# Improving Robustness Without Sacrificing Accuracy with Patch Gaussian Augmentation

**Raphael Gontijo Lopes**[1][*], **Dong Yin**[2][†], **Ben Poole**[1], **Justin Gilmer**[1], **Ekin D. Cubuk**[1]
[1]Google Brain, [2]UC Berkeley
{iraphael,pooleb,gilmer,cubuk}@google.com, dongyin@berkeley.edu

## Abstract

Deploying machine learning systems in the real world requires both high accuracy on clean data and robustness to naturally occurring corruptions. While architectural advances have led to improved accuracy, building robust models remains challenging. Prior work has argued that there is an inherent trade-off between robustness and accuracy, which is exemplified by standard data augment techniques such as `Cutout`, which improves clean accuracy but not robustness, and additive `Gaussian` noise, which improves robustness but hurts accuracy. To overcome this trade-off, we introduce `Patch Gaussian`, a simple augmentation scheme that adds noise to randomly selected patches in an input image. Models trained with `Patch Gaussian` achieve state of the art on the CIFAR-10 and ImageNet Common Corruptions benchmarks while also improving accuracy on clean data. We find that this augmentation leads to reduced sensitivity to high frequency noise (similar to `Gaussian`) while retaining the ability to take advantage of relevant high frequency information in the image (similar to `Cutout`). Finally, we show that `Patch Gaussian` can be used in conjunction with other regularization methods and data augmentation policies such as AutoAugment, and improves performance on the COCO object detection benchmark.

## 1 Introduction

Modern deep neural networks can achieve impressive performance at classifying images in curated datasets [31, 33, 27]. Yet their performance is not robust to corruptions that typically occur in real-world settings. For example, neural networks are sensitive to small translations and changes in scale [2], blurring and additive noise [10], small objects placed in images [41], and even different images from a similar distribution of the training set [40, 39]. For models to be useful in the real world, they need to be both accurate on a high-quality held-out set of images , which we refer to as "clean accuracy," and robust on corrupted images, which we refer to as "robustness." Most of the literature in machine learning has focused on architectural changes [43, 44, 23, 54, 45, 22, 55, 26, 36] to improve clean accuracy but has recently become interested in robustness as well.
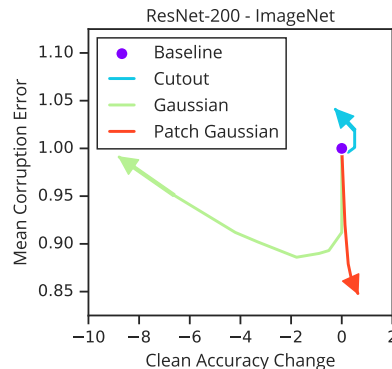


Figure 1: `Patch Gaussian` augmentation overcomes the accuracy/robustness trade-off observed in other augmentation strategies. Larger $\sigma$ of `Patch Gaussian` improves both mean corruption error (mCE) and clean accuracy, whereas larger $\sigma$ of `Gaussian` and patch size of `Cutout` hurt accuracy or robustness. More robust and accurate models are down and to the right.

---

[*]Work done as a member of the Google AI Residency program g.co/airesidency.
[†]Work completed during internship at Google Brain.

Research in neural network robustness has tried to quantify the problem by establishing benchmarks that directly measure it [24, 21] and comparing the performance of humans and neural networks [16]. Others have tried to understand robustness by highlighting systemic failure modes of current learning methods. For instance, networks exhibit excessive invariance to visual features [29], texture bias [15], sensitivity to worst-case (adversarial) perturbations [19], and a propensity to rely solely on non-robust, but highly predictive features for classification [11, 28]. Of particular relevance to our work, Ford et al [13] show that in order to get adversarial robustness, one needs robustness to noise-corrupted data.

Another line of work has attempted to increase model robustness performance, either by directly projecting out superficial statistics [48], via architectural improvements [5], pre-training schemes [25], or through the use of data augmentations. Data augmentation increases the size and diversity of the training set, and provides a simple method for learning invariances that are challenging to encode architecturally [6]. Recent work in this area includes learning better transformations [9, 52, 53], inferring combinations of transformations [6], unsupervised methods [49], theory of data augmentation [7], and applications for one-shot learning [1].

Despite these advances, individual data augmentation methods that improve robustness do so at the expense of reduced clean accuracy. Some have even claimed that there exists a fundamental trade-off between the two [47]. Because of this, many recent works focus on improving either one or the other [37, 15]. In this work we propose a data augmentation that overcomes this trade-off, achieving both improved robustness and clean accuracy. Our contributions are as follows:

- We characterize a trade-off between robustness and accuracy among two standard data augmentations: `Cutout` and `Gaussian` (Section 2.1).
- We devise a simple data augmentation method (which we term `Patch Gaussian`) that allows us to interpolate between the two augmentations above. (Section 3.1)
- We find that `Patch Gaussian` allows us to overcome the observed trade-off (Figure 1, Section 4.1), and achieves a new state of the art in the Common Corruptions benchmark [24] on CIFAR-C and ImageNet-C. (Section 4.2)
- We demonstrate that `Patch Gaussian` can be combined with other regularization strategies (Section 4.3) and data augmentation policies [6] (Section 4.4), and can improve COCO object detection performance as well (Section 4.5).
- We perform a frequency-based analysis [50] of models trained with `Patch Gaussian` and find that they can better leverage high-frequency information in lower layers, while not being too sensitive to them at later ones (Section 5.1).

## 2 Preliminaries

We start by considering two data augmentations: `Cutout` [9] and `Gaussian` [20]. The former sets a random patch of the input images to a constant (the mean pixel in the dataset) and is successful at improving clean accuracy. The latter works by adding independent Gaussian noise to each pixel of the input image, which can increase robustness to Gaussian noise directly.

To apply `Gaussian`, we uniformly sample a standard deviation $\sigma$ from 0 up to some maximum value $\sigma_{max}$, and add i.i.d. noise sampled from $\mathcal{N}(0, \sigma^2)$ to each pixel. To apply `Cutout`, we use a fixed patch size $W$, and randomly set a square region with size $W \times W$ to the constant mean of each RGB channel in the dataset. As in [9], the patch location is randomly sampled and can lie outside of the $32 \times 32$ CIFAR-10 (or $224 \times 224$ ImageNet) image but its center is constrained to lie within it. Patch sizes and $\sigma_{max}$ are selected based on the method described in Section 3.2.

### 2.1 Cutout and Gaussian exhibit a trade-off between accuracy and robustness

We compare the effectiveness of `Gaussian` and `Cutout` data augmentation for accuracy and robustness by measuring the performance of models trained with each on clean data, as well as data corrupted by various standard deviations of Gaussian noise. Figure 2 highlights an apparent trade-off in using these methods. In accordance to previous work [9], `Cutout` improves accuracy on clean test data. Despite this, we find it does not lead to increased robustness. Conversely, training with higher $\sigma$ of `Gaussian` can lead to increased robustness to Gaussian noise, but it also leads to decreased accuracy on clean data. Therefore, any robustness gains are offset by poor overall performance.
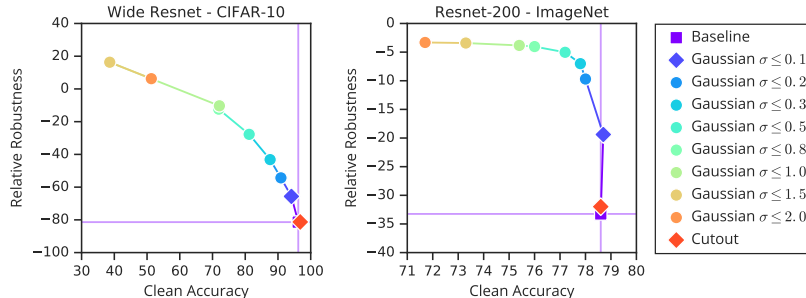
Figure 2: The robustness-accuracy trade-off between `Cutout` and `Gaussian` augmentations. Each dot represents a model trained with different augmentations and hyper-parameters. The y-axis is the change in accuracy when tested on data corrupted with Gaussian noise at various $\sigma$ (average corrupted accuracy minus clean accuracy). The diamond indicates augmentation hyper-parameters selected by the method in Section 3.2.

At first glance, these results seem to reinforce the findings of previous work [47], indicating that robustness comes at the cost of generalization. In the following sections, we will explore whether there exists augmentation strategies that do not exhibit this limitation.

## 3 Method

Each of the two methods seen so far achieves one half of our stated goal: either improving robustness or improving clean test accuracy, but never both. To overcome the limitations of existing data augmentation techniques, we introduce `Patch Gaussian`, a new technique that combines the noise robustness of `Gaussian` with the improved clean accuracy of `Cutout`.

### 3.1 Patch Gaussian

`Patch Gaussian` works by adding a $W$x$W$ patch of Gaussian noise to the image (Figure 3)[3]. As with `Cutout`, the center of the patch is sampled to be within the image. By varying the size of this patch and the maximum standard deviation of noise sampled $\sigma_{max}$, we can interpolate between `Gaussian` (which applies additive Gaussian noise to the whole image) and an approximation of Cutout (which removes all information inside the patch). See Figure 9 for more examples.



Figure 3: `Patch Gaussian` is the addition of Gaussian noise to pixels in a square patch. It allows us to interpolate between `Gaussian` and `Cutout`, approaching `Gaussian` with increasing patch size and `Cutout` with increasing $\sigma$.

All image transformations, including `Patch Gaussian` are performed on images with unnormalized pixel values in $[0, 1]$ range. For all images, standard random flipping and cropping is applied immediately *after* any augmentations mentioned on CIFAR-10 (*before*, on Imagenet). After noise-based augmentations, images are clipped to the $[0, 1]$ range.

### 3.2 Hyper-parameter selection

Our goal is to learn models that achieve both good clean accuracy and improved robustness to corruptions. When selecting hyper-parameters we need to decide how to weight these two metrics. Here, we focused on identifying the models that were most robust while still achieving a minimum accuracy (Z) on the clean test data. Values of Z vary per dataset and model, and can be found in the Appendix (Table 6). If no model has clean accuracy $\geq$Z, we report the model with highest clean accuracy, unless otherwise specified. We find that patch sizes around 25 on CIFAR ($\leq$250 on ImageNet, i.e.: uniformly sampled with maximum value 250) with $\sigma \leq 1.0$ generally perform the best. A complete list of selected hyper-parameters for all augmentations can be found in Table 6.

---

[3]A TensorFlow implementation of `Patch Gaussian` can be found in Appendix (Figure 10).

Here, robustness is defined as average accuracy of the model, when tested on data corrupted by various $\sigma$ (0.1, 0.2, 0.3, 0.5, 0.8, 1.0) of Gaussian noise, relative to the clean accuracy. This metric is correlated with mCE [13], so it ensures model rosbustness is generally useful beyond Gaussian corruptions. By picking models based on their Gaussian noise robustness, we ensure that our selection process does not overfit to the Common Corruptions benchmark [24].

$$\text{Relative Gaussian Robustness} = \mathbb{E}_{\sigma}(\text{Accuracy on Data Corrupted by } \sigma) - \text{Clean Accuracy}$$

### 3.3 Models, Datasets, & Implementation Details

We run our experiments on CIFAR-10 [32] and ImageNet [8] datasets. On CIFAR-10, we use the Wide-ResNet-28-10 model [51], as well as the Shake-shake-112 model [14], trained for 200 epochs and 600 epochs respectively. The Wide-ResNet model uses a initial learning rate of 0.1 with a cosine decay schedule. Weight decay is set to be 5e-4 and batch size is 128. We train all models, including the baseline, with standard data augmentation of horizontal flips and pad-and-crop. Our code uses the same hyper parameters as [6] [4].

On ImageNet, we use the ResNet-50 and Resnet-200 models [23], trained for 90 epochs. We use a weight decay rate of 1e-4, global batch size of 512 and learning rate of 0.2. The learning rate is decayed by 10 at epochs 30, 60, and 80. We use standard data augmentation of horizontal flips and crops. All CIFAR-10 and ImageNet experiments use the listed hyper-parameters above, unless specified otherwise. Our code uses the same hyper parameters as open-sourced implementations[5].

## 4 Results

We show that models trained with `Patch Gaussian` can overcome the trade-off observed in Fig. 2 and learn models that are robust while maintaining their generalization accuracy (Section 4.1). In doing so, we establish a new state of the art in CIFAR-C and ImageNet-C Common Corruptions benchmark [24] (Section 4.2). We then show that `Patch Gaussian` can be used in complement to other common regularization strategies (Section 4.3), data augmentation policies [6] (Section 4.4), and that it can also improve training of object detection models (Section 4.5)

### 4.1 Patch Gaussian overcomes this trade-off and improves both accuracy and robustness

We train models on various hyper-parameters of `Patch Gaussian` and find that the model selected by the method in Section 3.2 leads to improved robustness to Gaussian noise, like `Gaussian`, while also improving clean accuracy, much like `Cutout`. In Figure 4, we visualize these results in an ablation study, either varying patch sizes and fixing $\sigma$ to the selected value (1.0) or varying noise level $\sigma$ and fixing the patch size to the selected value (350).
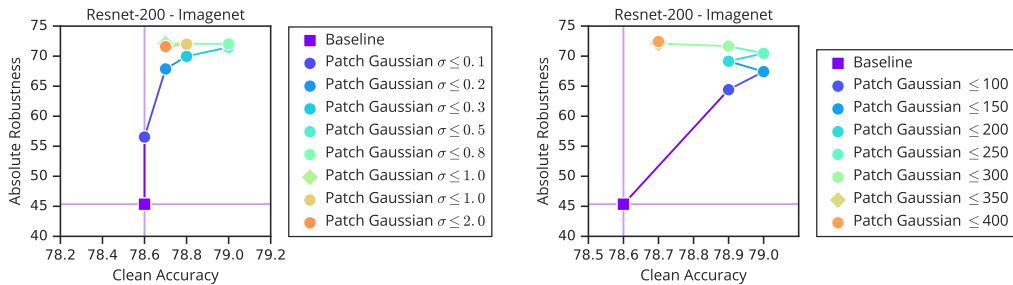


Figure 4: Training with `Patch Gaussian` improves clean data accuracy and robustness simultaneously. Each dot represents a model trained with various $\sigma$ (left) or patch sizes (right), while keeping the other fixed at the value indicated by the diamond. The y-axis is the average absolute accuracy when tested on data corrupted by Gaussian noise at various $\sigma$. The diamond indicates the augmentation hyper-parameters selected by the method in Section 3.2.

---

[4] Available at `https://github.com/tensorflow/models/tree/master/research/autoaugment`

[5] Available at `https://github.com/tensorflow/tpu/tree/master/models/official/resnet`

## 4.2 Training with Patch Gaussian leads to improved Common Corruption robustness

In this section, we look at how our augmentations impact robustness in a more realistic setting beyond robustness to additive Gaussian noise. Rather than focusing on adversarial examples that are worst-case bounded perturbations, we focus on a more general set of corruptions [18] that models are likely to encounter in real-world settings: the Common Corruptions benchmark [24]. This benchmark, also referred to as CIFAR-C and ImageNet-C, is composed of images transformed with 15 corruptions, at 5 severities each. The corruptions are designed to model those commonly found in real-world settings, such as brightness, different weather conditions, and different kinds of noise.

Tables 1 and 2 show that `Patch Gaussian` achieves state of the art on both of these benchmarks in terms of mean Corruption Error (mCE). However, ImageNet-C was released in compressed JPEG format [12], which alters the corruptions applied to the raw pixels. Therefore, we report results on the benchmark as-released ("Original mCE") as well as a version of 12 corruptions without the extra compression ("mCE") [6].

Table 1: `Patch Gaussian` achieves state of the art in the CIFAR-C benchmark [24] while improving clean accuracy. Augmentation hyper-parameters were selected based on the method in Section 3.2 and can be found in Appendix. *`Cutout 16` is presented for direct comparison with [9, 14].

|  | Augmentation | Test Accuracy | mCE | mCE (-noise) |
|---|---|---|---|---|
| Wide Resnet-28-10 | Adversarial | 87.3% | 1.049 | 1.157 |
|  | Baseline | 96.2% | 1.000 | 1.000 |
|  | Cutout | 96.8% | 1.265 | 1.185 |
|  | Cutout 16* | **97.0%** | 1.002 | 0.953 |
|  | Gaussian | 94.1% | 0.887 | 0.995 |
|  | Patch Gaussian | 96.6% | **0.797** | **0.858** |
| Shake 112 | Baseline | 96.8% | 1.000 | 1.000 |
|  | Cutout | 97.1% | 0.946 | 0.930 |
|  | Cutout 16* | **97.5%** | 0.912 | 0.872 |
|  | Gaussian | 94.6% | 0.977 | 1.111 |
|  | Patch Gaussian | 97.2% | **0.713** | **0.776** |

Table 2: `Patch Gaussian` achieves state of the art in the ImageNet-C benchmark [24] while improving uncorrupted test accuracy. "SIN+IN fIN" is the shape-biased model from [15]. "Original mCE" refers to the jpeg-compressed benchmark, as used in [15, 24]. "mCE" is a version of it without the extra jpeg compression. Note that `Patch Gaussian` improves robustness even in corruptions that aren't noise-based. Augmentation hyper-parameters were selected based on the method in Section 3.2 and can be found in Appendix. For Resnet-200, we also present `Gaussian` at a higher $\sigma$ to highlight the accuracy-robustness trade-off.

|  | Augmentation | Test Accuracy | Original mCE | Original mCE (-noise) | mCE | mCE (-noise) |
|---|---|---|---|---|---|---|
| Resnet-50 | SIN+IN ftIN | 76.7% | 0.738 | 0.731 | - | - |
|  | Baseline | **76.4%** | 0.753 | 0.763 | 1.00 | 1.00 |
|  | Cutout | 76.2% | 0.758 | 0.766 | 1.007 | 1.005 |
|  | Gaussian | 75.6% | 0.739 | 0.754 | 0.898 | 0.991 |
|  | Patch Gaussian | 76.0% | **0.714** | **0.736** | **0.872** | **0.969** |
| Resnet-200 | Baseline | 78.6% | 0.675 | 0.686 | 0.881 | 0.883 |
|  | Cutout | 78.6% | 0.671 | 0.687 | 0.874 | 0.884 |
|  | Gaussian | **78.7%** | 0.658 | 0.678 | 0.795 | 0.881 |
|  | Gaussian ($\sigma \leq 0.2$) | 78.1% | 0.644 | 0.665 | 0.784 | 0.874 |
|  | Patch Gaussian | **78.7%** | **0.604** | **0.634** | **0.736** | **0.818** |

---

[6]Available at `https://github.com/tensorflow/datasets` under `imagenet2012_corrupted`

To compute mCE, we normalize the corruption error for each model and dataset to the baseline with only flip and crop data augmentation. The one exception is Original mCE ImageNet, where we use the AlexNet baseline to be directly comparable with previous work [24, 15].

Because `Patch Gaussian` is a noise-based augmentation, we wanted to verify whether its gains on this benchmark were solely due to improved performance on noise-based corruptions (Gaussian Noise, Shot Noise, and Impulse Noise). To do this, we also measure the models' average performance on all *other* corruptions, reported as "Original mCE (-noise)", and "mCE (-noise)". We observe that `Patch Gaussian` outperforms all other models, even on corruptions like fog where `Gaussian` hurts performance [13]. Scores for each corruption can be found in the Appendix (Tables 7 and 8).

Comparing the lower capacity ResNet-50 and Wide ResNet models to higher-capacity ResNet-200 and Shake 112 models, we find diminished gains in clean accuracy and robustness. Still, `Patch Gaussian` achieves a substantial increase in mCE relative to other augmentation strategies.

### 4.3 Patch Gaussian can be combined with other regularization strategies

Since `Patch Gaussian` has a regularization effect on the models trained above, we compare it with other regularization methods: larger weight decay, label smoothing, and dropblock (Table 3). We find that while label smoothing improves clean accuracy, it weakens the robustness in all corruption metrics we have considered. This agrees with the theoretical prediction from [5], which argued that increasing the confidence of models would improve robustness, whereas label smoothing reduces the confidence of predictions. We find that increasing the weight decay from the default value used in all models does not improve clean accuracy or robustness.

Here, we focus on analyzing the interaction of different regularization methods with `Patch Gaussian`. Previous work indicates that improvements on the clean accuracy appear after training with Dropblock for 270 epochs [17], but we did not find that training for 270 epochs changed our analysis. Thus, we present models trained at 90 epochs for direct comparison with other results. Due to the shorter training time, Dropblock does not improve clean accuracy, yet it does make the model more robust (relative to baseline) according to all corruption metrics we consider.

We find that using label smoothing in addition to `Patch Gaussian` has a mixed effect, it improves clean accuracy while slightly improving robustness metrics except for the Original mCE. Combining Dropblock with `Patch Gaussian` reduces the clean accuracy relative to the `Patch Gaussian`-only model, as Dropblock seems to be a strong regularizer when used for 90 epochs. However, using Dropblock and `Patch Gaussian` together leads to the best robustness performance. These results indicate that `Patch Gaussian` can be used in conjunction with existing regularization strategies.

Table 3: `Patch Gaussian` can be used with other regularization methods for improved robustness. "Original mCE" refers to the jpeg-compressed benchmark, as used in [15, 24]. "mCE" is a version of it without the extra jpeg compression. All of the models are ResNet-50 trained on ImageNet with same hyperparameters for 90 epochs.

| | Regularization | Test Accuracy | Original mCE | Original mCE (-noise) | mCE | mCE (-noise) |
|---|---|---|---|---|---|---|
| Resnet-50 | Label Smoothing | **76.7%** | 0.760 | 0.765 | 1.01 | 1.01 |
| | Larger Weight Decay (0.001) | 74.9% | 0.766 | 0.777 | 1.02 | 1.03 |
| | Dropblock | 76.3% | 0.734 | 0.743 | 0.971 | 0.974 |
| | Patch Gaussian + Label Smoothing | 76.5% | 0.720 | 0.734 | **0.868** | 0.966 |
| | Patch Gaussian + Dropblock | 75.7% | **0.708** | **0.726** | 0.870 | **0.961** |

## 4.4 Patch Gaussian can be combined with AutoAugment policies for improved results

Knowing that `Patch Gaussian` can be combined with other regularizers, it's natural to ask whether it can also be combined with other data augmentation policies. Table 4 highlights models trained with AutoAugment [6]. For fair comparison of mCE scores, we train all models with the best AutoAugment policy, but without contrast and Inception color pre-processing, as those are present in the Common Corruptions benchmark. This process is imperfect since AutoAugment has *many* operations, some of which could still be correlated with corruptions. Regardless, we find that `Patch Gaussian` improves accuracy and robustness over simply using AutoAugment.

Because AutoAugment leads to state of the art accuracies, we are interested in seeing how far it can be combined with `Patch Gaussian` to improve results. Therefore, and unlike previous experiments, models are trained for 180 epochs to yield best results possible.

Table 4: `Patch Gaussian` can be combined with AutoAugment [6] data augmentation policy for improved results. "Original mCE" refers to the jpeg-compressed benchmark, as used in [15, 24]. "mCE" is a version of it without the extra jpeg compression. All of the models are ResNet-50 trained on ImageNet with best AutoAugment policy for 180 epochs, to highlight improvements.

| | Model (trained with AutoAugment) | Test Accuracy | Original mCE | Original mCE (-noise) | mCE | mCE (-noise) |
|---|---|---|---|---|---|---|
| ResNet 50 | Baseline | 77.0% | 0.674 | 0.697 | 0.855 | 0.882 |
| | Patch Gaussian ($W=150, \sigma \leq 0.5$) | **77.3%** | **0.656** | **0.682** | **0.779** | **0.863** |

## 4.5 Patch Gaussian improves performance in object detection

Since `Patch Gaussian` can be combined with both regularization strategies as well as data augmentation policies, we want to see if it is generally useful beyond classification tasks. We train a RetinaNet detector [35] with ResNet-50 backbone [23] on the COCO dataset [34]. Images for both baseline and `Patch Gaussian` models are horizontally flipped half of the time, after being resized to $640 \times 640$. We train both models for 150 epochs using a learning rate of 0.08 and a weight decay of $1e-4$. The focal loss parameters are set to be $\alpha = 0.25$ and $\gamma = 1.5$.

Despite being designed for classification, `Patch Gaussian` improves detection performance according to all metrics when tested on the clean COCO validation set (Table 5). On the primary COCO metric mean average precision (mAP), the model trained with `Patch Gaussian` achieves a 1% higher accuracy over the baseline, whereas the model trained with `Gaussian` suffers a 2.9% loss.

Next, we evaluate these models on the validation set corrupted by i.i.d. Gaussian noise, with $\sigma = 0.25$. We find that model trained with `Gaussian` and `Patch Gaussian` achieve the highest mAP of 26.1% on the corrupted data, whereas the baseline achieves 11.6%. It is interesting to note that `Patch Gaussian` model achieves a better result on the harder metrics of small object detection and stricter intersection over union (IOU) thresholds, whereas the `Gaussian` model achieves a better result on the easier tasks of large object detection and less strict IOU threshold metric.

Table 5: Mean average precision (mAP) on COCO with baseline augmentation of horizontal flips and `Patch Gaussian`. $mAP_S$, $mAP_M$, and $mAP_L$ refer to mAP for small, medium, and large objects, respectively. $mAP_{50}$ and $mAP_{75}$ refer to mAP at intersection over union values of 50 and 75, respectively. mAP in the final column is the averaged mAP over IoU=0.5:0.05:0.95.

| Tested on | | $mAP_S$ | $mAP_M$ | $mAP_L$ | $mAP_{50}$ | $mAP_{75}$ | mAP |
|---|---|---|---|---|---|---|---|
| Clean Data | Baseline | 15.6 | 36.9 | 48.3 | 50.8 | 35.6 | 33.2 |
| | Gaussian ($\sigma \leq 1.0$) | 13.1 | 32.6 | 44.0 | 45.7 | 31.2 | 29.3 |
| | Patch Gaussian ($W=200, \sigma \leq 1.0$) | **16.1** | **37.9** | **50.3** | **51.9** | **36.5** | **34.2** |
| Gaussian Noise (=0.25) | Baseline | 4.5 | 12.7 | 17.6 | 19.3 | 11.7 | 11.6 |
| | Gaussian ($\sigma \leq 1.0$) | 9.9 | 28.1 | **41.0** | **41.7** | 26.8 | **26.1** |
| | Patch Gaussian ($W=200, \sigma \leq 1.0$) | **10.1** | **28.2** | 40.4 | 41.3 | **27.2** | **26.1** |

Overall, as was observed for the classification tasks, training object detection models with `Patch Gaussian` leads to significantly more robust models without sacrificing clean accuracy.

# 5   Discussion

In an attempt to understand `Patch Gaussian`'s performance, we perform a frequency-based analysis of models trained with various augmentations using the method introduced in [50].

First, we perturb each image in the dataset with noise sampled at each orientation and frequency in Fourier space. Then, we measure changes in the network activations and test error when evaluated with these Fourier-noise-corrupted images: we measure the change in $\ell_2$ norm of the tensor directly after the first convolution, as well as the absolute test error. This procedure yields a heatmap, which indicates model sensitivity to different frequency and orientation perturbations in the Fourier domain. Each image in Fig 5 shows first layer (or test error) sensitivity as a function of frequency and orientation of the sampled noise, with the middle of the image containing the lowest frequencies, and the edges of the image containing highest frequencies.

For CIFAR-10 models, we present this analysis for the entire Fourier domain, with noise sampled with norm 4. For ImageNet, we focus our analysis on lower frequencies that are more visually salient add noise with norm 15.7.

Note that for `Cutout` and `Gaussian`, we chose larger patch sizes and $\sigma$s than those selected with the method in Section 3.2 in order to highlight the effect of these augmentations on sensitivity. Heatmaps of other models can be found in the Appendix (Figure 12).

## 5.1   Frequency-based analysis of models trained with Patch Gaussian

We confirm findings by [50] that `Gaussian` encourages the model to learn a low-pass filter of the inputs. Models trained with this augmentation, then, have low test error sensitivity at high frequencies, which could help robustness. However, valuable high-frequency information [4] is being thrown out at low layers, which could explain the lower test accuracy.

We further find that `Cutout` encourages the use of high-frequency information, which could help explain its improved generalization performance. Yet, it does not encourage lower test error sensitivity, which explains why it doesn't improve robustness either.

`Patch Gaussian`, on the other hand, seems to allow high-frequency information through at lower layers, but still encourages relatively lower test error sensitivity at high frequencies. Indeed, when we measure accuracy on images filtered with a high-pass filter, we see that `Patch Gaussian` models can maintain accuracy in a similar way to the baseline and to `Cutout`, where `Gaussian` fails to. See Figure 5 for full results.

Understanding the impact of data distributions and noise on representations has been well-studied in neuroscience [3, 42, 30]. The data augmentations that we propose here alter the distribution of inputs that the network sees, and thus are expected to alter the kinds of representations that are learned. Prior work on efficient coding [30] and autoencoders [38] has shown how filter properties change with noise in the unsupervised setting, resulting in lower-frequency filters with `Gaussian`, as we observe in Fig. 5. Consistent with prior work on natural image statistics [46], we find that networks are least sensitive to low frequency noise where the spectral density is largest. Performance drops at higher frequencies when the amount of noise we add grows relative to the typical spectral density observed at these frequencies. In future work, we hope to better understand the relationship between naturally occuring properties of images and sensitivity, and investigate whether training with more naturalistic noise can yield similar gains in corruption robustness.
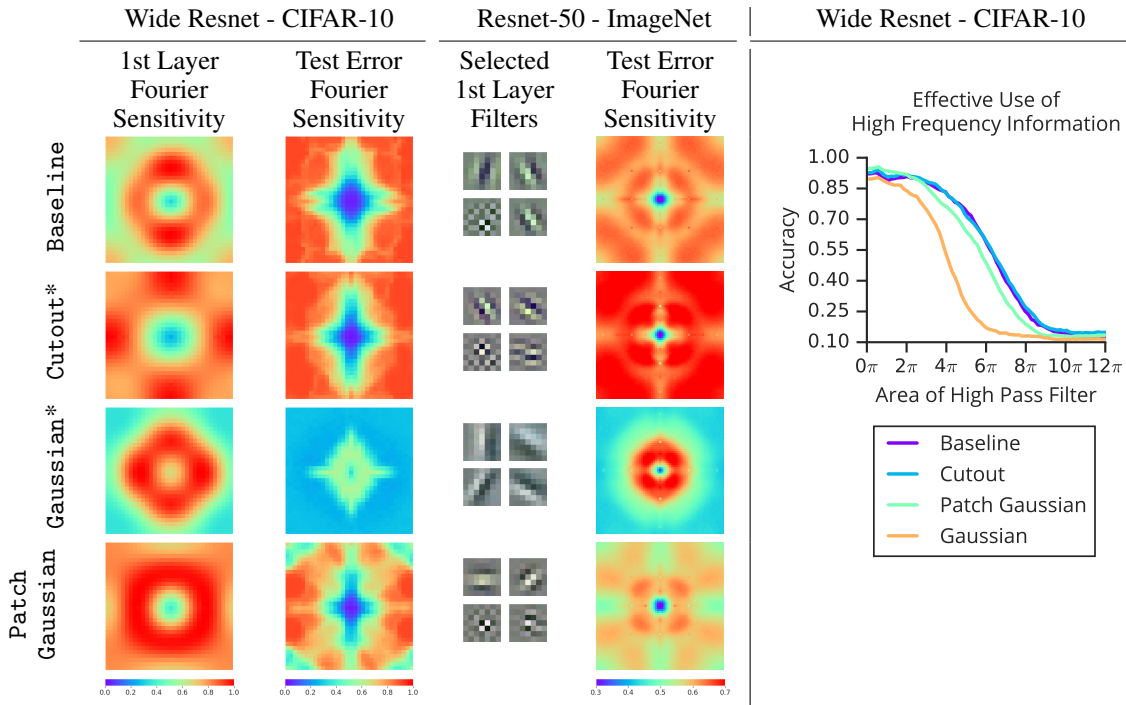
Figure 5: (left) Fourier analysis of various models, using method from [50]. Heatmaps depict model sensitivity to various sinusoidal gratings. `Cutout` encourages the use of high frequencies in earlier layers, but its test error remains too sensitive to them. `Gaussian` learns low-pass filtering of features, which increases robustness at later layers, but makes lower layers too invariant to high-frequency information (thus hurting accuracy). `Patch Gaussian` allows high frequencies to be used in lower layers, and its test error remains relatively robust to them. This can also be seen by the presence of high-frequency kernels in the first layer filters of the models (or lack thereof, in the case of `Gaussian`). (right) Indeed, `Patch Gaussian` models match the performance of `Cutout` and `Baseline` when presented with only the high frequency information of images, whereas `Gaussian` fails to effectively utilize this information (see Appendix Fig. 13 for experiment details). This pattern of reduced sensitivity of predictions to high frequencies in the input occurs across all augmentation magnitudes, but here we use larger patch sizes and $\sigma$ of noise to highlight the differences in models indicated by *. See text for details.

# 6 Conclusion

In this work, we introduced a single data augmentation operation, `Patch Gaussian`, which improves robustness to common corruptions without incurring a drop in clean accuracy. For models that are large relative to the dataset size (like ResNet-200 on ImageNet and all models on CIFAR-10), `Patch Gaussian` improves clean accuracy and robustness concurrently. We showed that `Patch Gaussian` achieves this by interpolating between two standard data augmentation operations `Cutout` and `Gaussian`. We also demonstrate that `Patch Gaussian` can be used in conjunction with other regularization and data augmentation strategies, and can also improve the performance of object detection models, indicating it is generally useful. Finally, we analyzed the sensitivity to noise in different frequencies of models trained with `Cutout` and `Gaussian`, and showed that `Patch Gaussian` combines their strengths without inheriting their weaknesses.

# References

[1] Asano, Y. M., Rupprecht, C., and Vedaldi, A. Surprising effectiveness of few-image unsupervised feature learning, 2019.

[2] Azulay, A. and Weiss, Y. Why do deep convolutional networks generalize so poorly to small image transformations? *arXiv preprint arXiv:1805.12177*, 2018.

[3] Barlow, H. B. et al. Possible principles underlying the transformation of sensory messages. *Sensory communication*, 1:217–234, 1961.

[4] Brendel, W. and Bethge, M. Approximating cnns with bag-of-local-features models works surprisingly well on imagenet. *arXiv preprint arXiv:1904.00760*, 2019.

[5] Cubuk, E. D., Zoph, B., Schoenholz, S. S., and Le, Q. V. Intriguing properties of adversarial examples. *arXiv preprint arXiv:1711.02846*, 2017.

[6] Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., and Le, Q. V. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018.

[7] Dao, T., Gu, A., Ratner, A. J., Smith, V., De Sa, C., and Ré, C. A kernel theory of modern data augmentation. *arXiv preprint arXiv:1803.06084*, 2018.

[8] Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.

[9] DeVries, T. and Taylor, G. W. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.

[10] Dodge, S. and Karam, L. A study and comparison of human and deep learning recognition performance under visual distortions. In *2017 26th international conference on computer communication and networks (ICCCN)*, pp. 1–7. IEEE, 2017.

[11] Doersch, C., Gupta, A., and Efros, A. A. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1422–1430, 2015.

[12] ECMA International. *JPEG Interchange Format (JFIF)*. 2009.

[13] Ford, N., Gilmer, J., Carlini, N., and Cubuk, D. Adversarial examples are a natural consequence of test error in noise. *arXiv preprint arXiv:1901.10513*, 2019.

[14] Gastaldi, X. Shake-shake regularization. *arXiv preprint arXiv:1705.07485*, 2017.

[15] Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F. A., and Brendel, W. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018.

[16] Geirhos, R., Temme, C. R., Rauber, J., Schütt, H. H., Bethge, M., and Wichmann, F. A. Generalisation in humans and deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 7538–7550, 2018.

[17] Ghiasi, G., Lin, T.-Y., and Le, Q. V. Dropblock: A regularization method for convolutional networks. In *Advances in Neural Information Processing Systems*, pp. 10727–10737, 2018.

[18] Gilmer, J., Adams, R. P., Goodfellow, I., Andersen, D., and Dahl, G. E. Motivating the rules of the game for adversarial example research. *arXiv preprint arXiv:1807.06732*, 2018.

[19] Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.

[20] Grandvalet, Y. and Canu, S. Noise injection for inputs relevance determination. *Advances in intelligent systems*, 41:378, 1997.

[21] Gu, K., Yang, B., Ngiam, J., Le, Q., and Shlens, J. Using videos to evaluate image model robustness. *arXiv preprint arXiv:1904.10076*, 2019.

[22] Han, D., Kim, J., and Kim, J. Deep pyramidal residual networks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6307–6315. IEEE, 2017.

[23] He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.

[24] Hendrycks, D. and Dietterich, T. G. Benchmarking neural network robustness to common corruptions and surface variations. *arXiv preprint arXiv:1807.01697*, 2018.

[25] Hendrycks, D., Lee, K., and Mazeika, M. Using pre-training can improve model robustness and uncertainty. *arXiv preprint arXiv:1901.09960*, 2019.

[26] Hu, J., Shen, L., and Sun, G. Squeeze-and-excitation networks. *arXiv preprint arXiv:1709.01507*, 2017.

[27] Huang, Y., Cheng, Y., Chen, D., Lee, H., Ngiam, J., Le, Q. V., and Chen, Z. Gpipe: Efficient training of giant neural networks using pipeline parallelism. *arXiv preprint arXiv:1811.06965*, 2018.

[28] Ilyas, A., Santurkar, S., Tsipras, D., Engstrom, L., Tran, B., and Madry, A. Adversarial examples are not bugs, they are features. *arXiv preprint arXiv:1905.02175*, 2019.

[29] Jacobsen, J.-H., Behrmann, J., Zemel, R., and Bethge, M. Excessive invariance causes adversarial vulnerability. *arXiv preprint arXiv:1811.00401*, 2018.

[30] Karklin, Y. and Simoncelli, E. P. Efficient coding of natural images with a population of noisy linear-nonlinear neurons. In *Advances in neural information processing systems*, pp. 999–1007, 2011.

[31] Karpathy, A. Lessons learned from manually classifying cifar-10. *Published online at http://karpathy. github. io/2011/04/27/manually-classifying-cifar10*, 2011.

[32] Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.

[33] Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, 2012.

[34] Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. In *European conference on computer vision*, pp. 740–755. Springer, 2014.

[35] Lin, T.-Y., Goyal, P., Girshick, R., He, K., and Dollár, P. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, 2017.

[36] Liu, H., Simonyan, K., and Yang, Y. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018.

[37] Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.

[38] Poole, B., Sohl-Dickstein, J., and Ganguli, S. Analyzing noise in autoencoders and deep networks, 2014.

[39] Recht, B., Roelofs, R., Schmidt, L., and Shankar, V. Do cifar-10 classifiers generalize to cifar-10? *arXiv preprint arXiv:1806.00451*, 2018.

[40] Recht, B., Roelofs, R., Schmidt, L., and Shankar, V. Do imagenet classifiers generalize to imagenet? *arXiv preprint arXiv:1902.10811*, 2019.

[41] Rosenfeld, A., Zemel, R., and Tsotsos, J. K. The elephant in the room. *arXiv preprint arXiv:1808.03305*, 2018.

[42] Simoncelli, E. P. and Olshausen, B. A. Natural image statistics and neural representation. *Annual review of neuroscience*, 24(1):1193–1216, 2001.

[43] Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. *Advances in Neural Information Processing Systems*, 2015.

[44] Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A., et al. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.

[45] Szegedy, C., Ioffe, S., Vanhoucke, V., and Alemi, A. A. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, 2017.

[46] Torralba, A. and Oliva, A. Statistics of natural image categories. *Network: computation in neural systems*, 14(3):391–412, 2003.

[47] Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., and Madry, A. Robustness may be at odds with accuracy. *stat*, 1050:11, 2018.

[48] Wang, H., He, Z., Lipton, Z. C., and Xing, E. P. Learning robust representations by projecting superficial statistics out. *arXiv preprint arXiv:1903.06256*, 2019.

[49] Xie, Q., Dai, Z., Hovy, E., Luong, M.-T., and Le, Q. V. Unsupervised data augmentation, 2019.

[50] Yin, D., Gontijo Lopes, R., Shlens, J., Cubuk, E. D., and Gilmer, J. A fourier perspective on model robustness in computer vision. *ICML Workshop on Uncertainty and Robustness in Deep Learning*, 2019.

[51] Zagoruyko, S. and Komodakis, N. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.

[52] Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.

[53] Zhong, Z., Zheng, L., Kang, G., Li, S., and Yang, Y. Random erasing data augmentation. *arXiv preprint arXiv:1708.04896*, 2017.

[54] Zoph, B. and Le, Q. V. Neural architecture search with reinforcement learning. In *International Conference on Learning Representations*, 2017.

[55] Zoph, B., Vasudevan, V., Shlens, J., and Le, Q. V. Learning transferable architectures for scalable image recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
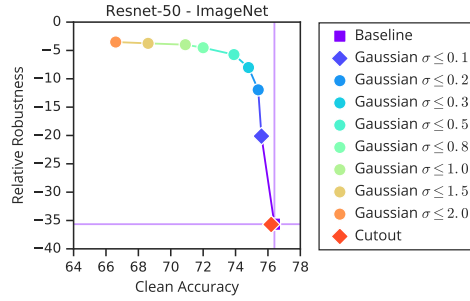
# Appendix



Figure 6: Accuracy/robustness trade-off observed for `Cutout` and `Gaussian` on Resnet-50 models. See Figure 2 for details.
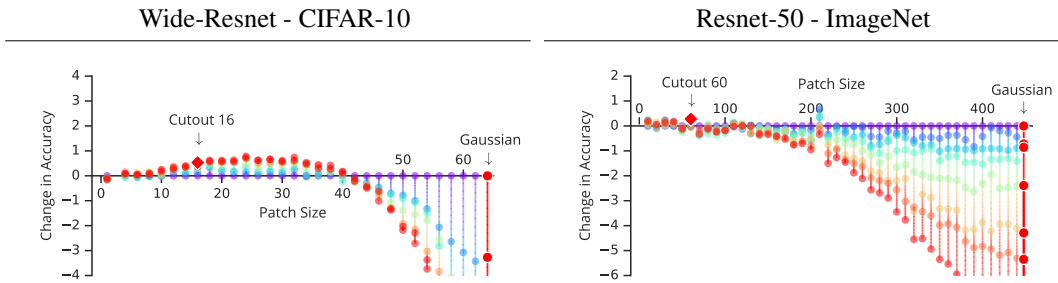


Figure 7: `Patch Gaussian` hyper-parameter sweep for Wide-Resnet on CIFAR-10 (left) and RN50 on Imagenet (right). `Patch Gaussian` approaches `Gaussian` with increasing patch size and `Cutout` with increasing $\sigma$. Each dot is a model trained with different hyper parameters. Colors indicate different $\sigma$.
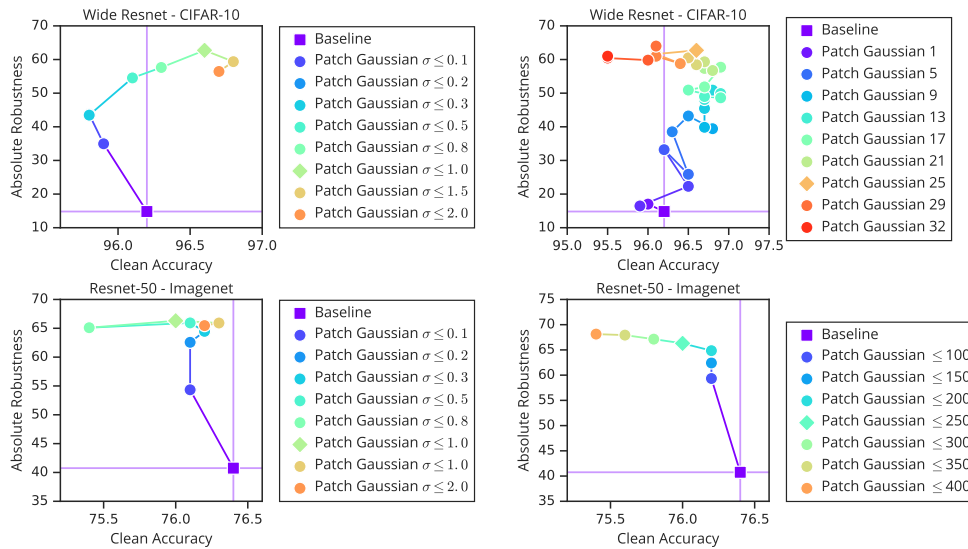


Figure 8: Overcoming the accuracy/robustness trade-off with `Patch Gaussian` for models trained on CIFAR-10 (top row) and Resnet-50 (bottom row). See figure 4 for details.
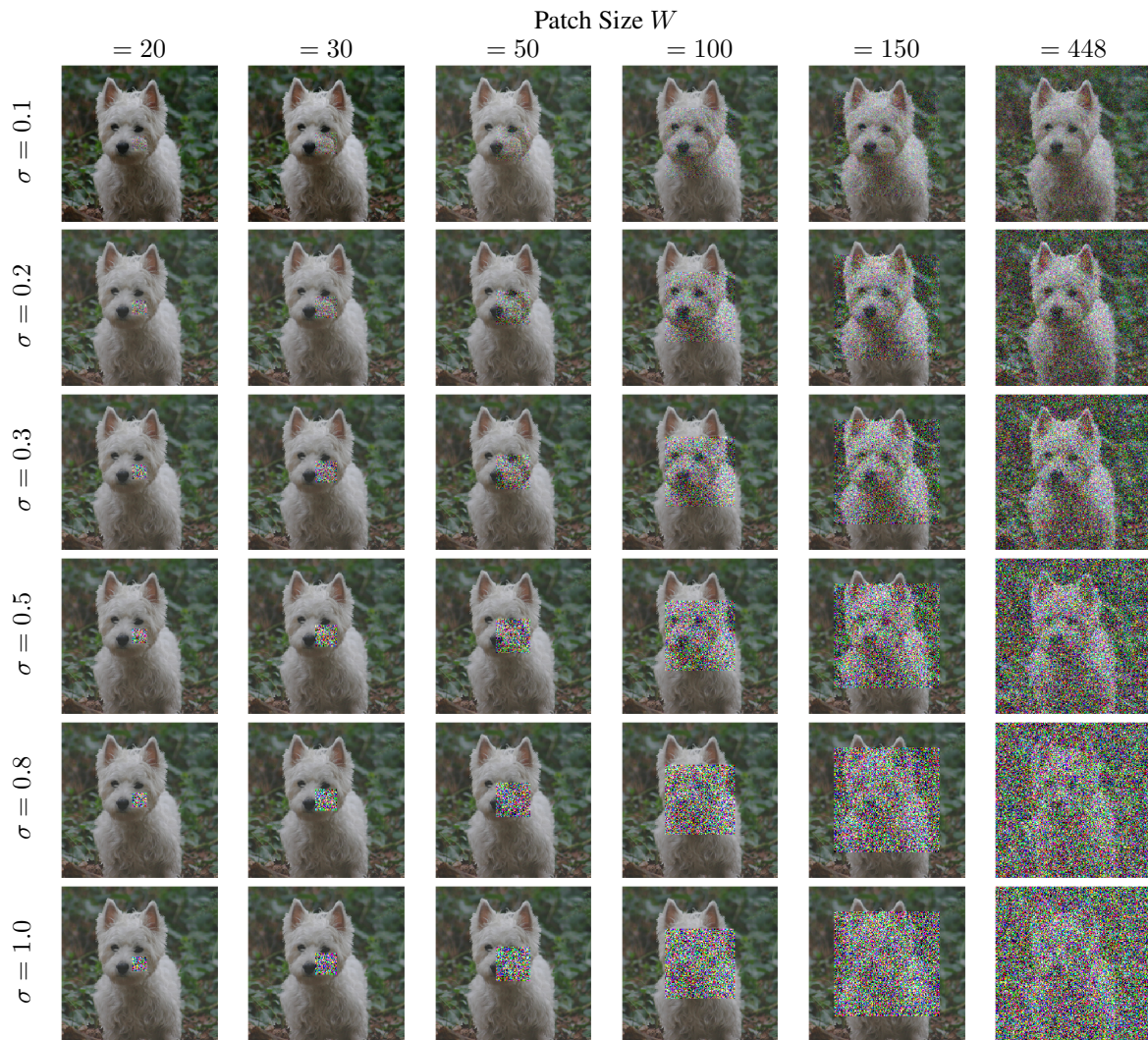
Patch Size $W$



Figure 9: Images modified with `Patch Gaussian`, with centered patch, at various $W$ & $\sigma$.

```python
def _get_patch_mask(patch_size):
  # randomly sample location in the image
  x = tf.random.uniform([], minval=0, maxval=224, dtype=tf.int32)
  y = tf.random.uniform([], minval=0, maxval=224, dtype=tf.int32)
  x, y = tf.cast(x, tf.float32), tf.cast(y, tf.float32)

  # compute where the patch will start and end
  startx, starty = x - tf.floor(patch_size/2), y - tf.floor(patch_size/2)
  endx, endy = x + tf.ceil(patch_size/2), y + tf.ceil(patch_size/2)
  startx, starty = tf.maximum(startx, 0), tf.maximum(starty, 0)
  endx, endy = tf.minimum(endx, 224), tf.minimum(endy, 224)

  # now let's convert these into how much we need to pad the patch
  lower_pad, upper_pad = 224 - endy, starty
  left_pad, right_pad = startx, 224 - endx
  padding_dims = [[upper_pad, lower_pad], [left_pad, right_pad]]

  # create mask
  mask = tf.pad(tf.zeros([endy - starty, endx - startx]),
                padding_dims, constant_values=1)
  mask = tf.expand_dims(mask, -1)
  mask = tf.tile(mask, [1, 1, 3])
  return tf.equal(mask, 0)


def patch_gaussian(image, patch_size, max_scale, sample_up_to):
  """Returns image with Patch Gaussian applied."""

  if sample_up_to:
    patch_size = tf.random.uniform([], 1, patch_size, tf.int32)
    # otherwise, patch_size is fixed.

  # make image (which is [0, 255]) be [0, 1]
  image = image / 255.0

  # uniformly sample scale from 0 to given scale
  scale = max_scale * tf.random.uniform([], minval=0, maxval=1)

  # apply gaussian to copy of image. Will be used to replace patch in image
  gaussian = tf.random.normal(stddev=scale, shape=image.shape)
  image_plus_gaussian = tf.clip_by_value(image + gaussian, 0, 1)

  # create mask and apply patch
  image = tf.where(_get_patch_mask(patch_size),
                   image_plus_gaussian, image)

  # scale back to [0, 255]
  return image * 255
```

Figure 10: TensorFlow implementation of Patch Gaussian

Table 6: Augmentation hyper-parameters selected with the method in Section 3.2 for each model/dataset. *Indicates manually-chosen stronger hyper-parameters, used to highlight the effect of the augmentation on the models. "$\leq$" indicates that the value is uniformly sampled up to the given maximum value.

| | | Z | Augmentation | W | $\sigma$ | Other |
|---|---|---|---|---|---|---|
| CIFAR-10 | Wide Resnet-28-10 | 96.5% | Cutout | = 12 | - | |
| | | | Gaussian | - | $\leq$ 0.1 | |
| | | | Patch Gaussian | = 25 | $\leq$ 1.0 | |
| | | | Cutout* | = 22 | - | |
| | | | Gaussian* | - | $\leq$ 1.0 | |
| | Shake 112 | 97.0% | Cutout | = 7 | - | |
| | | | Gaussian | - | $\leq$ 0.1 | |
| | | | Patch Gaussian | = 26 | $\leq$ 1.0 | |
| ImageNet | Resnet-50 | 76.0% | Baseline | - | - | includes weight decay = 0.0001 |
| | | | Cutout | = 60 | - | |
| | | | Gaussian | - | $\leq$ 0.1 | |
| | | | Patch Gaussian | $\leq$ 250 | $\leq$ 1.0 | |
| | | | Cutout* | = 200 | - | |
| | | | Gaussian* | - | $\leq$ 1.0 | |
| | | | Larger Weight Decay | - | - | 0.001 |
| | | | Dropblock | - | - | groups = 3,4; keep prob = 0.9 |
| | | | Label Smoothing | - | - | 0.1 |
| | Resnet 200 | 78.5% | Baseline | - | - | includes weight decay = 0.0001 |
| | | | Cutout | = 30 | - | |
| | | | Gaussian | - | $\leq$ 0.1 | |
| | | | Patch Gaussian | $\leq$ 350 | $\leq$ 1.0 | |

Table 7: Full original corruption errors (Original CEs) for ImageNet models trained with different augmentation strategies.

| | Augmentation | Noise | | | Blur | | | |
|---|---|---|---|---|---|---|---|---|
| | | Gaussian | Shot | Impulse | Defocus | Glass | Motion | Zoom |
| Resnet-50 | Baseline | 0.705 | 0.722 | 0.716 | 0.815 | 0.915 | 0.810 | 0.817 |
| | Cutout | 0.720 | 0.727 | 0.720 | 0.798 | 0.923 | 0.821 | 0.813 |
| | Gaussian | 0.677 | 0.681 | 0.677 | 0.781 | **0.864** | 0.813 | 0.808 |
| | Patch Gaussian | **0.623** | **0.633** | **0.624** | **0.751** | 0.898 | **0.782** | **0.783** |
| Resnet-200 | Baseline | 0.622 | 0.641 | 0.629 | 0.735 | 0.867 | 0.722 | 0.739 |
| | Cutout | 0.594 | 0.619 | 0.600 | 0.714 | 0.870 | 0.713 | 0.737 |
| | Gaussian | 0.573 | 0.583 | 0.575 | 0.723 | 0.814 | 0.737 | 0.741 |
| | Patch Gaussian | **0.486** | **0.498** | **0.478** | **0.649** | **0.805** | **0.693** | **0.687** |

| | Augmentation | Weather | | | | Digital | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Snow | Frost | Fog | Bright | Contrast | Elastic | Pixel | JPEG |
| Resnet-50 | Baseline | 0.827 | 0.756 | 0.589 | **0.582** | 0.748 | 0.753 | 0.799 | 0.747 |
| | Cutout | 0.839 | 0.764 | 0.599 | 0.586 | 0.747 | 0.752 | 0.803 | 0.752 |
| | Gaussian | 0.821 | **0.726** | 0.597 | 0.592 | 0.754 | **0.720** | 0.805 | 0.763 |
| | Patch Gaussian | **0.806** | 0.739 | **0.566** | 0.592 | **0.714** | 0.736 | **0.743** | **0.722** |
| Resnet-200 | Baseline | 0.754 | 0.694 | 0.497 | 0.520 | 0.658 | 0.669 | 0.696 | 0.681 |
| | Cutout | 0.741 | 0.684 | 0.507 | 0.516 | 0.671 | 0.670 | 0.751 | 0.672 |
| | Gaussian | 0.731 | 0.653 | 0.525 | 0.514 | 0.699 | 0.641 | 0.693 | 0.660 |
| | Patch Gaussian | **0.697** | **0.633** | **0.476** | **0.506** | **0.627** | **0.625** | **0.613** | **0.593** |

Table 8: Full corruption errors (CEs) for ImageNet models trained with different augmentation strategies.

|  | Augmentation | Noise | | | Blur | |
|---|---|---|---|---|---|---|
|  |  | Gaussian | Shot | Impulse | Defocus | Zoom |
| Resnet-50 | Baseline | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 |
|  | Cutout | 1.015 | 1.013 | 1.008 | 0.979 | 1.000 |
|  | Gaussian | 0.620 | 0.625 | 0.618 | 0.950 | 0.999 |
|  | Patch Gaussian | **0.585** | **0.577** | **0.577** | **0.922** | **0.963** |
| Resnet-200 | Baseline | 0.872 | 0.883 | 0.864 | 0.880 | 0.896 |
|  | Cutout | 0.841 | 0.862 | 0.833 | 0.866 | 0.892 |
|  | Gaussian | 0.533 | 0.538 | 0.538 | 0.855 | 0.910 |
|  | Patch Gaussian | **0.490** | **0.488** | **0.490** | **0.767** | **0.820** |

|  | Augmentation | Weather | | | Digital | | | |
|---|---|---|---|---|---|---|---|---|
|  |  | Frost | Fog | Bright | Contrast | Elastic | Pixel | JPEG |
| Resnet-50 | Baseline | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 |
|  | Cutout | 1.005 | 1.017 | 0.991 | 1.008 | 1.009 | 1.026 | 1.011 |
|  | Gaussian | **0.919** | 1.073 | 1.019 | 1.051 | **0.967** | 0.974 | **0.966** |
|  | Patch Gaussian | 0.976 | **0.978** | **0.990** | **0.956** | 0.982 | **0.957** | 0.998 |
| Resnet-200 | Baseline | 0.912 | 0.862 | 0.888 | 0.861 | 0.882 | 0.848 | 0.922 |
|  | Cutout | 0.915 | 0.868 | 0.877 | 0.877 | 0.871 | 0.875 | 0.911 |
|  | Gaussian | 0.830 | 0.948 | 0.889 | 0.960 | 0.848 | 0.836 | 0.855 |
|  | Patch Gaussian | **0.818** | **0.851** | **0.862** | **0.832** | **0.812** | **0.765** | **0.835** |

Table 9: Comparison with SIN+IN [15]. By using Z=74.6%, `Patch Gaussian` can match SIN+IN's og mCE and test accuracy. Understandably, however, our gains are more concentrated in noise-based corruptions, whereas shape-biased models get gains in other corruptions.

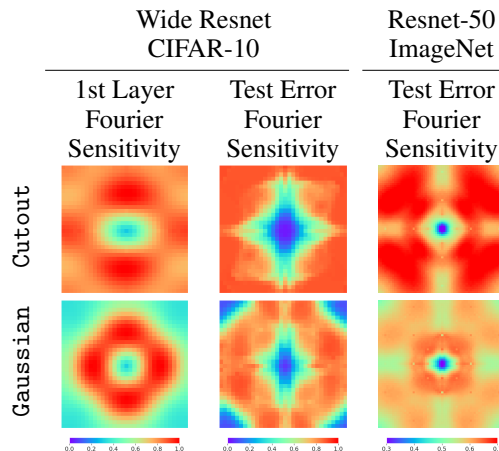|  | Augmentation | Test Accuracy | og mCE | og mCE (-noise) |
|---|---|---|---|---|
| Resnet-50 | SIN+IN | 74.6% | **0.693** | **0.699** |
|  | Patch Gaussian ($W \leq 400$, $\sigma \leq 0.8$) | **75.6%** | **0.693** | 0.718 |



Figure 11: Fourier analysis for `Cutout` and `Gaussian` models selected by the method in Section 3.2. See Figure 5 for details.
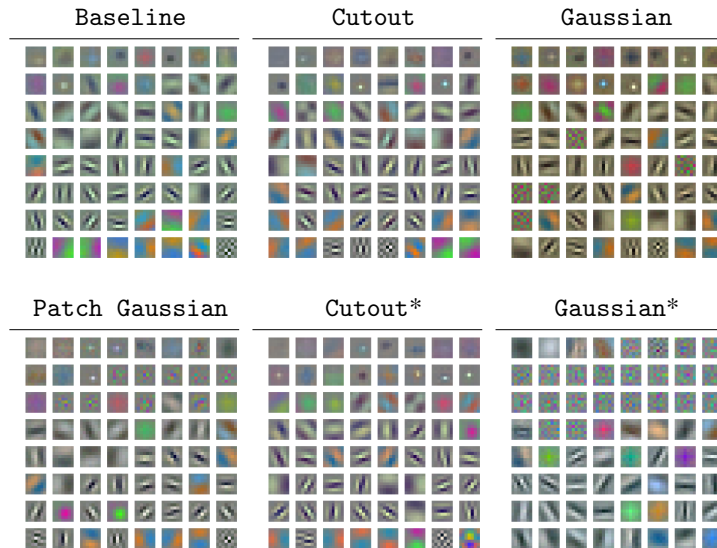
Figure 12: Complete filters for Resnet-50 models trained on ImageNet. * Indicates augmentations with larger patch sizes and $\sigma$. See Figure 5 for details. We again note the presence of filters of high fourier frequency in models trained with `Cutout*` and `Patch Gaussian`. We also note that `Gaussian*` exhibits high variance filters. We posit these have not been trained and have little importance, given the low sensitivity of this model to high frequencies. Future work will investigate the importance of filters on sensitivity.
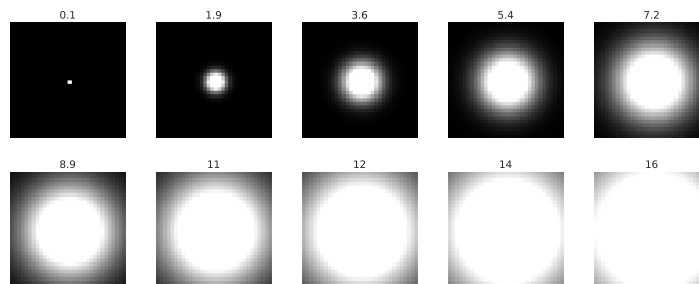


Figure 13: Examples of high pass filters at various radii, in fourier space centered at the zero-frequency component, used in the high-pass experiment of Figure 5.