

CNN Filter DB: An Empirical Investigation of Trained Convolutional Filters

Paul Gavrikov
IMLA, Offenburg University
paul.gavrikov@hs-offenburg.de

Janis Keuper
IMLA, Offenburg University & Fraunhofer ITWM
keuper@imla.ai

Abstract

Currently, many theoretical as well as practically relevant questions towards the transferability and robustness of Convolutional Neural Networks (CNNs) remain unsolved. While ongoing research efforts are engaging these problems from various angles, in most computer vision related cases these approaches can be generalized to investigations of the effects of distribution shifts in image data.

In this context, we propose to study the shifts in the learned weights of trained CNN models. Here we focus on the properties of the distributions of dominantly used 3×3 convolution filter kernels. We collected and publicly provide a data set with over 1.4 billion filters from hundreds of trained CNNs, using a wide range of data sets, architectures, and vision tasks. In a first use case of the proposed data set, we can show highly relevant properties of many publicly available pre-trained models for practical applications: I) We analyze distribution shifts (or the lack thereof) between trained filters along different axes of meta-parameters, like visual category of the data set, task, architecture, or layer depth. Based on these results, we conclude that model pre-training can succeed on arbitrary data sets if they meet size and variance conditions. II) We show that many pre-trained models contain degenerated filters which make them less robust and less suitable for fine-tuning on target applications.

Project website: <https://github.com/paulgavrikov/cnn-filter-db>

1. Introduction

Despite their overwhelming success in the application to various vision tasks, the practical deployment of convolutional neural networks (CNNs) is still suffering from several inherent drawbacks. Two prominent examples are I) the dependence on very large amounts of annotated training data [43], which is not available for all target domains and is expensive to generate; and II) still widely unsolved problems with the robustness and generalization abilities of CNNs [1] towards shifts of the input data distributions. One can argue that both problems are strongly related, since a

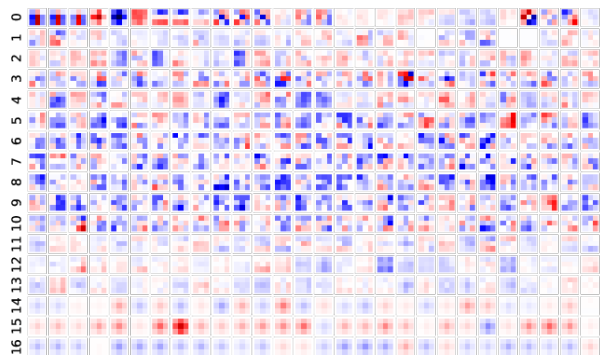


Figure 1. First learned filters extracted of each convolution layer with a 3×3 kernel of a *ResNet-18* trained on *CIFAR-10*. The filters show a clear loss of diversity and increasing sparsity with depth. The colormap range is determined layer-wise by the absolute peak weight of all filters in that layer.

common practical solution to I) is the fine-tuning [46] of pre-trained models by small data sets from the actual target domain. This results in the challenge to find suitable pre-trained models based on data distributions that are “as close as possible” to the target distributions. Hence, both cases (I+II) imply the need to model and observe distribution shifts in the contexts of CNNs.

In this paper, we propose not to investigate these shifts in the input (image) domain, but rather in the 2D filter-kernel distributions of the CNNs themselves. We argue that e.g. the distributions of trained convolutional filters in a CNN, which implicitly reflect the sub-distributions of the input image data, are more suitable and easier accessible representations for this task. In order to foster systematic investigations of learned filters, we collected and publicly provide a data set of over 1.4 billion filters with meta data from hundreds of trained CNNs, using a wide range of data sets, architectures, and vision tasks. To show the scientific value of this new data source, we conduct a first analysis and report a series of novel insights into widely used CNN models. Based on our presented methods we show that many publicly provided models suffer from degeneration.

We show that overparameterization leads to sparse and/or non-diverse filters (Fig. 1), while robust training increases filter diversity, and reduces sparsity. Our results also show that learned filters do not significantly differ across models trained for various tasks, except for extreme outliers such as GAN-Discriminators. Models trained on data sets of different visual categories do not significantly drift either. Most shifts in studied models are due to degeneration, rather than an actual difference in structure. Therefore, our results imply that pre-training can be performed independent of the actual target data, and only the amount of training data and its diversity matters. This is inline with recent findings that models can be pre-trained even with abstract images [23]. For classification models we show that the most variance in learned filters is found in the beginning and end of the model, while object/face detection models only show significant variance in early layers. Also, the most specialized filters are found in the last layers. We summarize our key contributions as follows:

- Publication of a diverse database of over 1.4B 3×3 convolution filters alongside with relevant meta information of the extracted filters and models.
- Presentation of a data-agnostic method based on sparsity and entropy of filters to find “degenerated” convolution layers due to overparameterization or non-convergence of trained CNN models.
- Showing that publicly available models often contain degenerated layers and can therefore be questionable candidates for transfer tasks.
- Analysis of distribution shifts in filters over various groups, providing insights that formed filters are fairly similar across a wide-range of examined groups.
- Showing that the model-to-model shifts that exist in classification models are, contrary to the predominant opinion, not only seen in deeper layers but also in the first layers.

Paper organization. We frame our work in the context of related works in Sec. 2. Next, we give an overview of our data set and its collection process in Sec. 3. We follow up by introducing methods to study filter structure, distributions shifts, and layer degeneration such as randomness, low variance in filter structure, and high sparsity of filters. Then in Sec. 4 we apply these methods to our collected data. We show the impact of overparameterization and robust training on filter degeneration and provide intuitions for threshold finding. Then we analyze filter structures by determining a suitable filter basis and looking into reproducibility of filters in training, filter formation during training, and an analysis of distribution shifts for various dimensions of the collected meta-data. We discuss limitations of our approach in Sec. 5 and, finally, draw conclusions in Sec. 6.

2. Related Work

We are unaware of any systematic, large scale analysis of learned filters across a wide range of data sets, architectures and task such as the one performed in this paper. However, there are of course several partially overlapping aspects of our analysis that have been covered in related works:

Filter analysis. An extensive analysis of features, connections, and their organization extracted from trained *InceptionV1* [44] models was presented in [5, 6, 34–36, 39, 41, 49, 50]. The authors claim different CNNs will form similar features and circuits even when trained for different tasks.

Transfer learning. A survey on transfer learning for image classification CNNs can be found in [21] and general surveys for other tasks and domains are available in [37, 56]. The authors of [54] studied learned filter representations in *ImageNet1k* classification models and presented the first moves towards transfer learning. They argued that different CNNs will form similar filters in early layers which will mostly resemble gabor-filters and color-blobs, while deeper layers will capture specifics of the data set by forming increasingly specialized filters. [3] captured convolution filter pattern distributions with Gaussian Mixture Models to achieve cross-architecture transfer learning. [48] demonstrated that convolutions filters can be replaced by a fixed filter basis that 1×1 convolution layers blend.

Pruning criteria. Although we do not attempt pruning, our work overlaps with pruning techniques as they commonly rely on estimation criteria to understand which parameters to compress. These either rely on data-driven computation of a forward-pass [2, 18, 30–32], or backward-propagation [11, 55], or estimate importance solely based on the numerical weight (typically any ℓ -norm) of the parameters [14, 16, 17, 28, 29].

CNN distribution shifts. A benchmark for distribution shifts that arise in real-world applications is provided in [24] and [47] measured robustness to natural distribution shifts of 204 *ImageNet1k* models. The authors concluded that robustness to real-world shifts is low. Lastly, [12] studied the correlation between transfer performance and distribution shifts of image classification models and find that increasing training set and model capacity increases robustness to distribution shifts.

3. Methods

3.1. Collecting filters

We collected a total of 647 publicly available CNN models from [8, 38, 52] and other sources that have been pre-trained for various 2D visual tasks¹. In order to provide a heterogeneous and diverse representation of convolution filters “in the wild”, we retrieved pre-trained models for

¹For more details refer to the supplementary materials.

11 different tasks e.g. such as *classification*, *segmentation* and *image generation*. We also recorded various meta-data such as depth and frequency of included operations for each model, and manually categorized the variety of used training sets into 16 visually distinctive groups like *natural scenes*, *medical ct*, *seismic*, or *astronomy*. In total, the models were trained on 71 different data sets. The dominant subset is formed by *image classification* models trained on *ImageNet1k* [9] (355 models).

All models were trained with full 32-bit precision² but may have been trained on variously scaled input data. Included in the data set are low-resolution variants of AlexNet [26], DenseNet-121/161/169 [20], ResNet-9/14/18/34/50/101/152 [15], VGG-11/13/16/19 [42], MobileNet v2 [40], Inception v3 [45] and GoogLeNet [44] image classification models that we have purposely trained on simple data sets such as CIFAR-10/100 [25], MNIST [10], Kuzushiji-MNIST (KMIST) [7] and Fashion-MNIST [53] in order to study the effect of overparameterization on learned filters.

All collected models were converted into the ONNX format [4] which allows a streamlined filter extraction without framework dependencies. Hereby, only the widely used filters from regular convolution layers with a kernel size of 3×3 were taken into account. Transposed (sometimes also called de-convolution or up-convolution) convolution layers were not included. In total, **1,464,797,156 filters** from **21,436 layers** have been obtained for our data set³.

3.2. Analyzing filter structures

We apply a full-rank principal component analysis (PCA) transformation implemented via a singular-value decomposition (SVD) to understand the underlying structure of the filters [22].

First, we stack the relevant set of n flattened filters into a $n \times 9$ matrix X . Thereupon, we center the matrix and perform a SVD into a $n \times 9$ rotation matrix U , a 9×9 diagonal scaling matrix Σ , and a 9×9 rotation matrix V^T . The diagonal entries $\sigma_i, i = 0, \dots, n - 1$ of Σ form the singular values in decreasing order of their magnitude. Row vectors $v_i, i = 0, \dots, n - 1$ in V^T then form the principal components. Every row vector $c_{ij}, j = 0, \dots, n - 1$ in U is the coefficient vector for f_i .

$$X^* = X - \bar{X} = U\Sigma V^T \quad (1)$$

Where \bar{X} denotes the vector of column-wise mean values of any matrix X . Then we obtain a vector \hat{a} of the explained variance ratio of each principal component. $\|\cdot\|_1$ denotes

²Although, initial experiments indicated that mixed/reduced precision training [33] does not affect distribution shifts beyond noise.

³The data set with all meta data, full citations, and code will be released publicly with the camera ready submission.

the ℓ_1 -norm.

$$\begin{aligned} \vec{a} &= (\Sigma \cdot I)^2 / (n - 1) \\ \hat{a} &= \vec{a} / \|\vec{a}\|_1 \end{aligned} \quad (2)$$

Finally, each filter f' is described by a linear, shifted sum of principal components v_i weighted by the coefficients c_i .

$$f' = \sum_i c_i v_i + \bar{X}_i \quad (3)$$

3.3. Measuring distribution shifts

All probability distributions are represented by histograms. The histogram range is defined by the minimum and maximum value of all coefficients. Each histogram is divided into 70 uniform bins. The divergence between two distributions is measured by the symmetric, non-negative variant of Kullback-Leibler (KL_{sym}) [27].

$$KL(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)} \quad (4)$$

$$KL_{\text{sym}}(P \parallel Q) = KL(P \parallel Q) + KL(Q \parallel P)$$

We define the drift D between two filter sets by the sum of the divergence of the coefficient distributions P_i, Q_i along every principal component index i . The sum is weighted by the ratio of variance \hat{a}_i explained by the i -th principal component.

$$D(P \parallel Q) = \sum_i \hat{a}_i \cdot KL_{\text{sym}}(P_i \parallel Q_i) \quad (5)$$

To avoid undefined expressions, all probability distributions F are set to hold $\forall x \in \mathcal{X} : F(x) \geq \epsilon$.

3.4. Measuring layer degeneration

Lottery Ticket Hypothesis [13] suggests that each architecture has a specific amount of convolution filters that saturate its ability to transform a given data set into a well separable feature-space. Exceeding this number will result in a partitioning of the model into multiple inter-connected sub-models. We hypothesize that these are seen in the form of degenerated filters in CNNs. In like manner, an insufficient amount of training samples or training epochs will also lead to degenerated filters. We characterize the following types of degeneration.

1. **High sparsity:** Filters are dominantly close to zero and therefore produce quasi-zero feature-maps [29]. These feature-maps carry no vital information and can be discarded.
2. **Low diversity in structure:** Filters are structurally similar to each other and therefore redundant. They produce similar feature-maps in different scales and could be represented by a subset of present filters.

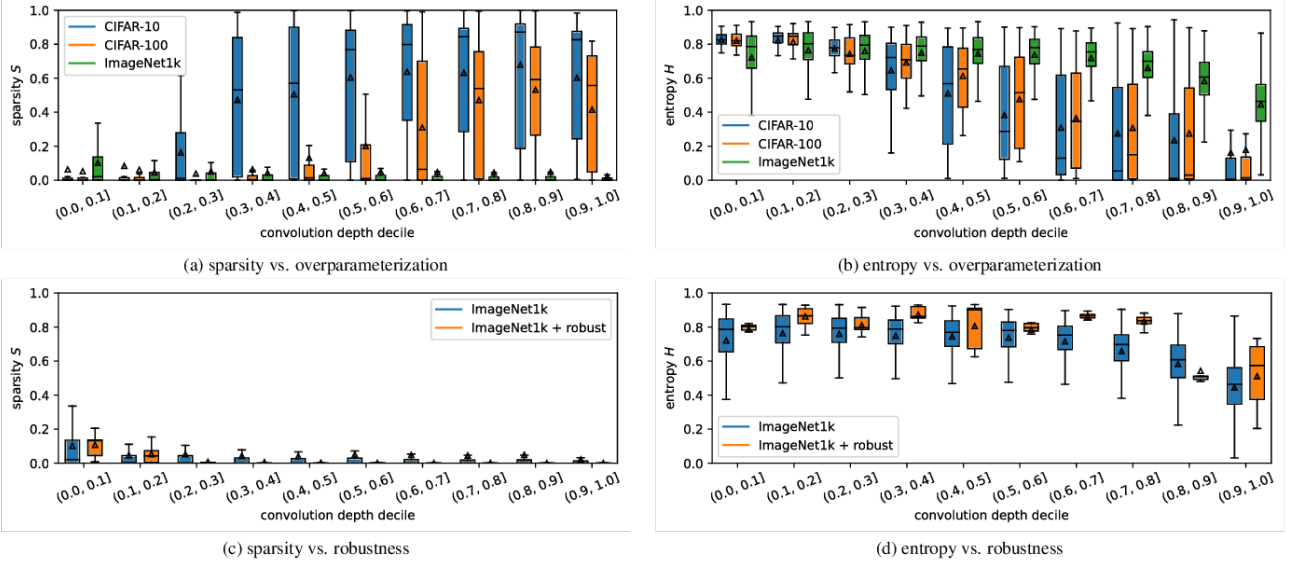


Figure 2. Comparison of layer entropy and sparsity of overparameterized, robust, and regular classification models. Outliers are hidden for clarity.

3. **Randomness:** Filter weights are conditionally independent of their neighbours. This indicates that no or not sufficient training was performed.

Sparsity degeneration is detectable by the share of sparse filters S in a given layer. We call a filter f sparse if all entries are near-zero. Consequently, given the number of input channels c_{in} , number of output channels c_{out} , and a set of filters in layer L , we can measure the layer sparsity by:

$$S(L) = \frac{|\{f|f \in L \wedge (\forall w \in f : -\epsilon_0 \leq w \leq \epsilon_0)\}|}{c_{in}c_{out}} \quad (6)$$

To detect the other types of degeneration we introduce a layer-wise metric based on the Shannon-Entropy of the explained variance ratio of each principal component obtained from a SVD of all filters in the examined layer (Sec. 3.2).

$$H = - \sum_i \hat{a}_i \log_{10} \hat{a}_i \quad (7)$$

If H is close to zero this indicates one strong principal component from which most of the filters can be reconstructed and is therefore a low filter diversity degeneration. On the other hand, a large entropy indicates a (close to) uniform distribution of the singular values and, thus, a randomness of the filters. Sparse layers are a specific form of low diversity degeneration and, generally both are correlated, whereas, sparsity and randomness are mutually exclusive. It should be noted, that $|\Sigma \cdot I| = \min(c_{in}c_{out}, 9)$ and therefore the entropy only becomes expressive if $c_{in}c_{out} \gg 9$.

4. Results: Analysis of trained CNN filters

4.1. Layer degeneration

In this section we study different causes of degeneration and aim provide thresholds for evaluation.

Overparameterization. The majority of the models that we have trained on our low resolution data sets are heavily overparameterized for these relatively simple problems. We base this argument on the fact that we have models with different depth for most architectures and already observe near perfect performance with the smallest variants. Therefore it is safe to assume that larger models are overparameterized especially given that the performance only increases marginally¹.

First we analyze layer sparsity and entropy for these models trained on *CIFAR-10/100* in comparison to all *ImageNet1k* classification models found in our data set. For each data set we have trained identical networks with identical hyperparameters. Both, *CIFAR-10* and *CIFAR-100*, consist of 60,000 32×32 px images, but *CIFAR-100* includes 10x more labels and thus fewer samples per class forming a more challenging data set.

Fig. 2a shows that the overparameterized models contain significantly more sparse filters on average, and that sparsity increases with depth. In particular, we see the most sparse filters for *CIFAR-10*. However, *ImageNet1k* classifiers also seem to have some kind of “natural” sparsity, even though we do not consider most of these models as overparameterized. Entropy, on the other hand, decreases with increasing layer depth for every classifier, but more rapidly in overpa-

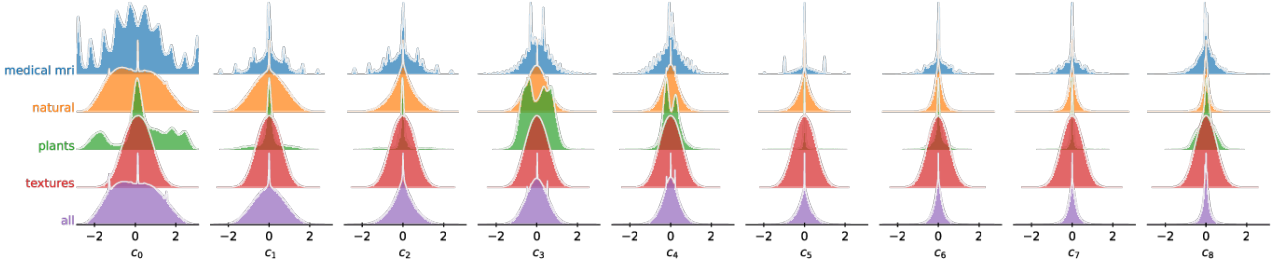


Figure 3. KDEs of the coefficient distributions along every principal component for selected¹ visual categories.

parameterized models (Fig. 2b). Again, the *CIFAR-10* models degrade faster and show more degeneration.

The overparameterized models contain layers that have an entropy close to 0 towards deeper layers which indicates that these models are “saturated” and only produce differently scaled variants of the same filters. In line with the oversaturation, these models also have increasingly sparse filters, presumably as an effect of regularization.

Filter degeneration and model robustness. Our data set also contains robust models from the *RobustBench leaderboard* [8]. When comparing robust models with non-robust models trained on *ImageNet1k*, it becomes clear that robust models form almost no sparse filters after in deeper convolution layers (Fig. 2c), while regular models show some sparsity there. The entropy of robust models is also higher throughout depth (Fig. 2d), indicating that robust models learn more diverse filters.

Thresholds. To obtain a threshold for randomness given a number of filters n per layer we perform multiple experiments in which we initialize convolution filters of different sizes from a standard normal distribution and fit a sigmoid T_H to the minimum results obtained for entropy.

$$T_H(n) = \frac{L}{1 + e^{-k(\log_2(n) - x_0)}} + b \quad (8)$$

We obtain the following values $L = 1.26$, $x_0 = 2.30$, $k = 0.89$, $b = -0.31$ and call any layer L with $H > T_H(n)$ random. On the opposite, defining a threshold for low diversity degeneration seems less intuitive and one can only rely on statistics: The average entropy H is 0.69 over all layers and continuously decreases from an average of 0.75 to 0.5 with depth. Additionally, the minimum of the 1.5 IQR also steadily decreases with depth.

The same applies to sparsity: the average sparsity S over all layers is 0.12 and only 56.5% of the layers in our data set hold $S < 0.01$ and 9.9% even show $S > 0.5$. In terms of convolution depth, the average sparsity varies between 9.9% and 14% with the largest sparsity found in the last 20% of the model depth. The largest outliers of the 1.5

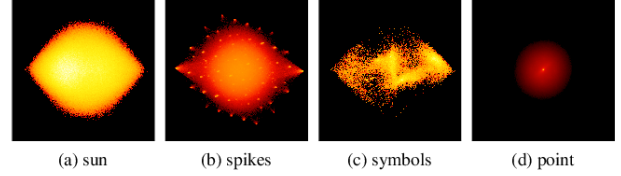


Figure 4. Bi-variate plots between component distributions showing the four phenotypes.

interquartile range (IQR) are, however, found in the first decile. In both cases we find it difficult to provide a meaningful general threshold and suggest to determine this value on a case-by-case basis¹.

4.2. Filter structure

In the next series of experiments, we analyze only the structure of 3×3 filters, neglecting their actual numerical weight in the trained models. Therefore, we normalize each filter f individually by the absolute maximum weight into f' .

$$d_i = \max_{i,j} |f_{ij}|$$

$$f'_{ij} = \begin{cases} f_{ij}/d_i, & \text{if } d_i \neq 0 \\ f_{ij}, & \text{else} \end{cases} \quad (9)$$

Then we perform a PCA transformation on the scaled filters. Fig. 5 shows some qualitative examples of obtained principal components, split by several meta-data dimensions. The images of the formed basis are often similar for all groups except for few outliers (such as *GAN-discriminators*). The explained variance however fluctuates significantly and sometimes changes the order of components. Consistently, we observe substantially higher variance on the first principal components. The explained variance does not necessarily correlate with the shift observed between models. Here, the biggest mean drift is also located in the first principal component ($\hat{D} = 0.90$), but is then followed by the sixth, third, second component ($\hat{D} = 0.78, 0.69, 0.58$). The coefficients of the sixth component also contain the strongest outliers (Fig. 6). We visualize the

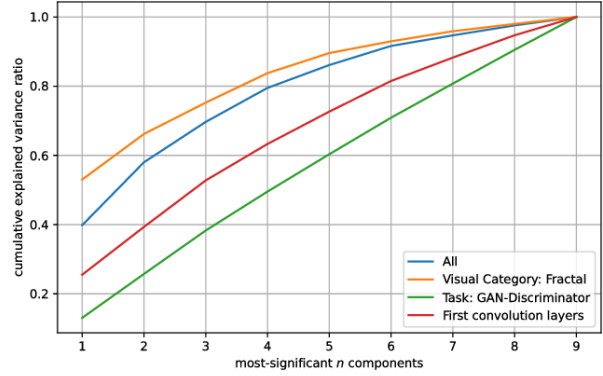
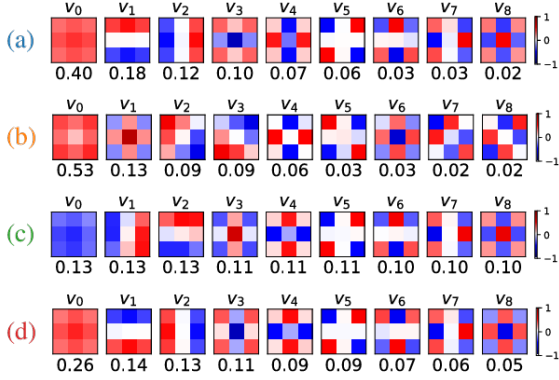


Figure 5. Selected¹ depiction of the filter basis and (cumulative) explained variance ratio per component for filters from (a) full data set, (b) models trained on images of *fractals*, (c) GAN discriminators, (d) first convolution layers.

distributions of PCA coefficients along every component for each group by plots of kernel density estimates (KDEs), e.g. Fig. 3 depicts the distributions of filters grouped by some selected visual categories in comparison to the distribution of coefficients for the full data set. Filters extracted from models with degenerated layers (as seen in *medical mri*) result in spiky/multi-modal KDEs. The distributions can alternatively be visualized by bi-variate scatter plots that may reveal more details than KDEs. For example, they let us categorize the distributions into phenotypes depending on their distribution characteristic in the PCA space (Fig. 4): *sun*: distributions where both dimensions are gaussian-like. These are to be expected coefficient distributions without significant sparsity/low diversity degeneration. Yet, this phenotype may also include non-converged filters; *spikes*: distributions suffering from a low variance degeneration resulting in local hotspots; *symbols*: at least one distribution is multi-modal, non-centered, highly sparse or otherwise non-normal (low variance degeneration); *point*: coefficients are primarily located in the center (sparsity degeneration).

Reproducibility of filters. We train low-resolution networks on *CIFAR-10* multiple times with identical hyperparameters except for random seeds and save checkpoints of each model at the best validation epoch. Most models are converging to highly similar coefficient distributions when retrained with different weight initialization (e.g. ResNet-9 with $D < 5.3 \cdot 10^{-4}$). However, some architectures such as *MobileNetv2* show higher shifts ($D < 2.6 \cdot 10^{-2}$). We assume that this is due to the structure of the loss surface, e.g. the residual skip connections found in *ResNets* smooth the surface, whereas other networks may contain more local minima due to noisy surfaces [51].

Formation of filter structures during training. Although our data set only includes trained convolutional filters we tried to understand how the coefficient distribution shifts

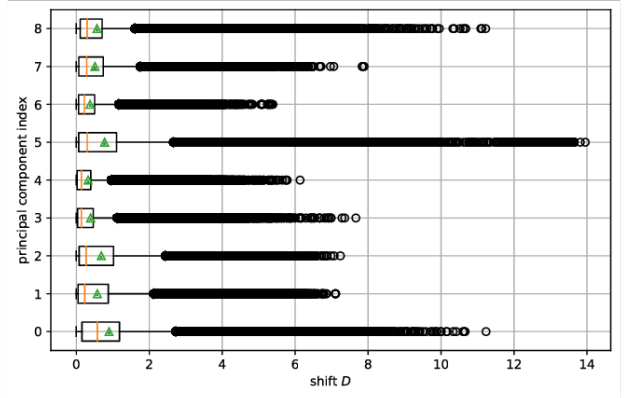


Figure 6. Distribution of the shift D along principal components computed on all possible pairings of models.

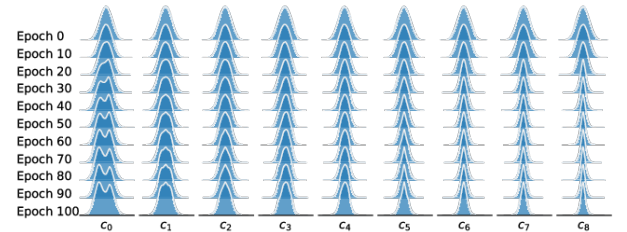


Figure 7. Coefficient distribution of a *ResNet-9* trained on *CIFAR-10* every 10 epochs.

during training. Therefore we recorded checkpoints of a *ResNet-9* trained on *CIFAR-10* every 10 training epochs beginning right after the weight initialization. Fig. 7 shows that the coefficient distributions along all principal components are gaussian-like distributed in the beginning and eventually shift during training. For this specific model, distributions along major principal components retain the standard deviation during training, while less-significant com-

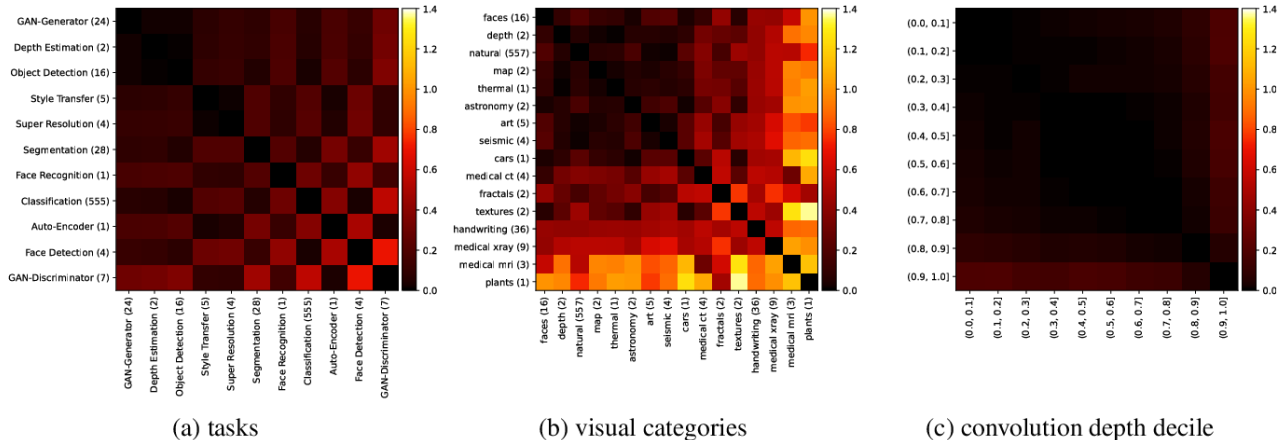


Figure 8. Heatmaps over the shift D for different filters groupings. The number in brackets denotes the number of models in this group. Low values/dark colors denote low shifts.

ponent distributions decrease. The initialization observation helped us removing models from our collection where we failed to load the trained parameters for any reason and is foundation for our provided randomness metric.

4.3. Distribution shifts between trained models

In this subsection we are investigating transfer distance in different meta-dimensions of pre-trained models. We compute the shift D and visualize this in the form of heatmaps (Fig. 8) that show shifts between all pairings.

Shifts between tasks. Unsurprisingly, *classification*, *segmentation*, *object detection*, and *GAN-generator* distributions are quite similar, since the non-*classification* models typically include a *classification* backbone. The smallest mean shift to other tasks is observed in *object detection*, *GAN-generators*, and *depth estimation* models. The least transferable distributions are *GAN-discriminators*. Their distributions do barely differ along principal components and can be approximated by a gaussian distribution. By our randomness metric this indicates a filter distribution that is close to random initialization, implying a “confused” *discriminator* that cannot distinguish between real and fake samples towards the end of (successful) training. It may be surprising to see a slightly larger average shift for *classification*. This is presumably due to many degenerated layers in our collected models, which are also visible in the form of spikes when studying the KDEs. An evaluation¹ of distributions including only non-degenerated *classifiers* actually shows a lower average shift due to the aforementioned similarity to other tasks.

Shifts between visual categories and training sets. We find that the distribution shift is well balanced across most visual categories and training sets. Notable outliers include

all *medical* types. They have visible spikes in the KDEs, once again indicating degenerated layers. Indeed, the average sparsity in these models is extreme in the last 80% of the model depth. Another interesting, albeit less significant outlier is the *fractal* category. It consists of models trained on *Fractal-DB*, which was proposed as a synthetic pre-training alternative to *ImageNet1k* [23]. The standard deviation of coefficient distributions tend to shrink towards the least significant principal components but this trend is not visible for this category indicating that sorting the basis by variance would yield a different order for this task and perhaps the basis itself is not well suited. Also notable is a remarkably high standard deviation on the distribution of the first principal component. Interestingly, we also observe sub-average degeneration for this category. Shifts in other categories can usually be explained by a biased representation. For example we only have one model for *plants*, our *handwriting* models consist exclusively of overparameterized networks that suffer from layer degeneration, and *textures* consists of only one *GAN-discriminator* which will naturally show a high randomness.

Shifts by filter/layer depth. The shift between layers of various depth deciles increases with the difference in depth, with distributions in the last decile of depth forming the most distinct interval, and outdistancing the second-to-last and first decile that follow next. An interesting aspect is also the model-to-model shift across deciles. This shift exemplifies the uniqueness of formed filters. Our observations overhaul the general recommendation for fine-tuning to freeze early layers in *classification* models, as the largest shifts are not only seen in deep layers but also in early vision (Fig. 9). *Segmentation*¹ models show the most drift in deeper layers. Contrary, *object/face detection* models only show drift in the early vision (*object detection* in the first,

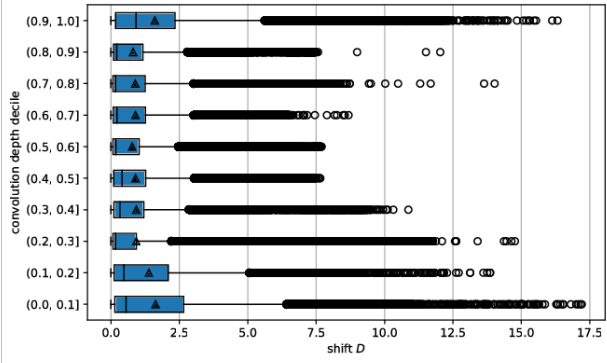


Figure 9. Boxplots showing the distribution of pair-wise model-to-model shift D of *classification* models per convolution depth decile (top to bottom in decreasing order). Our intentionally overparameterized models were left out of this analysis.

face detection in the first four depth deciles), but marginal drift in later convolution stages.

Shifts within model families. The shift between models of the same family trained for the same task is negligible (Fig. 10), indicating that every large enough data set is good enough and the common practice of pre-training models with *ImageNet1k* even for visually distant application domains is indeed a valid approach. *ResNet*-family outliers only consist out of models that show a high amount of sparsity. Additionally, this observation may be exploited by training small teacher networks and apply knowledge distillation [19] to initialize deeper models of the same family.

5. Limitations

Our data is biased against *classification* models and/or *natural* data sets such as *ImageNet1k*. Additionally, some splits will over-represent specific dimensions e.g. tasks may include exclusive visual categories and vice versa. Also, as previously shown, many of the collected models show a large amount of degenerated layers that significantly impact the distributions. This also biases measurements of the distribution shifts. We performed an ablation study by removing filters extracted from degenerated layer, but were unable to find a clear correlation between degeneration and distribution shifts¹, presumably due to a lack of justified thresholds.

6. Conclusions

Our first results support our initial hypothesis that the distributions of trained convolutional filters are a suitable and easy-to-access proxy for the investigation of image distributions in the context of transferring pre-trained models and

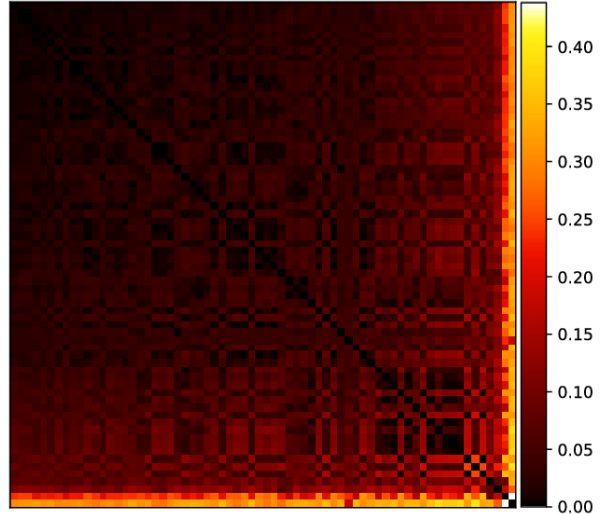


Figure 10. Heatmap over the shift D between different pairings of *ResNet-classifiers*. Each row/column depicts one model. Intentionally overparameterized models were not included.

robustness. While the presented results are still in the early stages of a thorough study, we report several interesting findings that could be explored to obtain better model generalizations and assist in finding suitable pre-trained models for fine-tuning. One finding is the presence of large amounts of degenerated (or untrained) filters in large, well-performing networks - resulting in the phenotypes *points*, *spikes*, and *symbols*. We assume that their existence is a symptom in line with the *Lottery Ticket Hypothesis* [13]. We conclude that ideal models should have relatively high entropy (but $H < T_H$) throughout all layers and almost no sparse filters. Models that show an increasing or generally high sparsity or a massive surge in entropy with depth are most likely overparameterized and could be pruned, which would benefit inference and training speed. Whereas, initialized but not trained models will have a constantly high entropy $H \geq T_H$ throughout all layers and virtually no sparsity.

Another striking finding is the observation of very low shifts in filter structure between different meta-groups: I) shifts inside a family of architectures are very low; II) shifts are mostly independent of the target image distribution and task; III) also we observe rather small shifts between convolution layers of different depths with the highest shifts in the first and last layers. Overall, the analysis of over 1.4 billion learned convolutional filters in the provided data set gives a strong indication that the common practice of pre-training CNNs is indeed a sufficient approach if the chosen model is not heavily overparameterized. Our first results indicate that the presented data set is a rich source for further research in transfer learning, robustness and pruning.

References

- [1] Naveed Akhtar and Ajmal Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *Ieee Access*, 6:14410–14430, 2018. 1
- [2] Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes, 2018. 2
- [3] Mehmet Aygun, Yusuf Aytar, and Hazim Kemal Ekenel. Exploiting convolution filter patterns for transfer learning. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV) Workshops*, Oct 2017. 2
- [4] Junjie Bai, Fang Lu, Ke Zhang, et al. Onnx: Open neural network exchange. <https://github.com/onnx/onnx>, 2019. 3
- [5] Nick Cammarata, Gabriel Goh, Shan Carter, Ludwig Schubert, Michael Petrov, and Chris Olah. Curve detectors. *Distill*, 2020. <https://distill.pub/2020/circuits/curve-detectors>. 2
- [6] Nick Cammarata, Gabriel Goh, Shan Carter, Chelsea Voss, Ludwig Schubert, and Chris Olah. Curve circuits. *Distill*, 2021. <https://distill.pub/2020/circuits/curve-circuits>. 2
- [7] Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha. Deep learning for classical japanese literature, 2018. 3
- [8] Francesco Croce, Maksym Andriushchenko, Vikash Sehwag, Edoardo Debenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial robustness benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. 2, 5
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 3
- [10] Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012. 3
- [11] Xiaohan Ding, Guiguang Ding, Yuchen Guo, and Jungong Han. Centripetal sgd for pruning very deep convolutional networks with complicated structure. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4938–4948, 2019. 2
- [12] Josip Djolonga, Jessica Yung, Michael Tschannen, Rob Romijnders, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver, Matthias Minderer, Alexander D’Amour, Dan Moldovan, Sylvain Gelly, Neil Houlsby, Xiaohua Zhai, and Mario Lucic. On robustness and transferability of convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16458–16468, June 2021. 2
- [13] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018. 3, 8
- [14] Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural networks, 2015. 2
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. 3
- [16] Yang He, Guoliang Kang, Xuanyi Dong, Yanwei Fu, and Yi Yang. Soft filter pruning for accelerating deep convolutional neural networks. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pages 2234–2240. International Joint Conferences on Artificial Intelligence Organization, 7 2018. 2
- [17] Yang He, Ping Liu, Ziwei Wang, Zhilan Hu, and Yi Yang. Filter pruning via geometric median for deep convolutional neural networks acceleration. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4335–4344, 2019. 2
- [18] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 1398–1406, 2017. 2
- [19] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015. 8
- [20] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 3
- [21] Mahbub Hussain, Jordan J. Bird, and Diego R. Faria. A study on cnn transfer learning for image classification. In Ahmad Lotfi, Hamid Bouchachia, Alexander Gegov, Caroline Langensiepen, and Martin McGinnity, editors, *Advances in Computational Intelligence Systems*, pages 191–202, Cham, 2019. Springer International Publishing. 2
- [22] I Jolliffe. *Principal Component Analysis*. Springer New York, New York, NY, 1986. 3
- [23] Hirokatsu Kataoka, Kazushige Okayasu, Asato Matsumoto, Eisuke Yamagata, Ryosuke Yamada, Nakamasa Inoue, Akio Nakamura, and Yutaka Satoh. Pre-training without natural images. 2020. 2, 7
- [24] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton Earnshaw, Imran Haque, Sara M Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. Wilds: A benchmark of in-the-wild distribution shifts. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 5637–5664. PMLR, 18–24 Jul 2021. 2
- [25] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). 3
- [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012. 3
- [27] S. Kullback and R. A. Leibler. On Information and Sufficiency. *The Annals of Mathematical Statistics*, 22(1):79 – 86, 1951. 3

- [28] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets, 2017. [2](#)
- [29] Y. Li, S. Lin, B. Zhang, J. Liu, D. Doermann, Y. Wu, F. Huang, and R. Ji. Exploiting kernel sparsity and entropy for interpretable cnn compression. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2795–2804, Los Alamitos, CA, USA, jun 2019. IEEE Computer Society. [2](#), [3](#)
- [30] Shaohui Lin, Rongrong Ji, Yuchao Li, Yongjian Wu, Feiyue Huang, and Baochang Zhang. Accelerating convolutional networks via global & dynamic filter pruning. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pages 2425–2432. International Joint Conferences on Artificial Intelligence Organization, 7 2018. [2](#)
- [31] Shaohui Lin, Rongrong Ji, Chenqian Yan, Baochang Zhang, Liujuan Cao, Qixiang Ye, Feiyue Huang, and David Doermann. Towards optimal structured cnn pruning via generative adversarial learning. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2785–2794, 2019. [2](#)
- [32] Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural network compression, 2017. [2](#)
- [33] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. Mixed precision training, 2018. [3](#)
- [34] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. An overview of early vision in inceptionv1. *Distill*, 2020. <https://distill.pub/2020/circuits/early-vision>. [2](#)
- [35] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. *Distill*, 5, 2020. [2](#)
- [36] Chris Olah, Nick Cammarata, Chelsea Voss, Ludwig Schubert, and Gabriel Goh. Naturally occurring equivariance in neural networks. *Distill*, 2020. <https://distill.pub/2020/circuits/equivariance>. [2](#)
- [37] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10):1345–1359, 2010. [2](#)
- [38] 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 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. [2](#)
- [39] Michael Petrov, Chelsea Voss, Ludwig Schubert, Nick Cammarata, Gabriel Goh, and Chris Olah. Weight banding. *Distill*, 2021. <https://distill.pub/2020/circuits/weight-banding>. [2](#)
- [40] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks, 2019. [3](#)
- [41] Ludwig Schubert, Chelsea Voss, Nick Cammarata, Gabriel Goh, and Chris Olah. High-low frequency detectors. *Distill*, 2021. <https://distill.pub/2020/circuits/frequency-edges>. [2](#)
- [42] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015. [3](#)
- [43] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, pages 843–852, 2017. [1](#)
- [44] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions, 2014. [2](#), [3](#)
- [45] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. [3](#)
- [46] Nima Tajbakhsh, Jae Y Shin, Suryakanth R Gurudu, R Todd Hurst, Christopher B Kendall, Michael B Gotway, and Jianming Liang. Convolutional neural networks for medical image analysis: Full training or fine tuning? *IEEE transactions on medical imaging*, 35(5):1299–1312, 2016. [1](#)
- [47] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring robustness to natural distribution shifts in image classification, 2020. [2](#)
- [48] Muhammad Tayyab and Abhijit Mahalanobis. Basisconv: A method for compressed representation and learning in cnns, 2019. [2](#)
- [49] Chelsea Voss, Nick Cammarata, Gabriel Goh, Michael Petrov, Ludwig Schubert, Ben Egan, Swee Kiat Lim, and Chris Olah. Visualizing weights. *Distill*, 2021. <https://distill.pub/2020/circuits/visualizing-weights>. [2](#)
- [50] Chelsea Voss, Gabriel Goh, Nick Cammarata, Michael Petrov, Ludwig Schubert, and Chris Olah. Branch specialization. *Distill*, 2021. <https://distill.pub/2020/circuits/branch-specialization>. [2](#)
- [51] Lifu Wang, Bo Shen, Ning Zhao, and Zhiyuan Zhang. Is the skip connection provable to reform the neural network loss landscape?, 2020. [6](#)
- [52] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019. [2](#)
- [53] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017. [3](#)
- [54] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? volume 4, 2014. [2](#)
- [55] Ruichi Yu, Ang Li, Chun-Fu Chen, Jui-Hsin Lai, Vlad I. Morariu, Xintong Han, Mingfei Gao, Ching-Yung Lin, and Larry S. Davis. Nisp: Pruning networks using neuron importance score propagation. In *2018 IEEE/CVF Conference*

on Computer Vision and Pattern Recognition, pages 9194–9203, 2018. 2

- [56] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning, 2020. 2