.94) leads to the solution of Eq. (A.90), given by

$$\min_{k \in \{1, \dots, r\}} \sum_{(i,j) \notin S_k} U_{i,j}^2 = \max_{k \in \{1, \dots, r\}} \sum_{(i,j) \in S_k} U_{i,j}^2 \quad (\text{A.95})$$

which is reached by choosing k such that $S_k = \text{TOP}_s(U)$. \square

Lemma A.3. *Let $m \leq M$, then for each $\Phi \in \mathbb{R}^{m \times n}$ there exists a $\mathbf{F} \in \mathbb{R}^{m \times M}$ and a $R \in \mathbb{R}^{M \times n}$ with $\|R\|_0 = \|\Phi\|_0$ and $\mathbf{F} \cdot R = \Phi$.*

Proof of Lemma A.3. Let $\Phi = (\Phi_{i,j})_{i,j} \in \mathbb{R}^{m \times n}$ be given. Now, we define $R \in \mathbb{R}^{M \times n}$ and $\mathbf{F} \in \mathbb{R}^{m \times M}$ via

$$R_{i,j} = \begin{cases} \Phi_{i,j}, & \text{if } i \leq m \\ 0, & \text{else} \end{cases} \quad (\text{A.96})$$

and

$$\mathbf{F}_{i,j} = \begin{cases} 1, & \text{if } i = j \text{ and } j \leq m \\ 0, & \text{else} \end{cases}. \quad (\text{A.97})$$

By construction of R and \mathbf{F} , it holds $\Phi = \mathbf{F} \cdot R$ and $\|R\|_0 = \|\Phi\|_0$. \square

Corollary A.4. Let $m \leq M$, $\varepsilon_{(1)}$ be the infimum of Eq. (A.83) and $\varepsilon_{(2)}$ be the minimum of Eq. (A.84), respectively. Then it holds $\varepsilon_{(1)} \leq \varepsilon_{(2)}$.

Proof of Corollary A.4. By Lemma A.2, Eq. (A.84) is always minimized by a $\Phi^* \in \mathbb{R}^{m \times n}$. By using Lemma A.3, since $m \leq M$, there exists $\mathbf{F}_0 \in \mathbb{R}^{m \times M}$ and $R_0 \in \mathbb{R}^{M \times n}$ with $\Phi^* = \mathbf{F}_0 \cdot R_0$ and $\|R_0\|_0 = \|\Phi^*\|_0 = s$. Thus, (\mathbf{F}_0, R_0) is feasible for Eq. (A.83) and consequently $\varepsilon_{(1)} \leq \varepsilon_{(2)}$. \square

Proof of Theorem A.1. First part of proof with $\text{supp } R = \text{supp } \Phi^*$. W.l.o.g. we assume $\varepsilon_{(1)} = \varepsilon_{(2)}$. By Corollary A.4, $\varepsilon_{(1)} \leq \varepsilon_{(2)}$ always holds. If $\varepsilon_{(1)} < \varepsilon_{(2)}$, we would be finished with the proof. Furthermore, we assume w.l.o.g. $m = M$, since otherwise we just fill the corresponding entries in $R_{i,j}$ and $\mathbf{F}_{k,i}$ for index values $m < i \leq M$ and arbitrary j, k with zeros.

Therefore, let $\varepsilon_{(1)} = \varepsilon_{(2)}$. Let $\Phi^* \in \mathbb{R}^{m \times n}$ solve Eq. (A.84). Then, there exists an equivalent feasible point $(\mathbf{F}_0, R_0) \in \mathbb{R}^{m \times m} \times \mathbb{R}^{m \times n}$ for Eq. (A.83), constructed according to Eqs. (A.96) and (A.97) in the proof of Lemma A.3. I.e., $\Phi^* = \mathbf{F}_0 \cdot R_0$ and $\|\Phi^*\|_0 = \|R_0\|_0$. By the assumption $\varepsilon_{(1)} = \varepsilon_{(2)}$, (\mathbf{F}_0, R_0) also solves Eq. (A.83). Especially, \mathbf{F}_0 defines a *global* minimum of the smooth, convex function

$$f : \mathbb{R}^{m \times m} \rightarrow \mathbb{R}, \mathbf{F} \mapsto \|U - \mathbf{F} \cdot R_0\|_F^2. \quad (\text{A.98})$$

Note, minimizing f is, contrarily to Eq. (A.83), a convex problem.

A necessary, and by convexity of f even sufficient, condition for \mathbf{F}_0 to minimize f is given by

$$\left. \frac{\partial f}{\partial \mathbf{F}} \right|_{\mathbf{F}=\mathbf{F}_0} = 0 \in \mathbb{R}^{m \times m}. \quad (\text{A.99})$$

It holds

$$\frac{\partial f}{\partial \mathbf{F}} = \frac{\partial}{\partial \mathbf{F}} \|U - \mathbf{F} \cdot R_0\|_F^2 = 2 \cdot (\mathbf{F} \cdot R_0 - U) \cdot R_0^T. \quad (\text{A.100})$$

By combining Eqs. (A.99) and (A.100), we get a necessary condition for \mathbf{F}_0 yielding a minimum for f , given by

$$(U - \mathbf{F}_0 \cdot R_0) \cdot R_0^T = 0. \quad (\text{A.101})$$

Consequently, Eq. (A.101) is a necessary condition for (\mathbf{F}_0, R_0) to define the minimum for Eq. (A.83). From the construction of \mathbf{F}_0 and R_0 we know that $\mathbf{F}_0 \cdot R_0 = \Phi^*$, $\mathbf{F}_0 = \text{id}_{\mathbb{R}^m}$ and $R_0 = \Phi^*$. By Lemma A.2, Φ^* is given by $\Phi_{i,j}^* = \chi_{\{(i,j) \in \text{TOP}_s(U)\}} \cdot U_{i,j}$ with the *characteristic function* $\chi_{\{\cdot\}}$. Combining this with Eq. (A.101) leads to the necessary condition

$$\hat{U} \cdot \check{U}^T = 0 \quad (\text{A.102})$$

with

$$\hat{U}_{i,j} = \begin{cases} U_{i,j}, & \text{if } (i,j) \notin \text{TOP}_s(U) \\ 0, & \text{if } (i,j) \in \text{TOP}_s(U) \end{cases} \quad (\text{A.103})$$

and

$$\check{U}_{i,j} = \begin{cases} U_{i,j}, & \text{if } (i,j) \in \text{TOP}_s(U) \\ 0, & \text{if } (i,j) \notin \text{TOP}_s(U) \end{cases}. \quad (\text{A.104})$$

In the following we will compute an upper bound δ for the probability $\mathbb{P}(\hat{U} \cdot \check{U}^T = 0)$. By using the fact that $\hat{U} \cdot \check{U}^T = 0$ is a necessary condition for (\mathbf{F}_0, R_0) being a minimizer to Eq. (A.83), which is equivalent to $\varepsilon_{(1)} = \varepsilon_{(2)}$, we finally get

$$\mathbb{P}(\varepsilon_{(1)} < \varepsilon_{(2)}) = 1 - \mathbb{P}(\varepsilon_{(1)} \geq \varepsilon_{(2)}) \quad (\text{A.105})$$

$$= 1 - \mathbb{P}(\varepsilon_{(1)} = \varepsilon_{(2)}) \quad (\text{A.106})$$

$$\geq 1 - \mathbb{P}(\hat{U} \cdot \check{U}^T = 0). \quad (\text{A.107})$$

Thus, the last step is to find an upper bound $\delta \geq \mathbb{P}(\hat{U} \cdot \check{U}^T = 0)$. In order to compute δ , we have a closer look on $\hat{U} \cdot \check{U}^T$. It holds

$$(\hat{U} \cdot \check{U}^T)_{i,j} = \sum_{k=1}^n \hat{U}_{i,k} \check{U}_{j,k} \quad (\text{A.108})$$

$$= \begin{cases} \sum_{k \in \mathcal{T}_{i,j}} U_{i,k} U_{j,k}, & \text{if } \mathcal{T}_{i,j} \neq \emptyset \\ 0, & \text{else} \end{cases}, \quad (\text{A.109})$$

where for each $(i,j) \in \{1, \dots, m\}^2$,

$$\mathcal{T}_{i,j} := \{k \in \{1, \dots, n\} : (i,k) \notin \text{TOP}_s(U) \text{ and } (j,k) \in \text{TOP}_s(U)\}. \quad (\text{A.110})$$

Now assume $S \subset (\{1, \dots, m\} \times \{1, \dots, n\})^2$ with $S \neq \emptyset$ to be given, then

$$\mathbb{P}\left(\sum_{(i_1,j_1),(i_2,j_2) \in S} U_{i_1,j_1} \cdot U_{i_2,j_2} = 0\right) = 0. \quad (\text{A.111})$$

This equality holds since for each $S \subset (\{1, \dots, m\} \times \{1, \dots, n\})^2$ with $S \neq \emptyset$, $\sum_{(i_1,j_1),(i_2,j_2) \in S} U_{i_1,j_1} \cdot U_{i_2,j_2}$ follows a continuous probability distribution.

Consequently,

$$\mathbb{P}(\hat{U} \cdot \check{U}^T = 0) \quad (\text{A.112})$$

$$= \mathbb{P}(\forall(i, j) : \mathcal{T}_{i,j} = \emptyset \vee (\mathcal{T}_{i,j} \neq \emptyset \wedge \sum_{k \in \mathcal{T}_{i,j}} U_{i,k} \cdot U_{j,k} = 0)) \quad (\text{A.113})$$

$$\leq \mathbb{P}(\forall(i, j) : \mathcal{T}_{i,j} = \emptyset \vee (\exists S \subset (\{1, \dots, m\} \times \{1, \dots, n\})^2 \setminus \emptyset : \sum_{(i_1, j_1), (i_2, j_2) \in S} U_{i_1, j_1} \cdot U_{i_2, j_2} = 0)) \quad (\text{A.114})$$

$$\leq \mathbb{P}(\forall(i, j) : \mathcal{T}_{i,j} = \emptyset + \left(\sum_{\substack{S \subset (\{1, \dots, m\} \times \{1, \dots, n\})^2 \\ S \neq \emptyset}} \mathbb{P} \left(\sum_{(i_1, j_1), (i_2, j_2) \in S} U_{i_1, j_1} \cdot U_{i_2, j_2} = 0 \right) \right) \quad (\text{A.115})$$

$$= \mathbb{P}(\forall(i, j) : \mathcal{T}_{i,j} = \emptyset), \quad (\text{A.116})$$

where the inequality (A.115) uses the subadditivity of probability measures and the final equality (A.116) is achieved by using Eq. (A.111). By looking at the definition of $\mathcal{T}_{i,j}$, we see that $\forall(i, j) : \mathcal{T}_{i,j} = \emptyset$ only happens if for each $k \in \{1, \dots, n\}$ either

$$\forall i \in \{1, \dots, m\} : (i, k) \in \text{TOP}_s(U) \quad (\text{A.117})$$

or

$$\forall i \in \{1, \dots, m\} : (i, k) \notin \text{TOP}_s(U) \quad (\text{A.118})$$

holds true. Otherwise, if $k \in \{1, \dots, n\}$, $i, j \in \{1, \dots, m\}$ exist with $(i, k) \notin \text{TOP}_s(U)$ and $(j, k) \in \text{TOP}_s(U)$, obviously $\mathcal{T}_{i,j} \neq \emptyset$. This shows, that $\varepsilon_{(1)} = \varepsilon_{(2)}$ is only possible in the trivial case, where each of the n filters (with m coefficients) is either completely pruned or not pruned at all.

Therefore, we need to compute the probability

$$\mathbb{P}(\forall k : (\forall i : (i, k) \in \text{TOP}_s(U)) \vee (\forall i : (i, k) \notin \text{TOP}_s(U))) \quad (\text{A.119})$$

Due to the i.i.d. assumption of the $U_{i,k}$, all (i, k) have the same probability of being in $\text{TOP}_s(U)$. Thus, deciding $(i, k) \in \text{TOP}_s(U)$ or $(i, k) \notin \text{TOP}_s(U)$ for all $(i, k) \in \{1, \dots, m\} \times \{1, \dots, n\}$ together can equivalently be modeled with choosing a subset of size s from a set of size $m \cdot n$, where each subset has the same probability of being sampled, *i.e.* with probability $\frac{1}{\binom{m \cdot n}{s}}$.

Furthermore, Eq. (A.119) is only possible if $s = \alpha \cdot m$ for some $\alpha \in \mathbb{N}$. Otherwise, there needs to exist at least one

k, i, j such that $(i, k) \notin \text{TOP}_s(U)$ and $(j, k) \in \text{TOP}_s(U)$. Assuming $s = \alpha \cdot m$ for some α , there exist exactly $\binom{n}{\alpha}$ different choices to find α many k that satisfy $\forall i : (i, k) \in \text{TOP}_s(U)$ which, by the discussion above, all have similar probability.

Altogether, the probability Eq. (A.119) is given by

$$\delta = \begin{cases} 0 & , \text{ if } s \not\equiv 0 \pmod{m} \\ \binom{\frac{n}{\alpha}}{\frac{m \cdot n}{s}} & , \text{ if } s \equiv 0 \pmod{m} \end{cases} \quad (\text{A.120})$$

Finally,

$$\mathbb{P}(\hat{U} \cdot \check{U}^T = 0) \leq \mathbb{P}(\forall(i, j) : \mathcal{T}_{i,j} = \emptyset) \quad (\text{A.121})$$

$$= \mathbb{P}(\forall k : (\forall i : (i, k) \in \text{TOP}_s(U)) \vee (\forall i : (i, k) \notin \text{TOP}_s(U))) \quad (\text{A.122})$$

$$\leq \delta. \quad (\text{A.123})$$

Using the estimation in Eq. (A.107), we finally get

$$\mathbb{P}(\varepsilon_{(1)} < \varepsilon_{(2)}) \geq 1 - \mathbb{P}(\hat{U} \cdot \check{U}^T) \geq 1 - \delta, \quad (\text{A.124})$$

which finishes the first part of the proof where $\text{supp } R$ is fixed to be equal to $\text{supp } \Phi^*$.

Second part of proof with arbitrary $\text{supp } R$. As shown in the first part of the proof, Eq. (A.119) is a necessary condition for $\varepsilon_{(1)} = \varepsilon_{(2)}$. This means that all columns of Φ^* are either $\Phi_{:,k}^* = 0 \in \mathbb{R}^m$ or $\|\Phi_{:,k}^*\|_0 = m$ which we therefore will assume from now on.^{A.6}

By assumption, $0 < s$ and therefore, there exists a k with $\|\Phi_{:,k}^*\|_0 = m$. Now, set $\hat{\mathbf{F}}_{:,1} = \Phi_{:,k}^*$ and $\hat{\mathbf{F}}_{:,j} = (\mathbf{F}_0)_{:,j}$ for all other j . Then, with $\mathbb{P} = 1$, $\hat{\mathbf{F}}$ still forms a basis.

Setting $\hat{R}_{i,k} = \delta_{i,1}$ for all $i \in \{1, \dots, m\}$ yields $(\hat{\mathbf{F}} \cdot \hat{R})_{:,k} = U_{:,k} = \Phi_{:,k}^*$. For all other $j \neq k$ with $\|\Phi_{:,j}^*\|_0 = m$ there exists a $\hat{R}_{:,j}$ with $(\hat{\mathbf{F}} \cdot \hat{R})_{:,j} = \Phi_{:,j}^* = U_{:,j}$ since $\hat{\mathbf{F}}$ forms a basis.

Setting the remaining $\hat{R}_{:,j} = 0$ leads to $\hat{\mathbf{F}} \cdot \hat{R} = \Phi^*$ and $\|\hat{R}\|_0 \leq \|\Phi^*\|_0 - (m - 1) < \|\Phi^*\|_0$. The last inequality holds, since $m > 1$ is assumed.

Finally, one of the (at least) remaining $m - 1$ coefficients which were not spend up to now can be used to better approximate one column of U which is completely zeroed in Φ^* . Such a column j_0 must fulfill $U_{:,j_0} \neq 0$ and $\Phi_{:,j_0}^* = 0$. Since $s < m \cdot n$ and $U_{i,j}$ i.i.d. $\mathcal{N}(0, 1)$, such a column j_0 exists with $\mathbb{P} = 1$. Since $\hat{\mathbf{F}}$ forms a basis, we can find some l, λ_l such that

$$\|U_{:,j_0} - \lambda_l \hat{\mathbf{F}}_{:,l}\|_2 < \|U_{:,j_0}\|_2. \quad (\text{A.125})$$

^{A.6}For a matrix $A \in \mathbb{R}^{d_1 \times d_2}$, the j^{th} column $A_{:,j}$ is given by $(A_{i,j})_i \in \mathbb{R}^{d_1}$.

Setting $\hat{R}_{l,j_0} = \lambda_l$ leads to

$$\|U - \hat{\mathbf{F}} \cdot \hat{R}\|_F^2 = \|U - \Phi^*\|_F^2 \quad (\text{A.126})$$

$$- (\|U_{:,j_0}\|_2^2 - \|U_{:,j_0} - \lambda_l \hat{\mathbf{F}}_{:,l}\|_2^2) \quad (\text{A.127})$$

$$< \|U - \Phi^*\|_F^2, \quad (\text{A.128})$$

which finishes the proof. \square