

Learning the *Essential* in less than 2k additional weights - a single large kernel convolution layer can improve prediction stability under unknown corruptions

Supplementary Material

952 Overview

953 In this supplementary material, we provide additional information and details to the experimental results in the main
 954 paper. In particular, we provide
 955

- 956 • More information and analysis on the employed datasets
 957 in section 7.
- 958 • More details in the experiment setup in section 8.
- 959 • Additional Experimental results including more details on
 960 the aggregated results given in the main paper, an ablation
 961 on kernel sizes, an ablation on thresholded kernel performance
 962 over thresholds and ablation on L_1 prior regularization
 963 strengths in section 9.
- 964 • An analysis of the optimized kernels in section 10.
- 965 • An extended discussion of our findings in section 11.

966 7. Datasets

967 To train the different models, subsets of the ImageNet [46]
 968 dataset are used. To have a widespread experiment set, the
 969 ImageNette [27] dataset with ten classes and the ImageNet-
 970 100 [55] with 100 classes from ImageNet are used.

971 Furthermore, to evaluate the robustness of these models,
 972 the subsets are corrupted via different corruptions from OpticsBench [39] and Common Corruptions [20].

974 7.1. Clean datasets

975 ImageNette provides 9,469 training and 3,925 validation
 976 images on ten diverse ImageNet classes. We use the ImageNette2 version in all experiments, publicly available
 977 from [27]. This allows extensive model training experiments to be carried out with limited computing resources.
 978 In addition, some of the experiments were carried out on ImageNet-100, consisting of 100 randomly selected classes
 979 from ImageNet-1k with a total of 128k training and 5k validation images. Each class has approximately 1300 training
 980 and 50 validation images. As there are no extra test labels, we test on the validation images. The validation images are
 981 kept from training. Finally, two experiments are conducted
 982 on ImageNet-1k, which consists of 1,000 classes with a total of 1,281k training and 50k validation images.

989 7.2. Common Corruptions

990 The diverse common corruptions [20] consist of digital,
 991 weather, noise and blur corruptions aimed at benchmarking
 992 DNNs in safety-critical applications. We use 16 different
 993 algorithmically generated corruptions, organized into five
 994 different severities. Besides the 15 benchmark corruptions

we also include the saturate corruption from [20] to test
 995 for more color corruptions. The individual corruptions are
 996 shown in Fig. 8 for severity 4 out of 5. Following [20] im-
 997 ages are saved using light JPEG compression, which means
 998 that other corruptions are also lightly corrupted by the JPEG
 999 compression.

1000 7.3. OpticsBench

1001 The recently published OpticsBench [39] provides addi-
 1002 tional blur corruptions obtained from optics, which are ob-
 1003 tained by convolving images with a given 3D kernel. The
 1004 blur kernels (x,y,color) are color dependent and have diverse
 1005 shapes. They are intended as base classes for primary aber-
 1006 rations and complement the blur corruptions from Common
 1007 Corruptions [20]. Fig. 8 shows the different OpticsBench
 1008 blur corruptions in the upper left along with the correspond-
 1009 ing blur kernel in linear scaling. Coma and astigmatism are
 1010 shown in one orientation. As the authors use the defocus
 1011 blur from [20] as a baseline, we also visualize the corre-
 1012 sponding defocus blur kernel in Fig. 8.

1013 7.4. Grouping Corruption Types

1014 Combining the two corrupted datasets, we evaluated all
 1015 models on 24 different corruptions, with each consisting of
 1016 five severities. The high number of corruptions provides
 1017 a broad and stable view on the robustness of our models
 1018 against such corruptions. However, to visualize the results
 1019 of our evaluation, we grouped corruptions in five super-
 1020 categories: noise, blur, compression, weather, and color.

1021 The five super-categories with the corresponding corrup-
 1022 tions can be examined in Fig. 8. As the first super-category,
 1023 we combined all noise generating corruptions from the
 1024 Common Corruptions into *noise*. Furthermore, we grouped
 1025 all blur corruptions from the Common Corruptions [20]
 1026 with all corruptions from OpticsBench [39] into the super-
 1027 category *blur*, as they all mimic different blur types. Similar
 1028 to [20] we combine all weather corruptions into the super-
 1029 category *weather*. However, we shifted the brightness cor-
 1030 ruption to the super-category *color*, as brightness, contrast,
 1031 and saturate behave similar. The last super-category com-
 1032 bines three different image transformation corruptions into
 1033 the super-category *compression*.

1034 Each of the evaluated corruptions and super-categories
 1035 have five severities. In addition to the spatial domain vi-
 1036 sualization in Fig. 8, the super-categories are visualized
 1037 in the frequency domain, in Fig. 9 to 13. Each of these
 1038 figures is divided into six columns, in which the first five

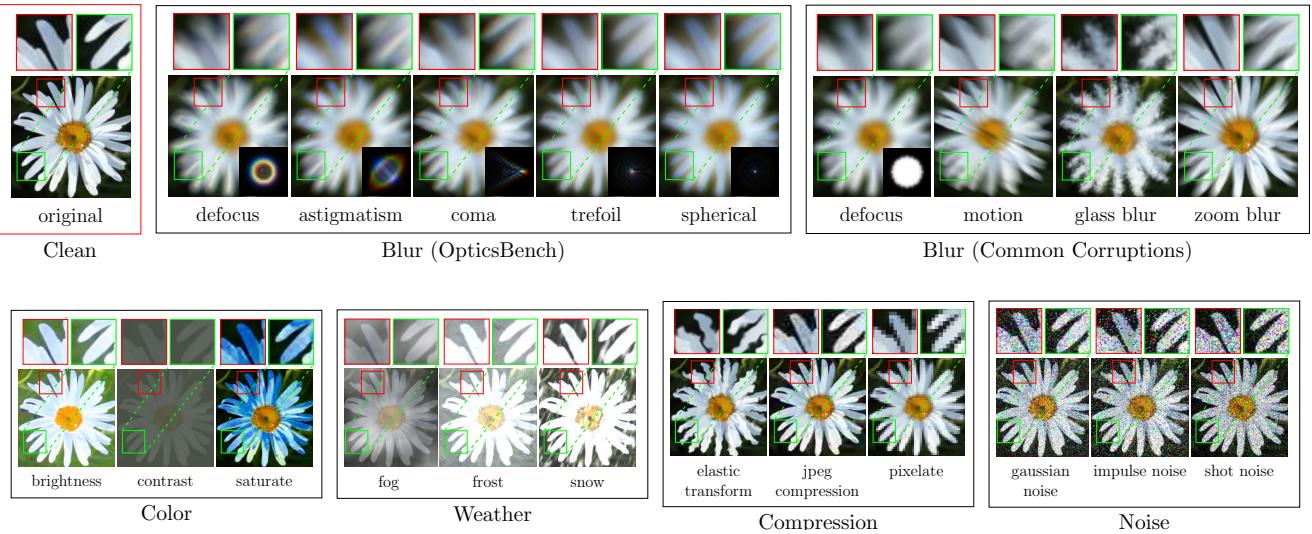


Figure 8. Overview of used corruptions applied to an ImageNet sample. The first row shows the unmodified image (upper left) and various blur corruptions from [39] (left box) and [20] (right box). The second row shows the remaining common corruptions from [20] grouped to color, weather, compression and noise (from left to right). All shown corrupted images represent severity 4 out of 5.

1040 columns represents one of the super-categories. For every
 1041 super-category, the average Fourier transformation magni-
 1042 tude over 100 images is displayed in the top row. The bot-
 1043 tom row corresponds to the absolute difference of corrupted
 1044 average Fourier transformation magnitude against the clean
 1045 data average Fourier transformation magnitude. The last
 1046 column represents the average frequency magnitude of 100
 1047 clean data images.

8. Experiment setup

1048 All models are trained on clean data, without any data aug-
 1049 mentation. To guarantee a comparability between all mod-
 1050 els from the same architecture type, we only added our
 1051 proposed input layer to the corresponding models. There-
 1052 fore, we did not change any hyperparameters in-between the
 1053 training of the same architecture type. The hyperparameters
 1054 for the training are used from [37].

9. Additional Experimental Results

1055 To have a more holistic view on the performance of mod-
 1056 els with our proposed pre-pend layer, multiple additional
 1057 experiment are conducted. This section contains multiple
 1058 detailed evaluation results from models trained on the Im-
 1059 ageNette [27] dataset. The next subsection (Sec. 9.1) in-
 1060 dicates the reasoning for the proposed kernel size. Section
 1061 9.2 presents more experiments conducted with a static fil-
 1062 ter (class II). In Sec. 9.3 and 9.4, more experiments can be
 1063 examined on class III kernels. More ImageNet-100 training
 1064

1065 results are introduced in Sec. 9.5. Additional to the cor-
 1066 ruption robustness, Sec. 9.6 indicates a higher robustness
 1067 against adversarial attack by adding our proposed trainable
 1068 pre-processing layer to existing models.

1069 In Fig. 14 the same pattern, as in Fig. 1, can be in-
 1070 spected. The convolutional preprocessing (ours) increases
 1071 the performance of the AlexNet model on all corruptions
 1072 and severities of the OpticsBench [39] dataset.

1073 An increase of accuracy over the most corruptions from
 1074 Common Corruptions [20] and OpticsBench [39], espe-
 1075 cially on higher severities are displayed in Fig. 16 to 22.

9.1. Kernel Sizes

1076 To explore a suitable kernel size for our proposed layer, we
 1077 conducted six experiments with slightly increasing kernel
 1078 sizes. The results over the clean data, the Common Corrup-
 1079 tion corrupted data, and the OpticsBench corrupted data, is
 1080 displayed in Table 5. On the clean data, the baseline, with-
 1081 out any prior layer to the mode, outperforms all kernel sizes.
 1082 Furthermore, the models with a prior trainable layer and a
 1083 small kernel size (7×7 and 15×15) achieve a higher accu-
 1084 racy on the clean data, than models with a large kernel size,
 1085 such as 31×31 and 35×35 . However, on the corrupted
 1086 datasets (OpticsBench and Common Corruptions), the mod-
 1087 els with a larger kernel size outperform models with smaller
 1088 kernel sizes and the baseline significantly. The best per-
 1089 formance on both corrupted datasets and a just slightly worse
 1090 performance on the clean dataset is achieved by the model
 1091

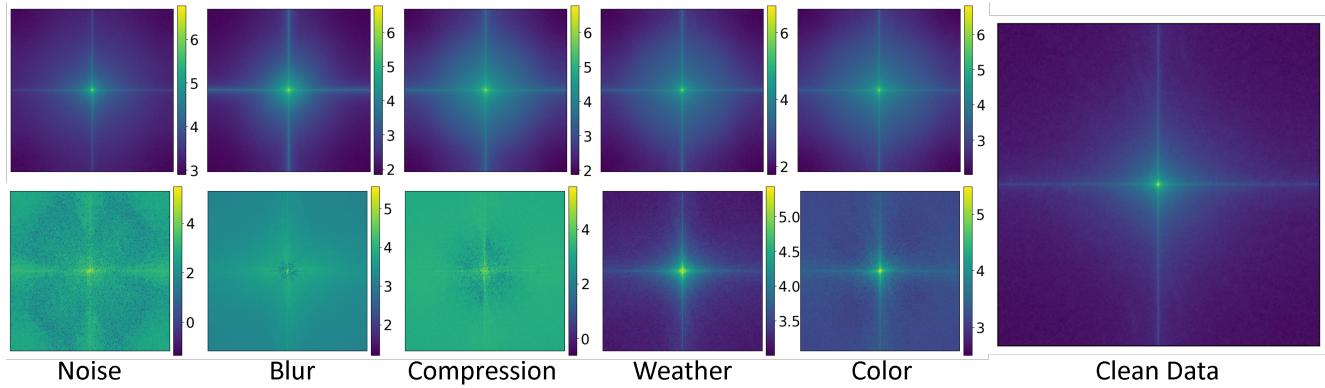


Figure 9. Corruption Severity 1 - Average Fourier transformation magnitude over 100 ImageNet images for each super-category (column 1–5) and the clean data (column 6). Top: The average Fourier transformation magnitude. Bottom: The difference between the corrupted average Fourier transformation magnitude and the clean average Fourier transformation magnitude.

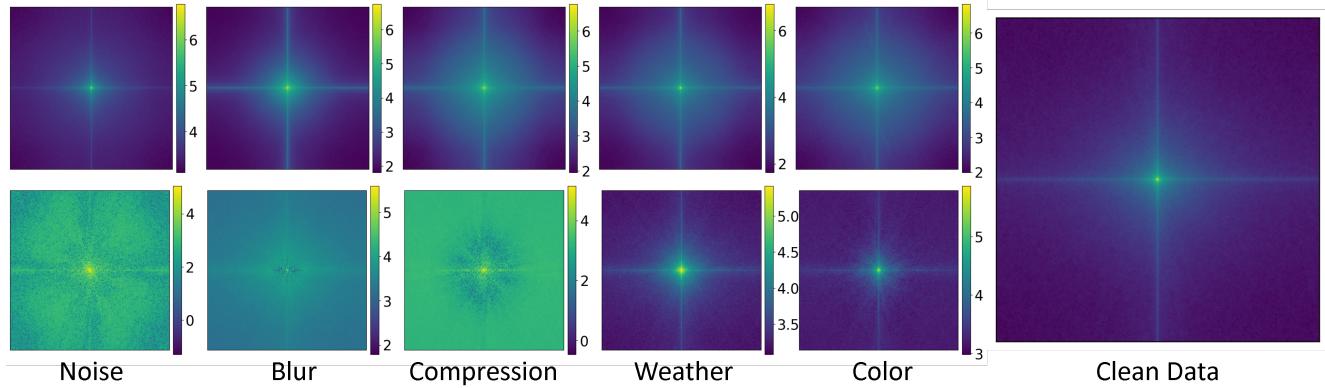


Figure 10. Corruption Severity 2 - Average Fourier transformation magnitude over 100 ImageNet images for each super-category (column 1–5) and the clean data (column 6). Top: The average Fourier transformation magnitude. Bottom: The difference between the corrupted average Fourier transformation magnitude and the clean average Fourier transformation magnitude.

1094 with the kernel size of 25×25 . Therefore, we used this ker-
1095 nel size for all other experiments, if not stated differently.
1096

1097 9.2. Static Kernels (Class II)

1098 To investigate the extent to which kernel training improves
1099 prediction stability, we freeze different convolutional layer
1100 initializations. These static filters are then compared to both
1101 the baseline and the fully trainable filters. We investigate
1102 two blur filters, one rotational-symmetric Gauss blur filter
1103 and a directional (horizontal coma) blur filter obtained from
1104 OpticsBench [39]. These have both lowpass characteristic
1105 and remove high frequency content. The color distortion
1106 filter aims to split color-specific information by translation.
1107 The results on clean data and the different corruption bench-
1108 marks on ImageNette are shown in Table 6. It is noticeable
1109 that the two blur filters have an in-domain accuracy compa-
1110 rable to the fully trainable or the baseline model, while the

Kernel size	CD	OB [39]	CC [20]
Base	*0.800	*0.592	*0.487
7	0.782	0.616	0.526
15	0.790	0.575	0.532
21	0.765	0.610	0.530
25	*0.775	*0.685	*0.565
31	0.707	0.681	0.556
35	0.713	0.678	0.549

Table 5. Results on ImageNette for ResNet50 and different kernel size for the proposed trainable convolutional preprocessing layer. CD = Clean Data, OB = OpticsBench, CC = Common corruptions. * = average from multiple seeds. The results on the two corruption benchmarks are averaged across severity and corruption

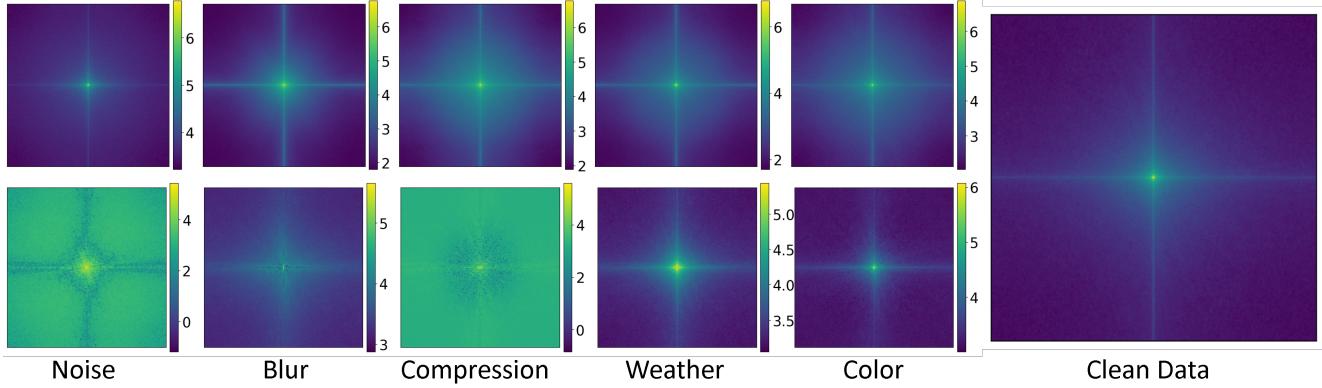


Figure 11. Corruption Severity 3 - Average Fourier transformation magnitude over 100 ImageNet images for each super-category (column 1-5) and the clean data (column 6). Top: The average Fourier transformation magnitude. Bottom: The difference between the corrupted average Fourier transformation magnitude and the clean average Fourier transformation magnitude.

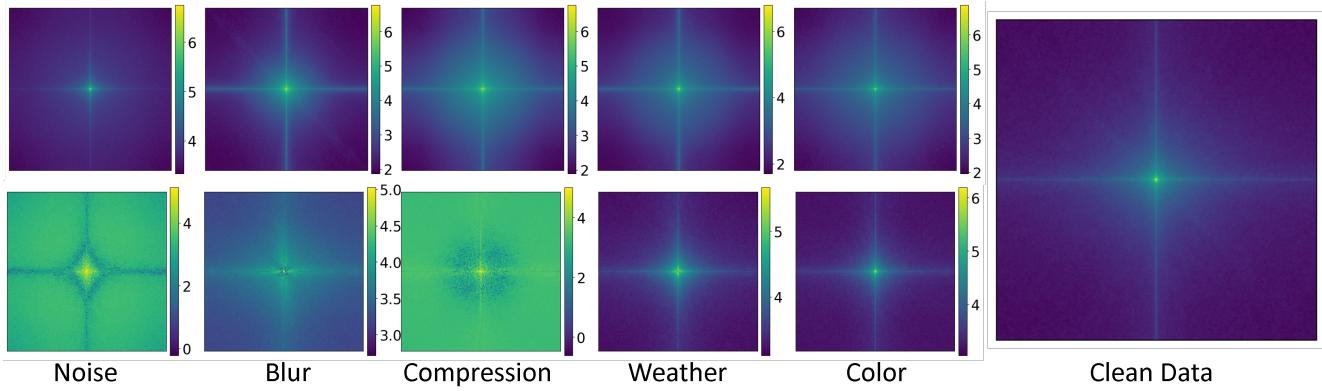


Figure 12. Corruption Severity 4 - Average Fourier transformation magnitude over 100 ImageNet images for each super-category (column 1-5) and the clean data (column 6). Top: The average Fourier transformation magnitude. Bottom: The difference between the corrupted average Fourier transformation magnitude and the clean average Fourier transformation magnitude.

color distortion kernel model is significantly lower. Overall, only the Gaussian blur and color distortion models perform better than the baseline on the two corruption benchmarks. The directional blur model performs worst on corruptions, but second best on clean data.

9.3. Detailed plots of Thresholded Kernel Performance

Fig. 7 shows, the nearly constant performance on the OpticsBench dataset, while reducing the frequencies, with small coefficients, in the trainable kernel. To show this effect in a little more detail, Fig. 23 presents the evaluation results of only one training run. To illustrate the model performance of class III kernels, we trained on MobileNet v3 large model with our proposed convolutional preprocessing layer and no restrictions, such as L_1 prior. Subsequent to the training, we evaluated the model multiple times with different shares of removed frequencies of the proposed kernel

Kernel type	class	CD	OB [39]	CC [20]
None (Base)	-	*0.800	*0.592	*0.487
Preserve content	I	0.754	0.476	0.513
Preserve content large	I	0.645	0.478	0.440
Fully trainable	II	*0.775	*0.685	*0.565
L_1 prior	III	0.712	0.699	0.567
Random initialization	II	0.711	0.673	0.550
Con2D KS=25	II	0.655	0.601	0.501
Directional blur filter	II	0.778	0.437	0.345
Gauss blur filter	II	0.764	0.668	0.541
Color distortion	II	0.677	0.660	0.533

Table 6. Results on ImageNette for ResNet50 and different input layer large kernel types. CD = Clean Data, OB = OpticsBench, CC = Common corruptions. * = average from multiple seeds. The results on the two corruption benchmarks are averaged across severity and corruption.

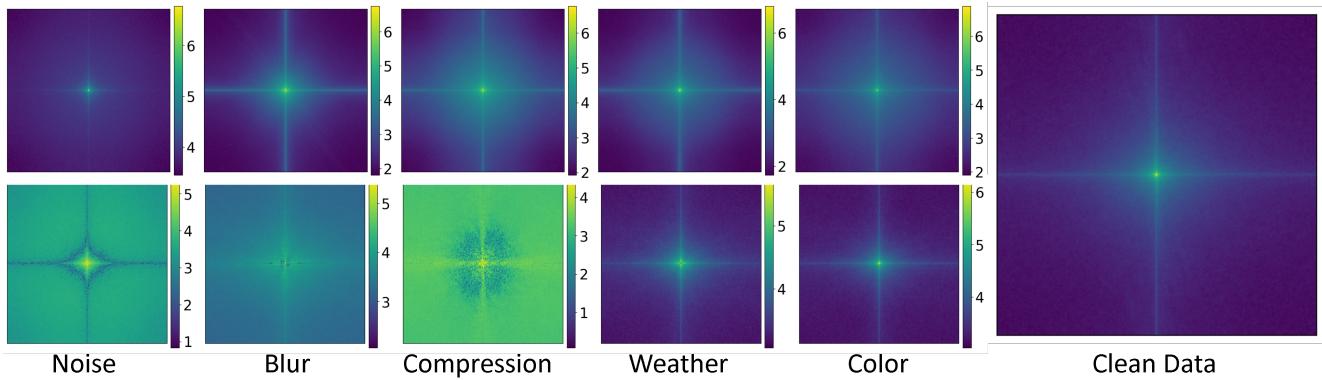


Figure 13. Corruption Severity 5 - Average Fourier transformation magnitude over 100 ImageNet images for each super-category (column 1-5) and the clean data (column 6). Top: The average Fourier transformation magnitude. Bottom: The difference between the corrupted average Fourier transformation magnitude and the clean average Fourier transformation magnitude.

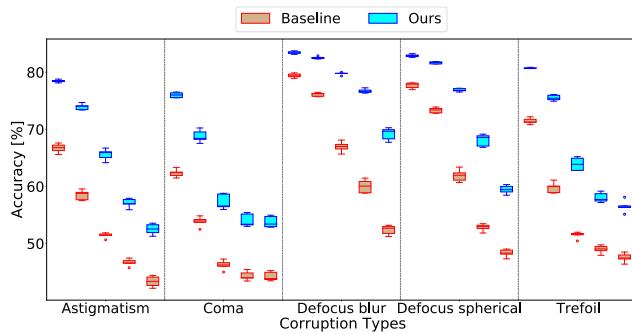


Figure 14. Convolutional preprocessing (ours) can increase the robustness of classification networks against unknown corruptions *without data augmentation*. AlexNet improved with a trainable preprocessing filter evaluated on ImageNette [27] blur and corruptions from OpticsBench [39]. For each corruption type, five levels of severity are shown from left to right. The variation, visualized via the box plots, results from five different seeds per model.

layer. In Fig. 23 the evaluation runs with a lower frequency removal share (0% & 90%) performs better for low severities, but decreases in a faster manner, than evaluation runs with a high frequency removal share (95% & 99%).

9.4. L1 prior different regularization strengths

Sec. 9.3 compares already trained class II kernels, which are adjusted to class III kernels, by removing frequencies from the proposed layer kernel. In this section, we further investigate the different class III kernels. Therefore, multiple ResNet50 models with our proposed fully trainable convolutional layer are trained with an L_1 prior on the frequency coefficients on the trainable kernel of our proposed layer. Each training run uses a different λ to shift the loss more on the frequency coefficients on the trainable kernel. Subsequently, these models are evaluated on the

OpticsBench [39] and Common Corruptions [20], which results in Fig. 24. This experiment shows similar results as in Fig. 23, as the results with a lower lambda are decreasing over higher severities faster. Especially, for noise and blur corruptions, this decreasing over severities effect is more significant. However, it seems like models, which are trained and optimized to also minimize the frequency coefficients on the trainable kernel, tend to have a higher overall accuracy for small λ s ($\lambda = 0.00001$ to $\lambda = 0.005$), while models with larger λ s ($\lambda = 0.001$ to $\lambda = 0.01$), seem to perform more constant over all severities of corruptions.

9.5. ImageNet-100 and ImageNet-1k

Additionally, to all the experiments on ImageNette [27], we also tested the scalability of the proposed convolutional prior layer. Therefore, Fig. 25 to 32 show comparisons of four different models (ResNet, AlexNet, EfficientNet, and MobileNet) on ImageNet-100 [55]. In each figure, we compare a model with our proposed fully trainable convolutional input layer against the baseline. The figures show the performance on the ImageNet-100 corruptions of OpticsBench [39] and Common Corruptions [20]. All the comparisons on ImageNet-100 show, that our proposed model outperforms the baseline on the evaluated corruptions. However, the performance gain is slightly lower than on ImageNette.

While training on the whole ImageNet-1k dataset [46], the performance differences between the baseline and our proposed model is even lower than on ImageNet-100. Only for severity 4 and 5 our model outperforms the baseline. This effect is visualized in Fig. 33 and 34. However, the performance on clean data is also just slightly in favor of the baseline, which is stated in Table 7.

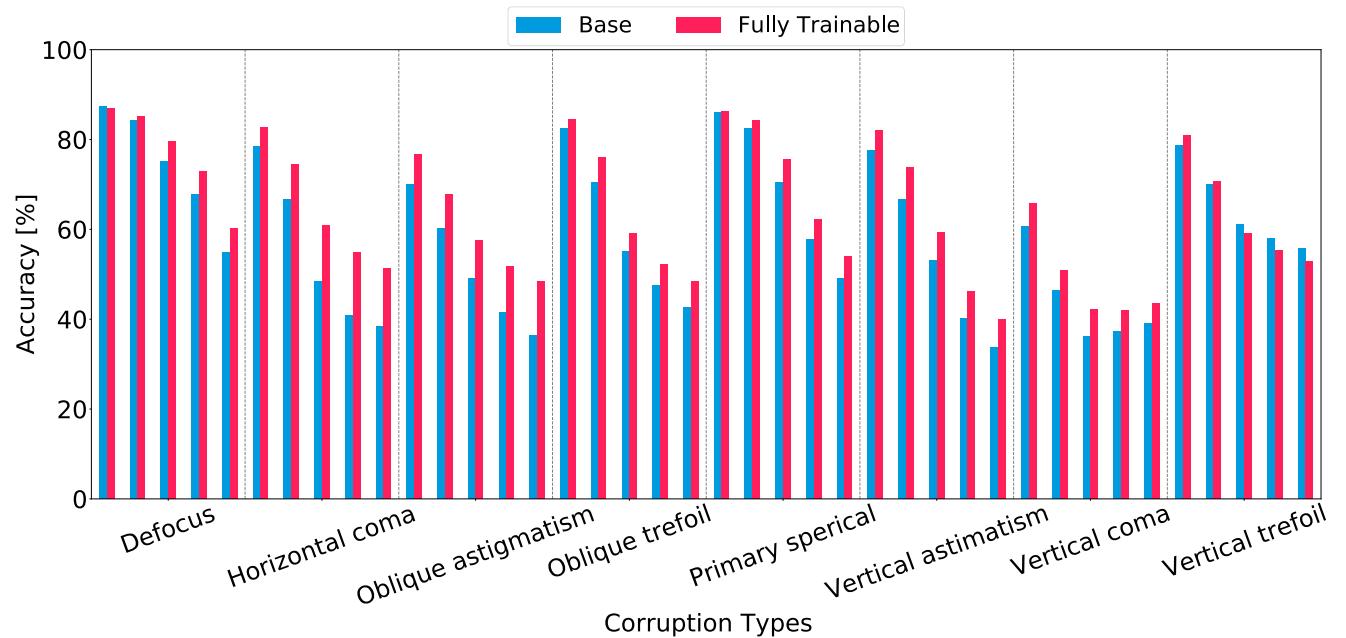


Figure 15. MobileNet v3 large - ImageNette OpticsBench. A comparison of the MobileNet v3 large (Base) and a MobileNet v3 large with our proposed layer (Fully Trainable), on ImageNette OpticsBench.

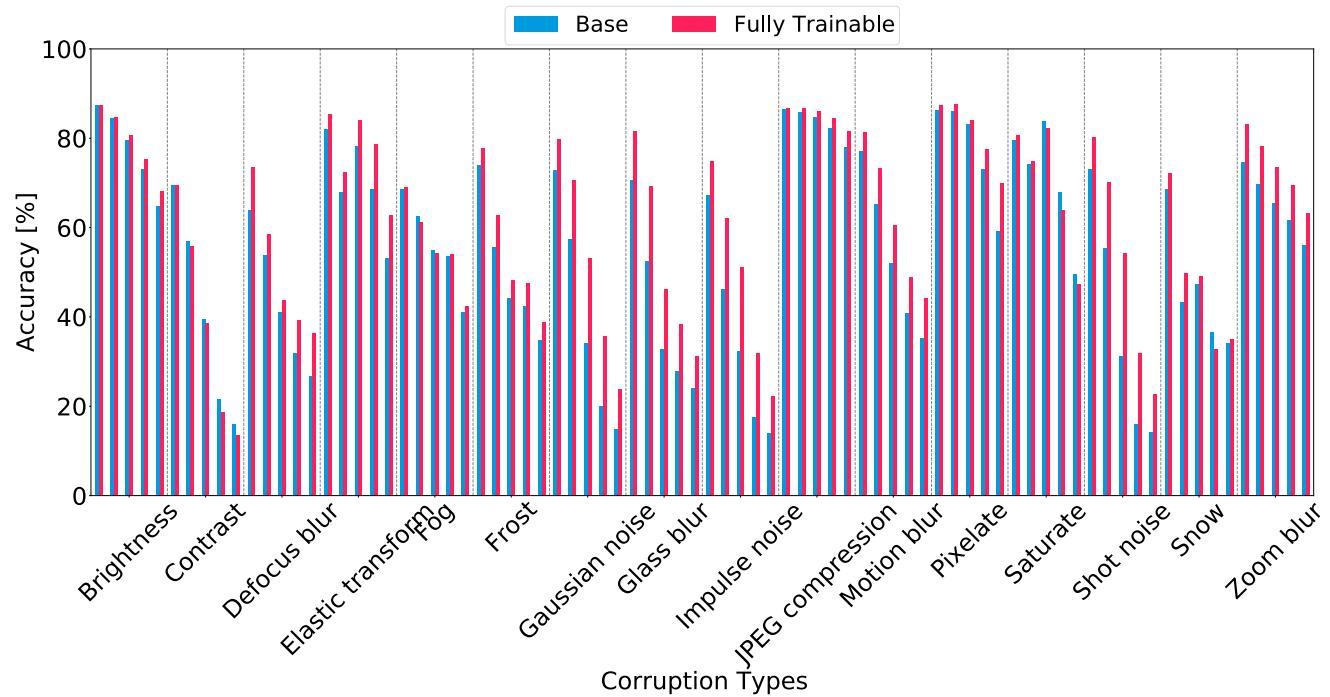


Figure 16. MobileNet v3 large - ImageNette Common Corruptions. A comparison of the MobileNet v3 large (Base) and a MobileNet v3 large with our proposed layer (Fully Trainable), on ImageNette Common Corruptions.

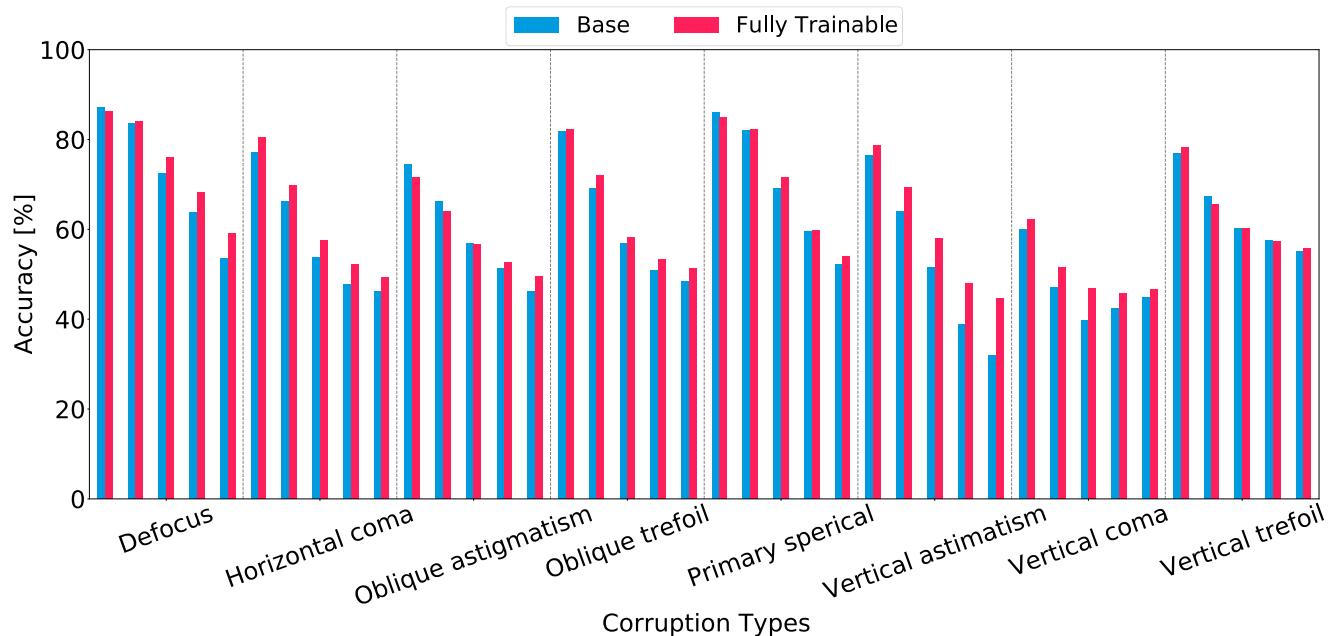


Figure 17. EfficientNet b0 - ImageNette OpticsBench. A comparison of the EfficientNet b0 (Base) and a EfficientNet b0 with our proposed layer (Fully Trainable), on ImageNette OpticsBench.

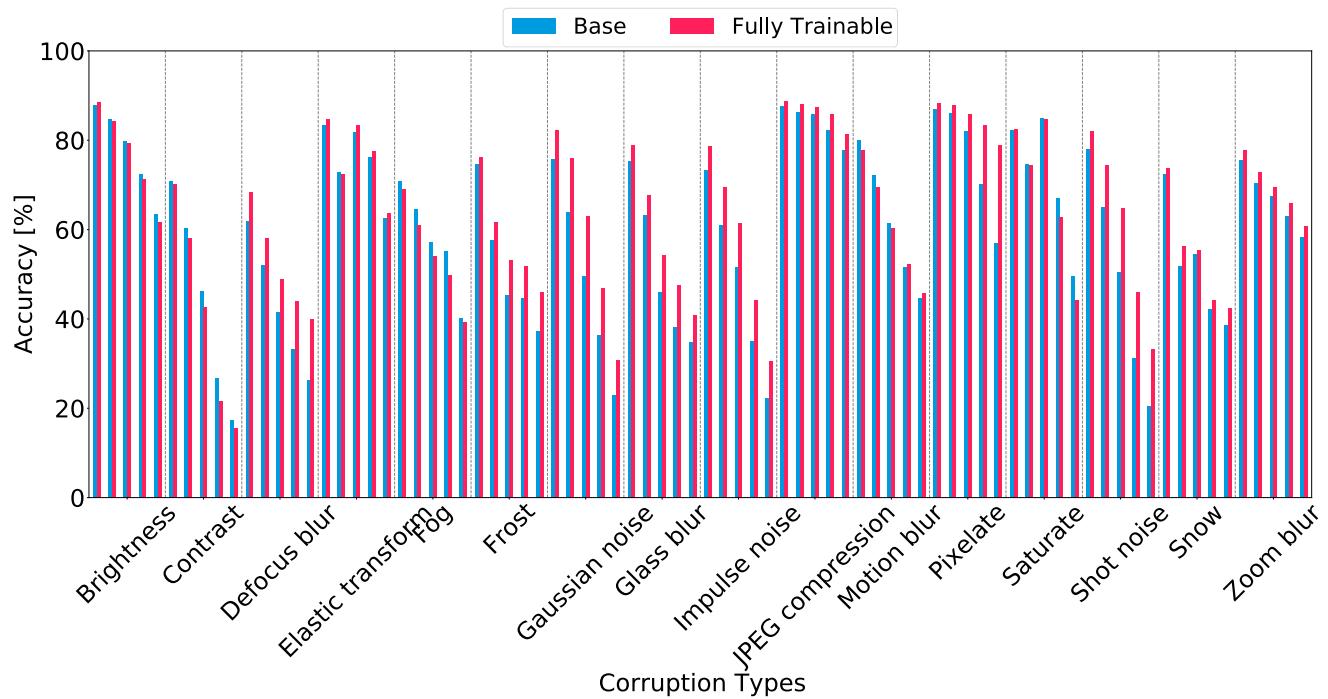


Figure 18. EfficientNet b0 - ImageNette Common Corruptions. A comparison of the EfficientNet b0 (Base) and a EfficientNet b0 with our proposed layer (Fully Trainable), on ImageNette Common Corruptions.

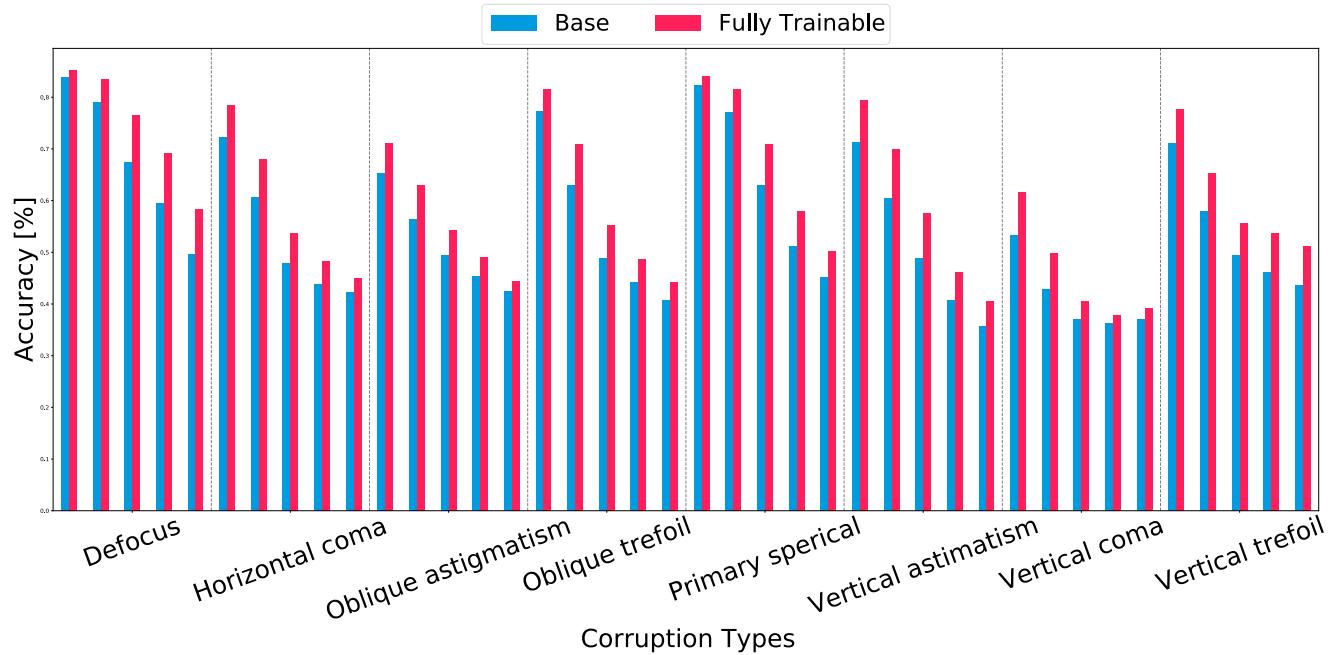


Figure 19. DenseNet161 - ImageNette OpticsBench. A comparison of the DenseNet161 (Base) and a DenseNet161 with our proposed layer (Fully Trainable), on ImageNette OpticsBench.

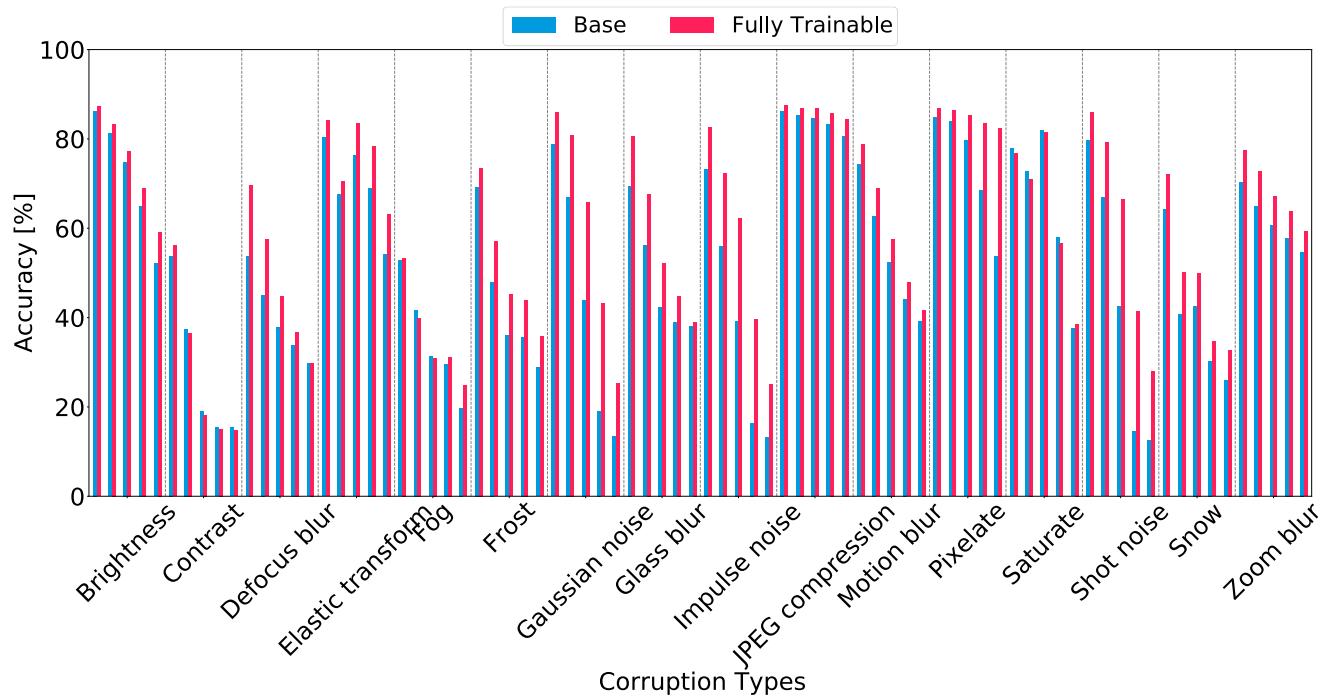


Figure 20. DenseNet161 - ImageNette Common Corruptions. A comparison of the DenseNet161 (Base) and a DenseNet161 with our proposed layer (Fully Trainable), on ImageNette Common Corruptions.

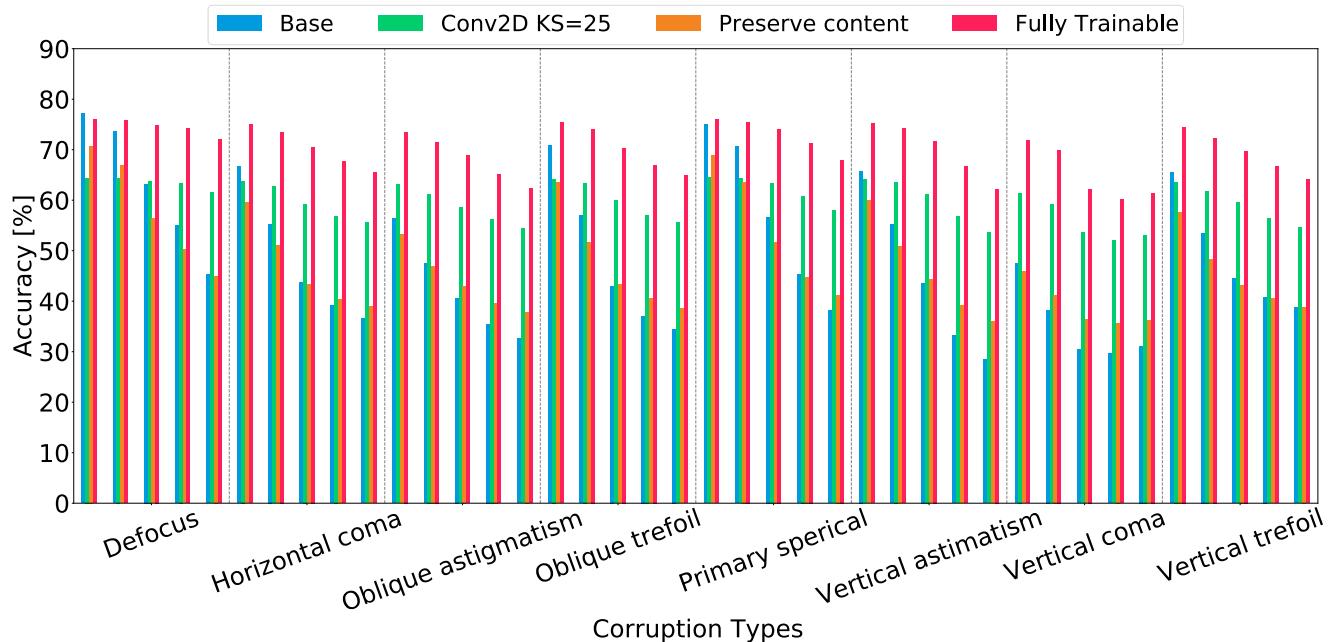


Figure 21. ResNet50 - ImageNette OpticsBench. A comparison of the ResNet50 (Base), the ResNet50 with an additional prior 2D convolutional layer (Conv2D KS=25), a ResNet50 with an additional trainable class I layer (Preserve Content), and a ResNet50 with our proposed layer (Fully Trainable), on ImageNette OpticsBench.

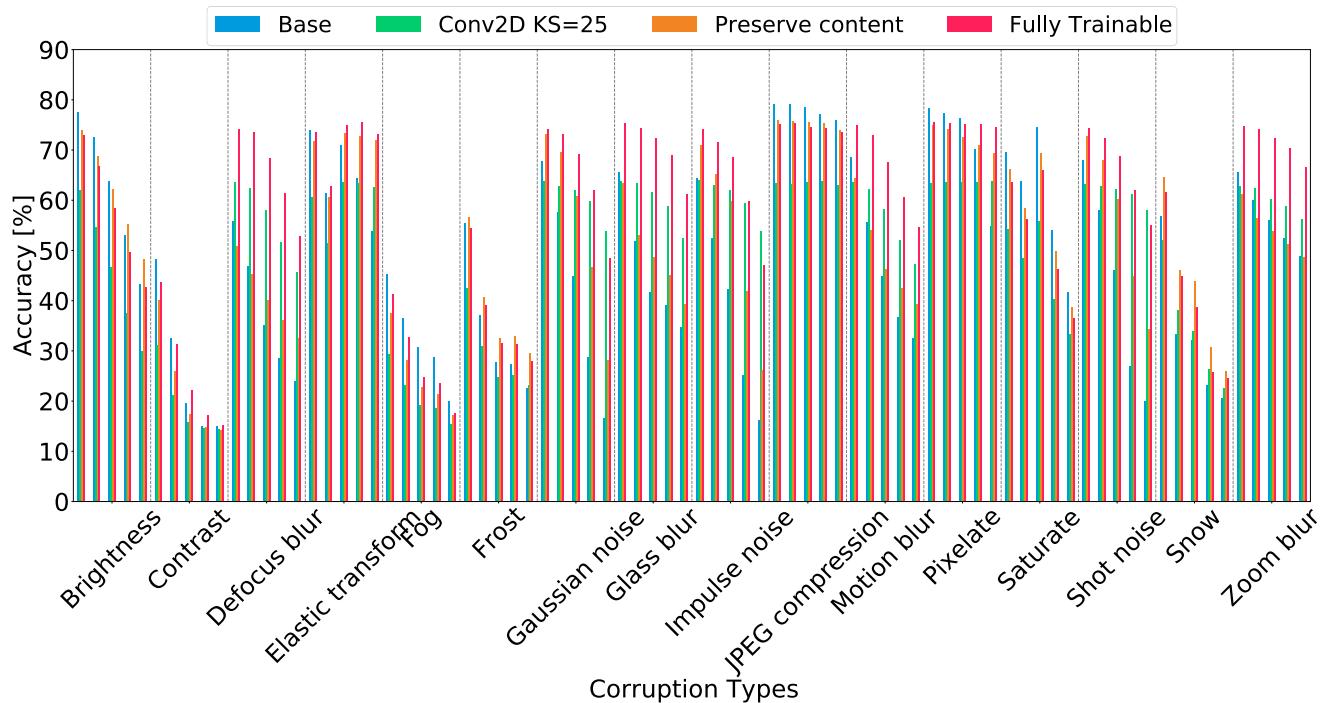


Figure 22. ResNet50 - ImageNette Common Corruptions. A comparison of the ResNet50 (Base), the ResNet50 with an additional prior 2D convolutional layer (Conv2D KS=25), a ResNet50 with an additional trainable class I layer (Preserve Content), and a ResNet50 with our proposed layer (Fully Trainable), on ImageNette Common Corruptions.

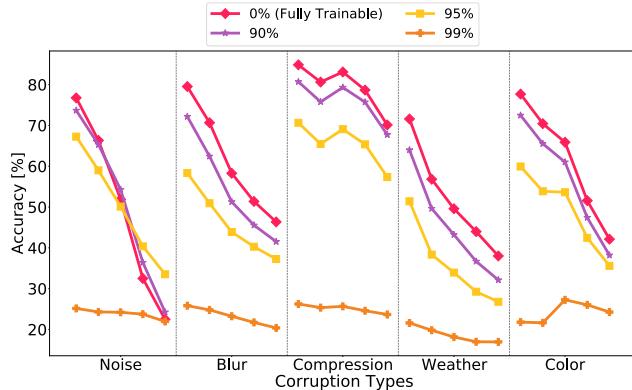


Figure 23. Comparison of frequency thresholded kernels on OpticsBench [39] and Common Corruptions [20]. All evaluation results from this diagram are from the same MobileNet v3 large trainings run, with different thresholding intervals.

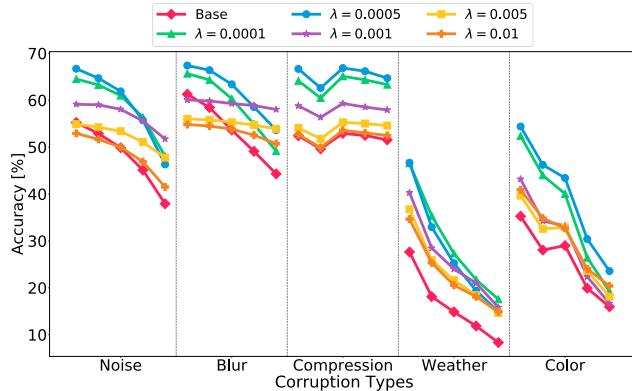


Figure 24. Comparison of models, which are trained with an L_1 prior. Each model has an ResNet50 architecture, with our additional proposed fully trainable convolutional layer. All these models are trained on ImageNette [27] and evaluated on corruptions from OpticsBench [39] and Common Corruptions [20].

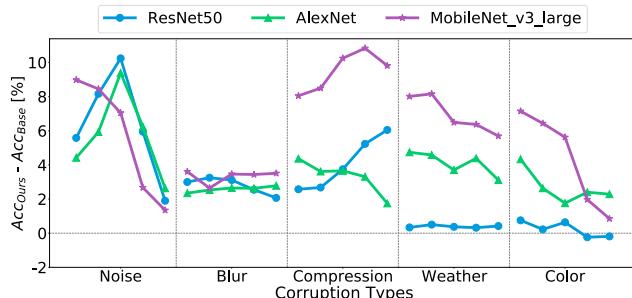


Figure 25. Improvement on ImageNet-100 OpticsBench and Common Corruptions. For each corruption, five levels of severity are shown from left to right. For readability, we summarize different corruption types and take the average.

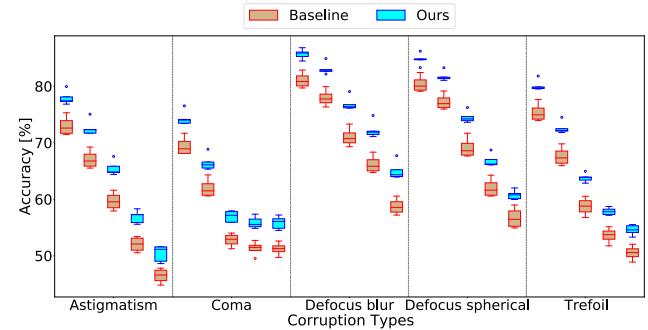


Figure 26. Convolutional preprocessing (ours) can increase the robustness of classification networks against unknown corruptions *without data augmentation*. ResNet50 improved with a trainable preprocessing filter evaluated on ImageNet-100 [55] blur and corruptions from OpticsBench [39]. For each corruption type, five levels of severity are shown from left to right. The variation, visualized via the box plots, results from five different seeds per model.

Model	Version	CD	OB	CC
ResNet50	Base	0.522	0.221	0.195
ResNet50	Fully trainable	0.503	0.226	0.197
MobileNet	Base	0.646	0.292	0.240
MobileNet	Fully trainable	0.639	0.296	0.233

Table 7. Results on ImageNet-1k for ResNet50 and MobileNet v3 large. CD=Clean Data, OB = OpricsBench, CC=Common Corruptions. The results on the two corruption benchmarks are averaged across severity and corruption.

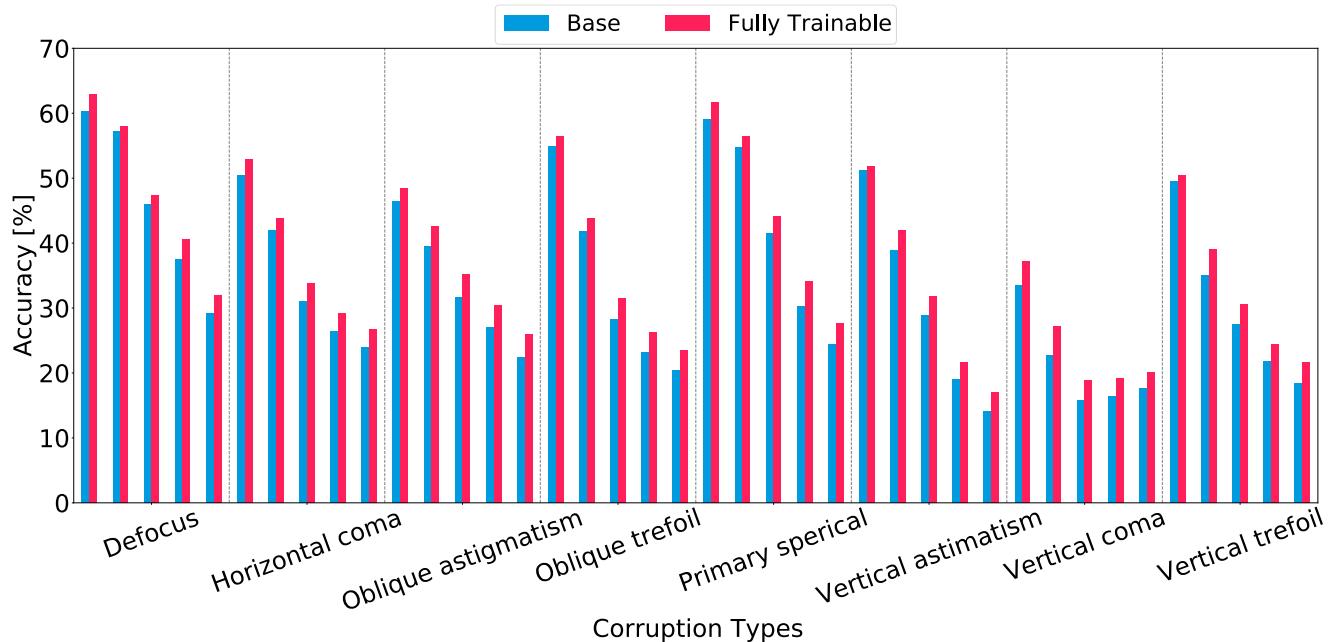


Figure 27. AlexNet - ImageNet-100 OpticsBench. A comparison of the AlexNet (Base) and a AlexNet with our proposed layer (Fully Trainable), on ImageNet-100 OpticsBench.

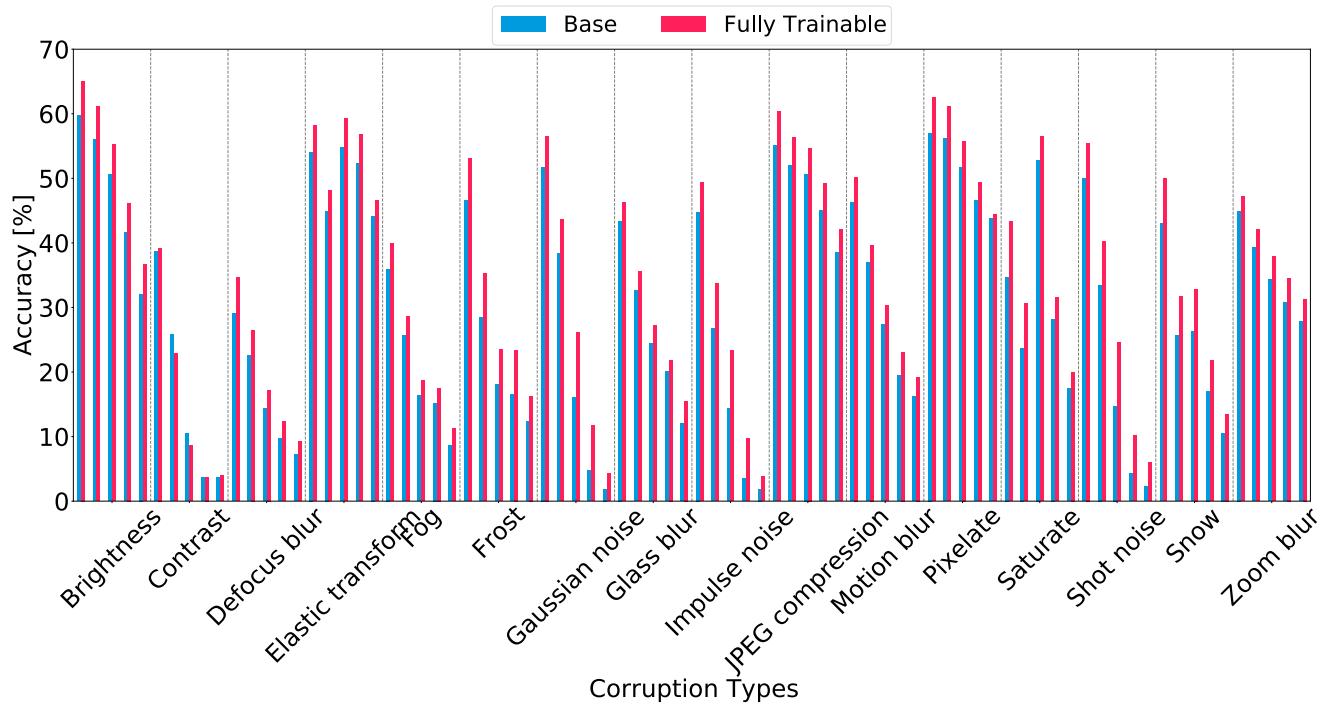


Figure 28. AlexNet - ImageNet-100 Common Corruptions. A comparison of the AlexNet (Base) and a AlexNet with our proposed layer (Fully Trainable), on ImageNet-100 Common Corruptions.

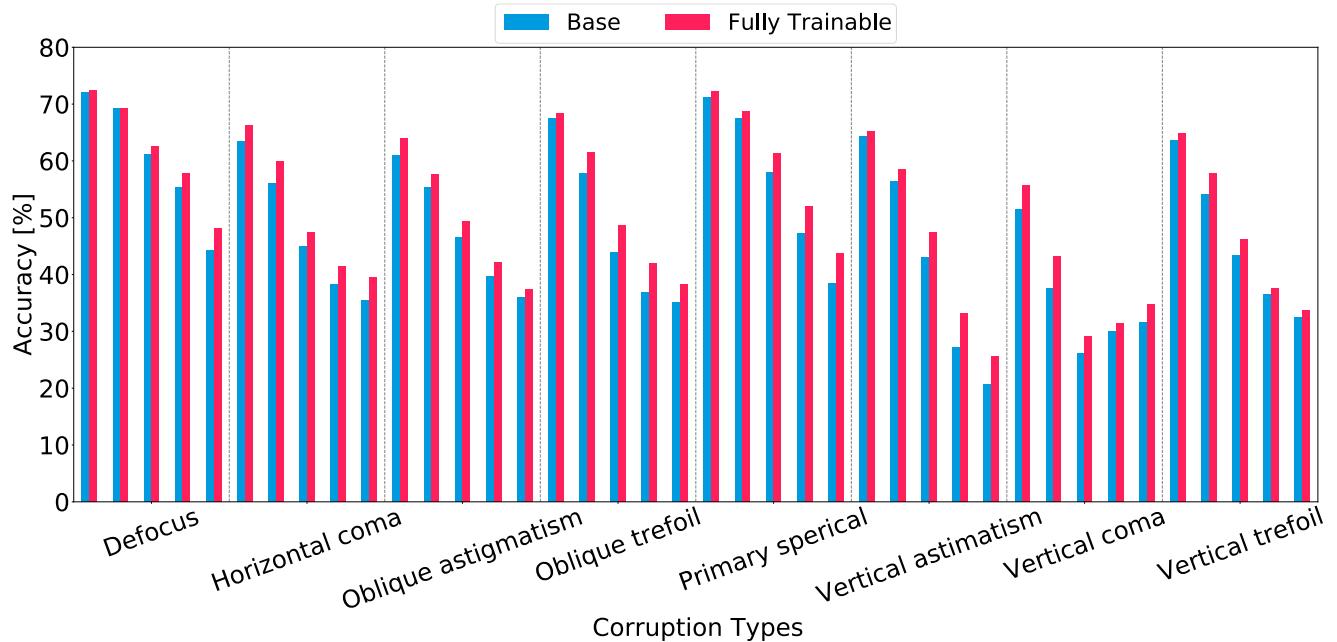


Figure 29. EfficientNet b0 - ImageNet-100 OpticsBench. A comparison of the EfficientNet b0 (Base) and a EfficientNet b0 with our proposed layer (Fully Trainable), on ImageNet-100 OpticsBench.

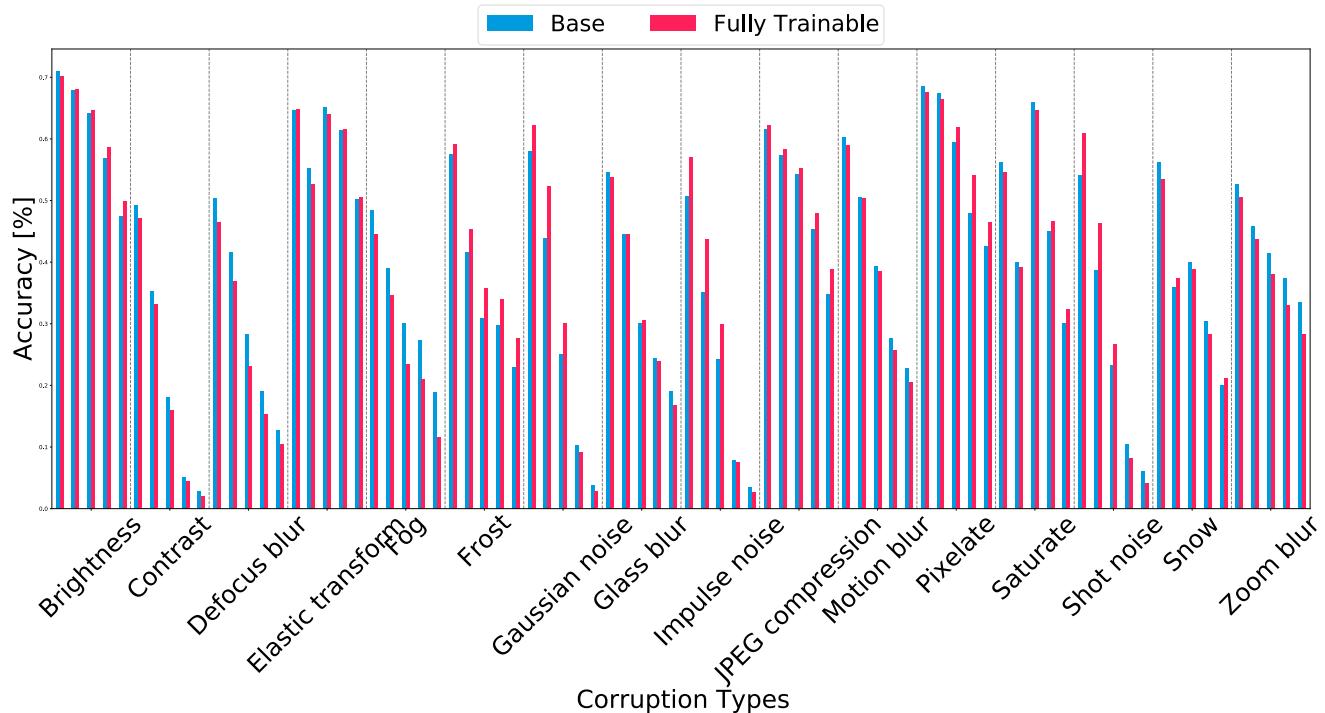


Figure 30. EfficientNet b0 - ImageNet-100 Common Corruptions. A comparison of the EfficientNet b0 (Base) and a EfficientNet b0 with our proposed layer (Fully Trainable), on ImageNet-100 Common Corruptions.

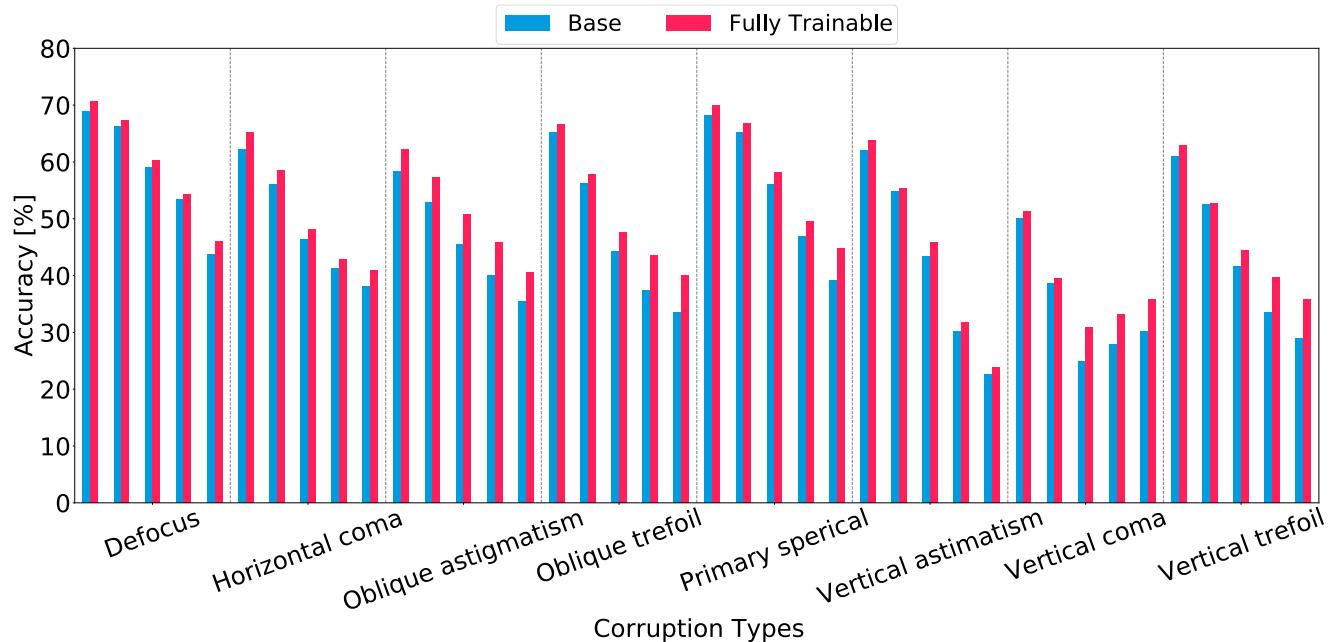


Figure 31. MobileNet v3 large - ImageNet-100 OpticsBench. A comparison of the MobileNet v3 large (Base) and a MobileNet v3 large with our proposed layer (Fully Trainable), on ImageNet-100 OpticsBench.

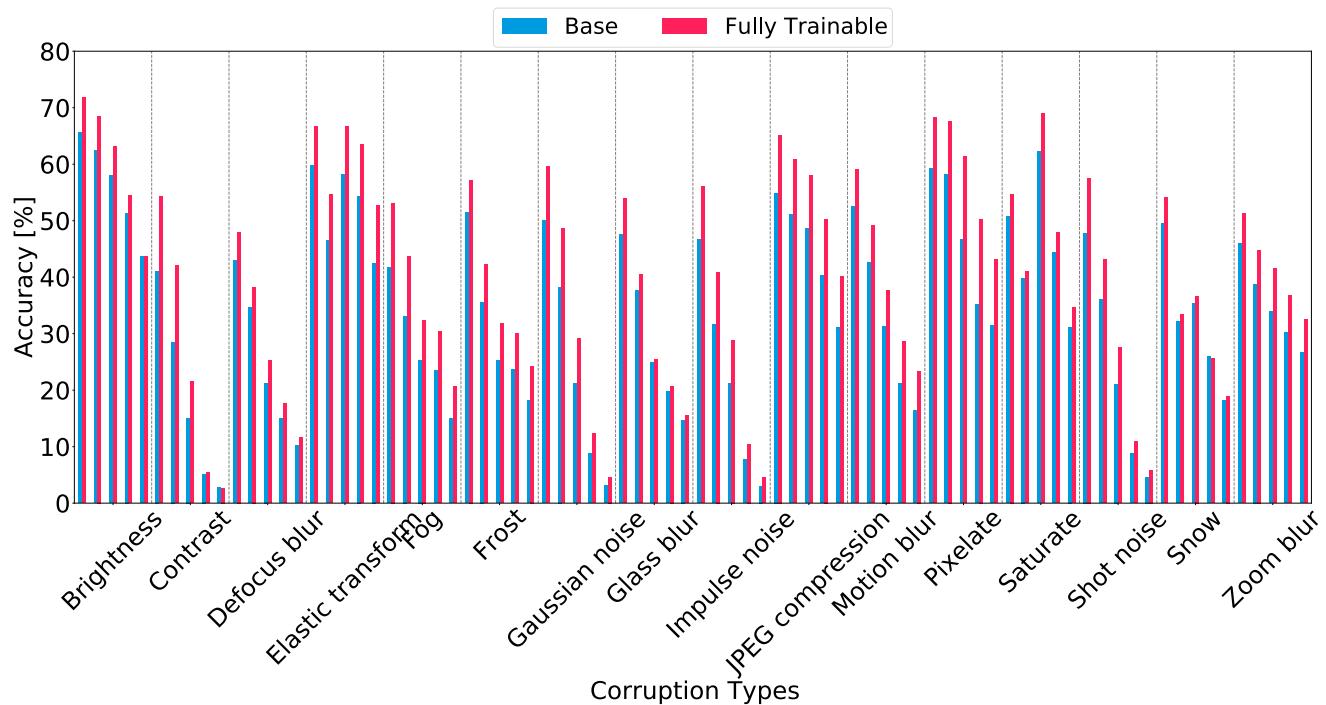


Figure 32. MobileNet v3 large - ImageNet-100 Common Corruptions. A comparison of the MobileNet v3 large (Base) and a MobileNet v3 large with our proposed layer (Fully Trainable), on ImageNet-100 Common Corruptions.

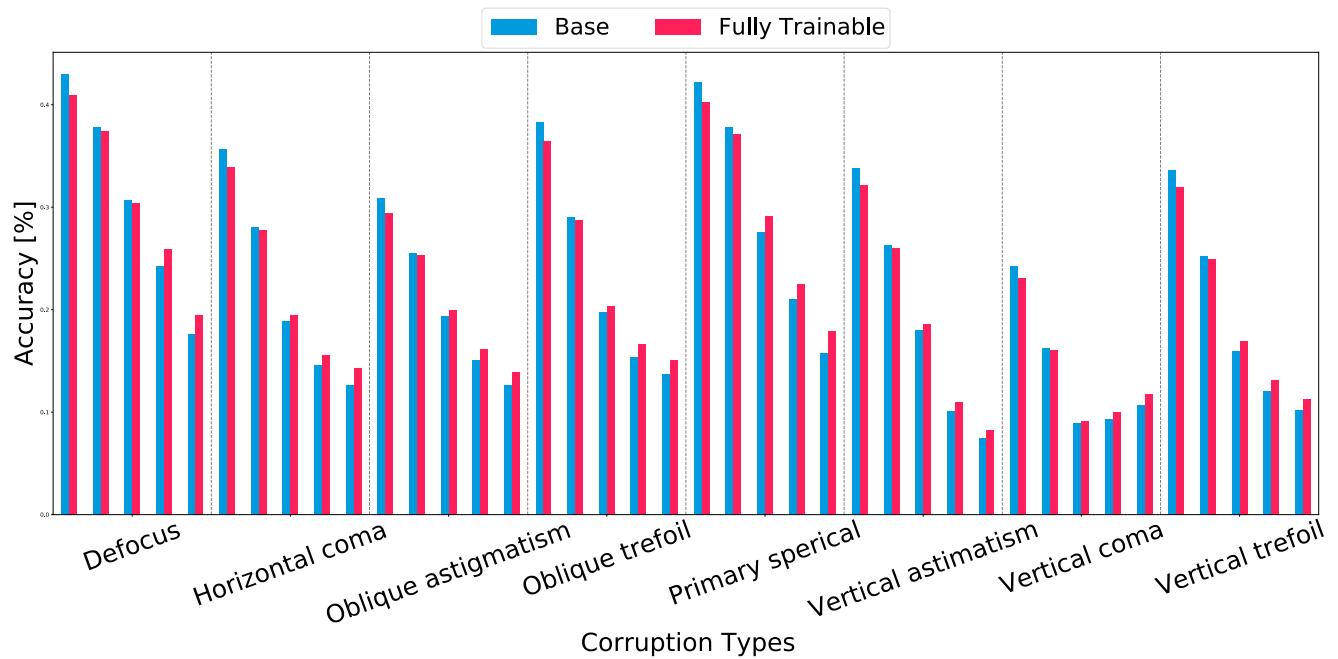


Figure 33. ResNet50 - ImageNet-1k OpticsBench. A comparison of the ResNet50 (Base) and a ResNet50 with our proposed layer (Fully Trainable), on ImageNet-1k OpticsBench.

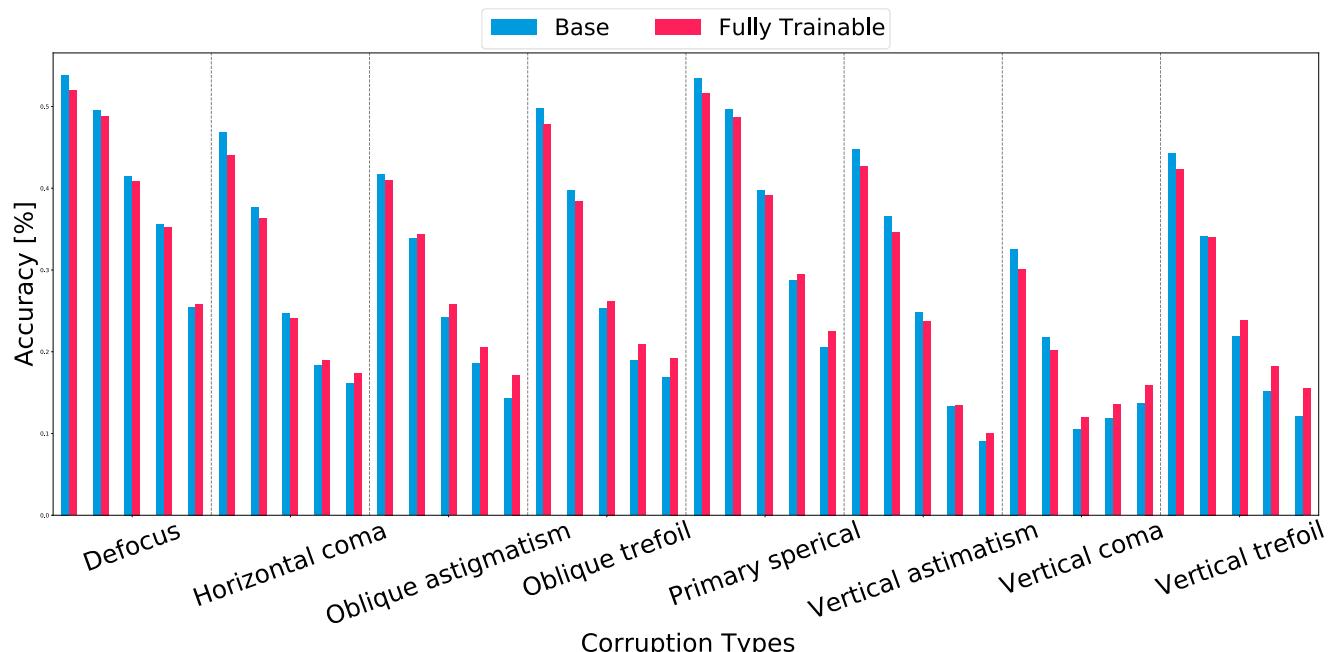


Figure 34. MobileNet v3 large - ImageNet-1k OpticsBench. A comparison of the MobileNet v3 large (Base) and a MobileNet v3 large with our proposed layer (Fully Trainable), on ImageNet-1k OpticsBench.

1175 **9.6. Adversarial training and robustness**

1176 We further investigate the performance of the proposed input
 1177 layer on white-box adversarial attacks. Therefore, we
 1178 trained a ResNet50 with and without our proposed trainable
 1179 convolutional input layer. The models are trained with the
 1180 fast gradient sign method (FGSM) [15] and an ϵ of $\frac{2}{255}$ and
 1181 $\frac{6}{255}$ and evaluated on clean data and under the FGSM ad-
 versarial attack, with two different ϵ of $\frac{2}{255}$ and $\frac{6}{255}$.

Model	Version	CD	FGSM $\epsilon = \frac{2}{255}$	FGSM $\epsilon = \frac{6}{255}$
ResNet50	Adv. Baseline	0.665	0.567	0.385
ResNet50	Adv. Trainable	0.776	0.667	0.389
ResNet50	Adv. Baseline	0.534	0.483	0.760
ResNet50	Adv. Trainable	0.581	0.527	0.776

Table 8. Adversarial Results on ImageNette [27] data. CD= Clean Data, FGSM = FGSM attack [15] with two different ϵ of $\frac{2}{255}$ and $\frac{6}{255}$. The two models in the top rows are trained with an ϵ of $\frac{2}{255}$, while the two model on the bottom are trained with an ϵ of $\frac{6}{255}$.

In Table 8 the results of these experiments, indicate the robustness of the models with our preprocessing layer: they are outperforming the baseline on all evaluated data. While both models have similar accuracy under FGSM attacks with ϵ of $\frac{2}{255}$, our model achieves better results on clean data (+11.1% & +4.7%) and with a lower ϵ value of $\frac{2}{255}$ (+10.0% & +4.4%). These increases are possible while only adding one layer with less than 2k of parameters.

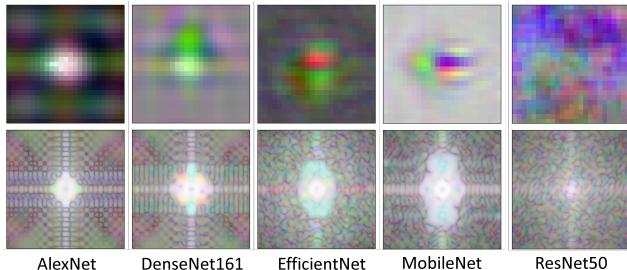
1191 **10. Analysis of optimized kernels**

Figure 35. Optimized kernels in the spatial domain (top) and their image-sized frequency magnitudes in \log_{10} scale (bottom) trained on ImageNette.

In this section, we perform an in-depth analysis of the pre-pended kernels. Fig. 35 visualizes optimized class II kernels and their corresponding frequency magnitudes. It can be observed that the optimized kernels for different neural networks are different. In the following, we quantitatively assess the similarities of these kernels, compare the evolution of their spectra, and discuss their condition numbers in finer detail.

1200 **10.1. Spectrum correlation**

To ensure that our results reflect performance gains due to the proposed method and not because of fortuitous initialization, we experiment with multiple seeds. The kernel is, without change, initialized with a specific coma filter, and the weights of the remaining network are generated pseudo-randomly. We find that the performance of our networks is unaffected. A note-worthy discovery, however, is that despite different initializations, the optimized kernels' spectra are found to be very similar, given the neural network architecture and the input dataset remain the same. Fig. 36 numerically assesses the structural similarity of the spectra of multiple kernels post-optimization. One reason for the observed similarity could be that the proposed kernel learns to map the input to an information-dense subspace most suited to aid the following network in its classification task. Given that the input dataset and the architecture remain invariant, the information-dense subspace should not change as well. Another observation is that given the same input dataset, changing the network changes the mapping learned by the kernel. This suggests that different neural networks predominantly make use of different information content of the unadulterated dataset.

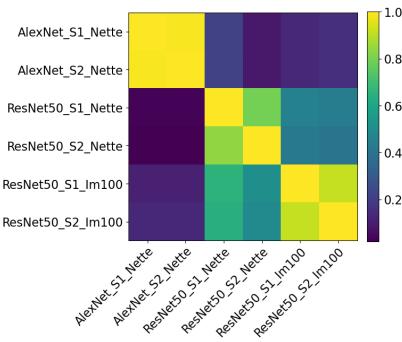


Figure 36. Structural similarity of the magnitudes of the spectra of different kernels. S1 and S2 refer to different seed values, and Im100 and Nette refer to ImageNet-100 and ImageNette, respectively.

1223 **10.2. Spectrum evolution**

We initialize our kernels with a specific coma filter, which can be observed in Fig. 38 (a) and (d). Moreover, these kernels do not mix channel information. We can observe the evolution of the kernels in the spatial and frequency domains in Fig. 38, when the networks are trained on ImageNet-100. From the figure, it is also apparent that the kernels learn different projections for each color channel.

Fig. 39 compares the evolution of class II and class III kernels. The frequency band axis represents the dimension of a square window, with its center being the center of the frequency-shifted spectrum. Each bar indicates the average

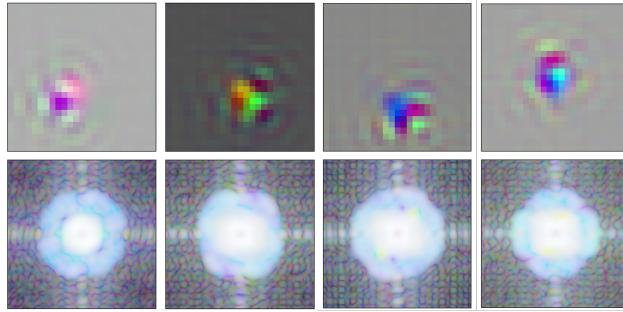


Figure 37. Optimized kernels in the spatial domain (top) and their image-sized frequency magnitudes in \log_{10} scale (bottom) prepended to ResNet50 initialized with different seeds and trained on ImageNet-100.

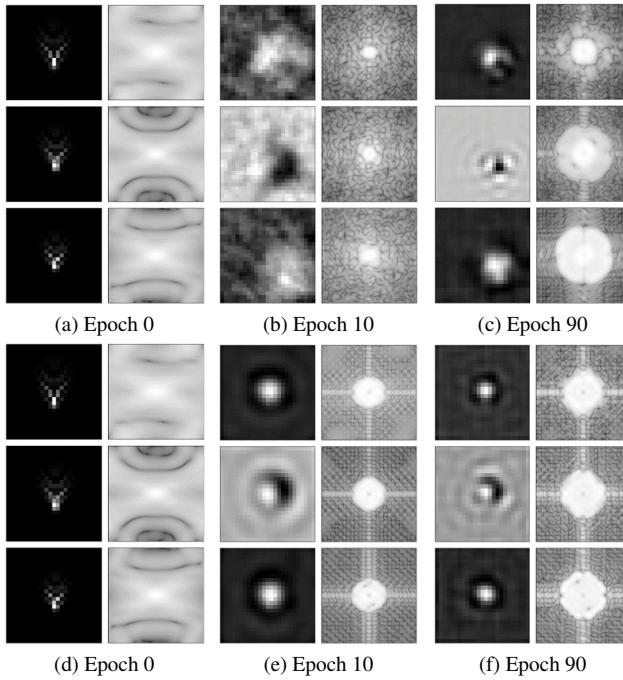
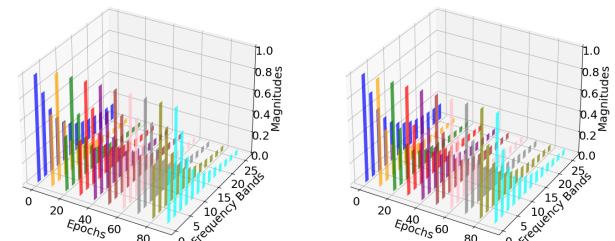


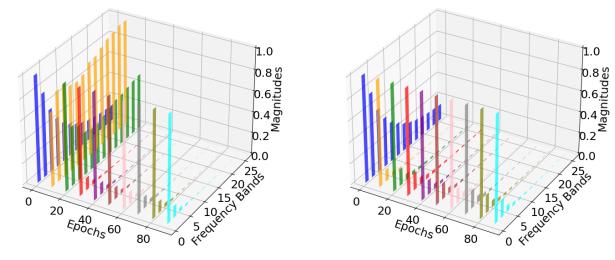
Figure 38. Evolution of kernels in the spatial domain (left) and their frequency magnitudes in \log_{10} scale (right). Kernels in (a-c) and (d-f) were pre-pended to ResNet50 and AlexNet, respectively. Each row corresponds to a color channel.

1235 magnitude of the frequencies in this window, which we refer to as a frequency band. Then, the first bar in each epoch
1236 is the magnitude of the DC component of the spectrum and the last is the mean of the magnitudes of the entire spectrum.
1237 From Fig. 39, we can observe that the spread of the bars depends on the presence of the $L1$ prior. The $L1$ prior
1238 reduces the higher frequency components of the spectrum.
1239 Barring the change in prior, the evolution of the spectrum is
1240 similar.
1241
1242
1243



(a) Alexnet w/o $L1$ Prior

(b) Resnet w/o $L1$ Prior



(c) Alexnet with $L1$ Prior

(d) Resnet with $L1$ Prior

Figure 39. Evolution of spectra of (a-b) Fully Trainable and (c-d) $L1$ Prior kernels. The bar height indicates the average of the absolute value of the Fourier coefficients in different frequency bands (DC component in the front). Each epoch is normalized separately.

10.3. Condition numbers

We use the condition number (CN) as one of the two defining characteristics of our classes of kernels. One significant feature of this description is that the convolution kernel, $\cdot * g$, and the corresponding deconvolution kernel, $\cdot * g^{-1}$, share the same CN. This is particularly important for class I kernels because it suggests that the neural network following $\cdot * g$ can perfectly restore the input without being stymied by noise amplification, which distinguishes class I kernels from the rest. In this section, we will derive the equation for CN provided in Sec. 3 and further inspect the CNs of class II kernels.

Derivation of the Condition Number: For an arbitrary invertible matrix A , CN is defined as:

$$CN(A) = \|A\| \cdot \|A^{-1}\| = \|A^{-1}\| \cdot \|A\| = CN(A^{-1}),$$

where $\|\cdot\|$ is the norm of a matrix. The sensitivity of A to noise or the noise amplification possible by the transformation is equal to that of its inverse. Moreover, $\|A\| = \sigma_{\max}(A)$ and $\|A^{-1}\| = 1/\sigma_{\min}(A)$, where σ_{\max} and σ_{\min} are the maximum and minimum singular values of the matrix, respectively. In our case, the transformation A is a circulant matrix describing the convolution operation. The eigenvectors of A form an orthonormal Fourier basis. Let

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1267 each basis be a column of Q , and let the corresponding
 1268 eigenvalues λ_i , the Fourier coefficients, form a diagonal
 1269 matrix, Λ . A can be decomposed as

1270
$$A = Q\Lambda Q^T,$$

1271 as it is a normal matrix. Then using the relation $\sigma_i =$
 1272 $\sqrt{\lambda_i(A^T A)}$,

$$\begin{aligned}\sigma_{\max}(A) &= \sqrt{\lambda_{\max}(A^T A)} \\ &= \sqrt{\lambda_{\max}((Q\Lambda Q^T)^T Q\Lambda Q^T)} \\ &= \sqrt{\lambda_{\max}(Q\Lambda^2 Q^T)} \\ &= |\lambda_{\max}(A)|.\end{aligned}$$

1274 Similarly, $\sigma_{\min}(A) = |\lambda_{\min}(A)|$. Therefore,

1275
$$CN(A) = \frac{|\lambda_{\max}(A)|}{|\lambda_{\min}(A)|} = \frac{1/|\lambda_{\min}(A)|}{1/|\lambda_{\max}(A)|} = CN(A^{-1}).$$

1276 Note, the definition of CN in the main paper does not include $|.|$. This is equivalent to using the magnitudes of the
 1277 Fourier coefficients instead of the coefficients themselves.
 1278 However, all the numbers presented are correctly computed.
 1279

1280 **Class II Kernels:** The CN of class I kernels is by their
 1281 definition equal to unity and the CNs of class III kernels
 1282 are numerically very high (approaching infinity) due to the
 1283 presence of zeros in the Fourier domain. However, class II
 1284 kernels have high finite CNs, around 10^4 . We can observe
 1285 these in Fig. 40. The range of CNs of the optimized kernels
 1286 is within one order of magnitude despite being pre-pended
 1287 to different neural networks and trained on different data
 1288 sets.

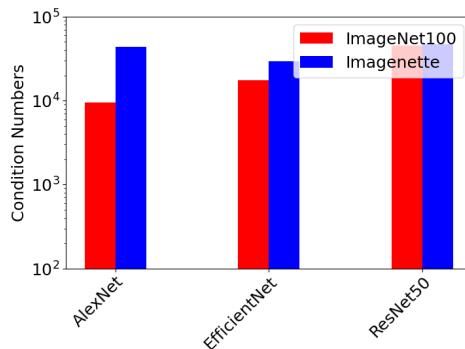


Figure 40. Condition numbers of post-optimization class II kernels pre-pended to different neural networks and trained on ImageNet-100 and ImageNette. Each bar is a mean of up to 5 runs.

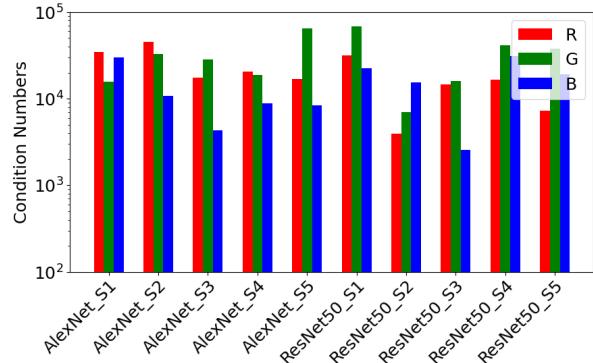


Figure 41. Condition numbers of post-optimization class II kernels pre-pended to networks specified in the xlabel and trained on ImageNette. S1-S5 refers to different seed values, and R, G, and B refer to the color channels.

11. Extended discussion

In the main paper, we only briefly visit two interesting discussion points, which we pick up here to extend the discussion.

11.1. Static Low Pass Filtering

The first point arises when considering Table 4, in particular the results for the non-trainable, static Gaussian kernel. This kernel defines a particular subspace projection, which one would expect to perform well under high-frequency noise. Yet, it also removes information in a fixed and predetermined way which leads to a decay in the clean classification accuracy, which becomes more severe as the hardness of the considered classification task increases. While in ImageNette, the static Gaussian kernel only performs one percent point below a fully learned filter of the same size of 25×25 , the gap is severe when looking at ImageNet-100 in Figure 42. As the classes become more difficult to discriminate, the simple Gaussian blurring yield very unsatisfactory results whereas the proposed fully trainable filter increases model robustness across various corruptions.

11.2. Sparse Coding

The second point is an extended discussion on the relationship of our approach to sparse coding. In the main paper, we started the discussion of why the standard initial layers of the usual networks, do not learn similar patterns as our proposed dimension-preserving large kernel convolution. These layers are known to learn an overcomplete representation of the data in a sparse coding sense, *i.e.* in every channel of every feature map, only very few features are "active". When considering for example a convolution layer with 3×3 kernels, it is also clear that the usual increase in feature map channels from 3 in the input to *e.g.* 64 likely leads to redundant information (a 3×3 filter can

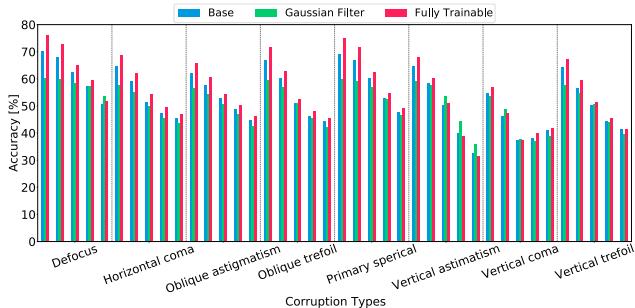


Figure 42. ResNet ImageNet100 OpticsBench. This evaluation demonstrates the benefit of our proposed fully trainable large kernel filter over a fixed Gaussian filter. While the performance gap on ImageNette is smaller, on ImageNet-100, the Gaussian pre-
filtered images can not be reliably classified.

1322 be represented by 9 basis vectors, *i.e.* more than 9 feature
1323 map channels imply redundancy, not only with respect to
1324 the input information but also in the amount of different
1325 (*i.e.* linearly independent) features that can be extracted).
1326 Consequently, such usual early layers are likely to represent
1327 sparse codes, representing every kind of data that is given
1328 in the input. Such behavior is likely to impact a model’s
1329 robustness, as exemplified in the below toy example:

1330 Suppose the input data lives in an approximate sub-
1331 space: the clean data forms (noisy) samples of this subspace
1332 since it is only approximate. If this approximate subspace
1333 is mapped to higher dimensions via a convolutional input
1334 layer with many redundant output channels, more degrees
1335 of freedom are available to fit the data points. The optimizer
1336 may therefore choose to use additional redundant dimen-
1337 sions to marginally increase the fit of the noisy data points.
1338 While this improves performance on the clean data, the ex-
1339 trapulation capabilities to unseen samples may decrease,
1340 giving rise to overfitting and by implication worse perfor-
1341 mance on corrupted data. Forcing the optimizer to stay
1342 in the original subspace and choosing a more constrained
1343 model therefore may aid extrapolation, *i.e.* generalization
1344 capabilities of the resulting model.

1345 In our layer, such behavior is avoided, since we are map-
1346 ping every channel of the input to exactly one channel in
1347 the output of our layer. If a particular piece of informa-
1348 tion from the input is to be preserved in the output of our
1349 layer, the large kernel with its $3 \times 25 \times 25$ weights has to
1350 be learned appropriately. As discussed in section 3 **Con-**
1351 **tent Preserving Filters**, only a limited set of special filters
1352 can fully preserve the input content. Yet, for the kernel to
1353 be learned, it uses gradient signal that is provided from the
1354 classification task at hand. In particular, if an input fea-
1355 ture is not contributing to the discrimination between given
1356 classes, there will be no gradient signal provided to the
1357 input that encourages to learn this particular feature, and since

the representation learned is not over-complete, it will also
1358 not be learned by chance. A large kernel makes it partic-
1359 ularly hard to learn local patterns by chance. It learns to
1360 predominantly represent signal that is *needed* for the task at
1361 hand, *i.e.* the *Essential*.

One potential interpretation is to understand the trained
layer as an information bottleneck, through which only such
information is passed, that is explicitly needed for the clas-
sifier to perform well on the training set. However, remov-
ing the additional parts of the signal that are not helpful for
classification makes the model more robust under input cor-
ruptions.

1362
1363
1364
1365
1366
1367
1368
1369