# Robust Contrastive Learning Using Negative Samples with Diminished Semantics

**Songwei Ge**
Univeristy of Maryland
songweig@cs.umd.edu

**Shlok Mishra**
Univeristy of Maryland
shlokm@cs.umd.edu

**Haohan Wang**
Carnegie Mellon University
haohanw@cs.cmu.edu

**Chun-Liang Li**
Google Cloud AI
chunliang@google.com

**David Jacobs**
Univeristy of Maryland
dwj@cs.umd.edu

## Abstract

Unsupervised learning has recently made exceptional progress because of the development of more effective contrastive learning methods. However, CNNs are prone to depend on low-level features that humans deem non-semantic. This dependency has been conjectured to induce a lack of robustness to image perturbations or domain shift. In this paper, we show that by generating carefully designed negative samples, contrastive learning can learn more robust representations with less dependence on such features. Contrastive learning utilizes positive pairs that preserve semantic information while perturbing superficial features in the training images. Similarly, we propose to generate negative samples in a reversed way, where only the superfluous instead of the semantic features are preserved. We develop two methods, texture-based and patch-based augmentations, to generate negative samples. These samples achieve better generalization, especially under out-of-domain settings. We also analyze our method and the generated texture-based samples, showing that texture features are indispensable in classifying particular ImageNet classes and especially finer classes. We also show that model bias favors texture and shape features differently under different test settings. Our code, trained models, and ImageNet-Texture dataset can be found at https://github.com/SongweiGe/Contrastive-Learning-with-Non-Semantic-Negatives.

## 1 Introduction

Recent studies on self-supervised learning have shown great success in learning visual representations without human annotations. The gap between unsupervised and supervised learning has been progressively closed by contrastive learning [52, 47, 7, 18, 48, 6, 16, 55]. In the meantime, CNNs trained in the supervised setting are known to learn correlations between labels and superfluous features such as local patches [4, 3], texture [14], high-frequency components [50], and even artificially added features [26], which has raised concerns about deploying these models in a real scenario [30, 15]. CNNs trained by contrastive learning methods are no exception [23]. In this paper, we propose to construct negative samples that only preserve non-semantic features. We show that using contrastive learning methods trained with these negative samples can mitigate these concerns.

Contrastive learning methods exploit carefully designed augmentations to construct positive pairs and pull their representations together. These augmentations are crucial to contrastive learning [7, 6]. A common assumption behind these augmentations is to preserve the semantics of the input images while perturbing other superficial signals. This inspires us to generate negative samples and inject additional implicit biases on the visual features learned by the models. Specifically, we utilize

Semantic positive samples · Non-semantic negative samples

(a) Input image · (b) Query sample · (c) Positive sample · (d) Texture-based NS · (e) Patch-based NS
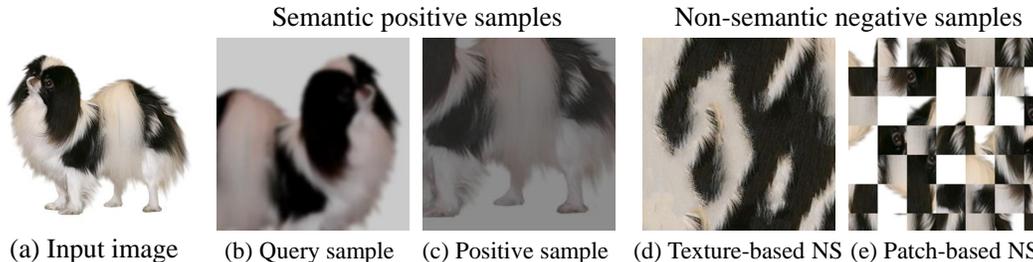
Figure 1: We propose to construct negative samples (NS) from input images for contrastive learning with augmentations that only preserve non-semantic information such as texture and local features.

augmentations that diminish the semantic features while keeping the undesired features such as texture. By pushing apart the representations of such negative samples and input images, the models are expected to rely more on the semantics of the images and less on superficial features.

Inspired by the non-semantic features, we propose two methods to craft negative samples. The first method relies on texture synthesis tools from classic approaches [12, 51]. It generates realistic texture images based on two patches extracted from input images, as shown in Figure 1(d). For each image in ImageNet, we generate its texture version and form a dataset which we call ImageNet-Texture. The second method constructs non-semantic images by tiling randomly sampled patches of different sizes from the input image, as shown in Figure 1(e). Comparing the non-semantic negative samples with two semantic positive samples in Figure 1(b) and Figure 1(c), the dog from the input image is still recognizable in the positive samples but hard to understand from negative samples. Instead, local statistics such as the fur and color of the dog are preserved in the negative samples.

The generated non-semantic samples can be readily used with existing contrastive learning methods that distinguish positive pairs from negative pairs such as MoCo [18] and SimCLR [7]. Despite their simplicity, we show that these non-semantic negative samples are actually harder than the standard negative samples used by these methods, which are inefficient at leveraging hard negatives [13, 31]. Further, our negative pairs can also be used by contrastive learning methods that do not explicitly use negative samples, such as BYOL [16]. We evaluate our methods with two contrastive learning methods, MoCo [18, 8] and BYOL [16], on three datasets, ImageNet-100, ImageNet-1K and STL-10. When using our proposed augmentations to generate negative samples and minimize their representation similarity to the input images, we notice a consistent improvement on the generalization performance over backbone methods [8, 16] and previous negative example generation strategies [31, 43], especially under out-of-distribution (OOD) settings.

We conduct a systematic analysis of how the shape-texture trade-off influences model performance based on the proposed ImageNet-Texture dataset. We control the penalty on similarities between the non-semantic negative examples and the query samples. This impacts the trade-off between using shape and texture features. We find that the relative importance of texture and shape features varies across different datasets. For example, shape bias benefits ImageNet-Sketch [49] more than the original ImageNet validation set. On the other hand, texture bias benefits finer-grained classification more, such as dog breed classification included in ImageNet. These results complement previous evidence showing the effectiveness of shape features in classifying 16 coarse classes [14]. Such preference for one feature over the other is also observed intra-dataset: the texture is more important for some classes such as dishrag and plaque. These observations make us question the relationship between shape and texture features as the implicitly necessary bias of CNNs and advocate for an adaptive combination of both when deploying the model in real scenarios. In summary:

- We propose texture-based and patch-based augmentations to generate negative samples from input images, and show that these negative samples improve the generalization of contrastive learning.

- We introduce the ImageNet-Texture dataset, which contains texture versions of ImageNet images generated by texture synthesis tools.

- We provide fine-grained analysis on the shape-texture trade-off of CNNs, and show different scenarios when one is preferred over the other.

2

# 2  Negative Samples with Diminished Semantics

CNNs are apt to learn low level features such as texture under supervised settings [3, 14, 50]; this has been recently witnessed under the contrastive learning setting as well [23]. To mitigate this problem, we propose two methods, texture-based and patch-based augmentations, to generate negative samples for contrastive learning. Texture-based augmentation generates realistic images based on texture synthesis and patch-based augmentation exploits more comprehensive local features by sampling patches from input images. By penalizing learned similarities between the representations of images and their non-semantic counterparts, the model is encouraged to rely less on the undesired features and focus more on the semantics. In practice, we find the two negative samples play similar roles and the patch-based method works slightly better. In this section, we start with an overview of contrastive learning and show how non-semantic negatives are used in these frameworks. Then we elaborate on the two approaches to generate negative samples with diminished semantics.

## 2.1  Contrastive learning with non-semantic negatives

Given an encoder network $f$ and an image $\mathbf{x}$, we denote the output of the network as $\mathbf{z} = f(\mathbf{x})$. We use $z_i$ and $z_p$ to denote the representations of the query sample $x_i$ and a positive sample $x_p$ generated from the same input image with augmentations that preserve semantics. For contrastive learning methods like MoCo [18] and SimCLR [7], $z_n$ denotes the representation of the standard negative sample $x_n$ extracted from the memory bank (MoCo) or other images in the current batch (SimCLR). $z_{ns}$ is the representation of the proposed negative sample $x_{ns}$ which contains particular non-semantic features of the input image with the semantic part weakened. We extend the noise-contrastive estimation (NCE) loss as below:

$$\mathcal{L}_{\text{NCE}} = -\sum_{i \in I} \log \frac{\exp\left(z_i^T z_p / \tau\right)}{\exp\left(z_i^T z_p / \tau\right) + \exp\left(\alpha z_i^T z_{ns} / \tau\right) + \sum_{n \in \mathcal{N}} \exp\left(z_i^T z_n / \tau\right)}, \tag{1}$$

where $\tau$ is a temperature parameter and $\alpha$ is an additional scaling parameter for non-semantic negatives. A larger $\alpha$ implies a stronger penalty on the similarity between the representations of the query image and its non-semantic version. In Appendix B.1 we discuss other possible ways to apply $\alpha$.

Methods like BYOL [16] do not explicitly rely on negative samples. Nevertheless, BYOL adapts the loss to maximize the agreement of positive pairs. Therefore, we explicitly use the non-semantic negative sample with their loss to minimize its similarity to the query sample:

$$\mathcal{L}_{\text{BYOL}} = \|z_i - z_p\| - \alpha \|z_i - z_{ns}\| = 2 - 2\alpha - 2z_i^T z_p + 2\alpha z_i^T z_{ns}. \tag{2}$$

We overload $\alpha$ to be the parameter that controls the penalty on the similarity between the representations of input image and its non-semantic version under BYOL, with similar intention as MoCo and SimCLR above. To minimize either $\mathcal{L}_{\text{NCE}}$ or $\mathcal{L}_{\text{BYOL}}$, the encoder must learn features from $x_i$ that are not contained in $x_{ns}$ but shared with $x_p$.

## 2.2  Texture-based negative sample generation

We use texture synthesis tools to generate negative samples. Texture synthesis aims to generate realistic images that preserve as much local structure as possible from an example image [19, 11, 42]. For instance, as shown in Figure 1(d), the texture of the input dog image preserves the fur and colors of the dog. Notably, in previous discussion of robustness [3, 49], such local structure has often been recognized as highly correlated with the labels yet superfluous to generalization. For example, under large domain shift due to lighting, motion, and even modality, the texture is more apt to change than the semantic features, such as the shape. Furthermore, CNNs trained on ImageNet are more likely to classify images based on the texture features rather than the shape features which are instead preferred by humans due to their transferability [14]. To encourage the model to rely more on the shape features, we propose a two-step method to generate the texture image of input images as the negative samples for contrastive learning.

To be specific, we first sample two patches from given images as the input to the texture synthesis algorithms. One patch is extracted from the center of the image. This patch is expected to reflect the texture of the object according to the implicit bias contained in the ImageNet dataset that most

of the objects are center-oriented in the images [2]. The other patch is extracted from a random location to reflect other possible textures of the image (e.g. background, peripheral region of the object). In this work, we extract patches with size $96 \times 96$ when image size allows, otherwise $48 \times 48$ patches are extracted. Second, we adopt off-the-shelf texture synthesis algorithms [12, 51, 1] to generate texture images based on the two patches. These non-parametric algorithms iteratively sample pixels from given patches that share a similar neighborhood with the current pixel. Specifically, we use the open-source software built on these methods [12, 51, 1] with multi-threaded CPU support implemented in Rust [1]. For each sample in the ImageNet dataset, we generate one $224 \times 224$ texture image to construct a dataset that has the same training and validation size as the ImageNet dataset. We call this dataset ImageNet-Texture. More examples can be found in Appendix A.1.

### 2.3    Patch-based negative sample generation

To simulate the local information contained in the images [4, 3], we propose an efficient patch-based method to generate non-semantic images. Given an image and a patch size $d$, we sample patches of size $d$ from $(\lceil \frac{224}{d} \rceil)^2$ non-overlapping random locations that lie entirely in the image. The patches are then tiled and cropped into $224 \times 224$ as negative samples. Compared with the texture-based method, this generation process takes negligible time, therefore it can be implemented as part of the data loading process in parallel with training. By doing so, each training sample can be paired with different negative samples generated from different patches every time it is used, compared with the two fixed patches selected when generating texture images.

Different from the texture-based method that generates realistic images, the patch-based method generates images with artificial lines as shown in Figure 1e. One might be concerned with possible degenerate solution where the model outputs a low similarity whenever it detects the repeated sharp changes in the horizontal or vertical directions, which could be done with a single layer of convolution. However, interestingly, we find that the model does not find such a simple solution in practice. This is also noticed in a previous study where the image and its copy with a patch cut out are non-trivially distinguished by the model [34]. To mitigate this potential issue, we randomly sample patch size $d$ from a prior distribution instead of using a fixed $d$ in practice, which allows the model to look at texture at different scales.

### 2.4    How hard are the texture-based and patch-based negative samples?

Constrastive learning methods are known to struggle with finding hard negative samples [13, 31] and researchers have proposed several ways to better leverage hard negatives [31, 43]. An intermediate question is how hard are our proposed negative samples compared with those standard negative samples used in previous constrastive learning methods [7, 18], i.e. random training samples.



(a) Positive Sample    (b) Standard NS    (c) Texture-based NS    (d) Patch-based NS
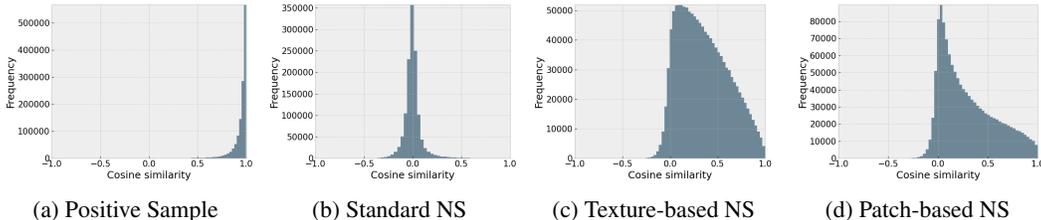
Figure 2: The histogram of cosine similarity between the representations of query sample and its paired samples, namely positive sample and standard, texture-based, and patch-based negative samples (NS), using MoCo-v2 model trained on the ImageNet-1K dataset for 200 epochs.

We use the official MoCo-v2 model pretrained on the ImageNet-1K dataset for 200 epochs to calculate the cosine similarities between different kinds of pairs across the ImageNet training set. We plot the histogram of these similarities in Figure 2. As shown in the Figures 2a and 2b, most positive pairs and negative pairs have similarity close to $1$ and $0$ respectively. Specifically, the average similarities across the training samples are $0.94257$ and $0.00018$ for positive and negative pairs. As shown in

the Figures 2d and 2c, the distributions of patch-based and texture-based negative samples are very different from those of standard negative samples; their similarity distributions have heavy tails in the positive region. Specifically, the distribution of patch-based and texture-based negative samples have average similarity 0.29503 and 0.35248 across the dataset, which shows that they remain difficult after training with standard negative examples.

## 3 Experiments

In this section, we evaluate the two kinds of non-semantic negative samples with two contrastive learning methods, MoCo and BYOL, on the ImageNet dataset. We also experiment using its subset, the ImageNet-100 dataset [47, 31], which allows us to perform more comprehensive experiments. We report accuracy on out-of-domain (OOD) datasets including the ImageNet-C(orruption) [22], ImageNet-S(ketch) [49], Stylized-ImageNet [14], and ImageNet-R(endition)[21] datasets as an evaluation of the model's robustness to domain shifts. ImageNet-C and Stylized-ImageNet contain images transformed from the images in the ImageNet validation set with common corruption and transferred style. ImageNet-S and ImageNet-R are collected independently from the ImageNet dataset and share all or a subset of the classes in the ImageNet dataset with a focus on sketch and other rendering modalities. We show that with our proposed non-semantic negatives, contrastive learning generalizes better under domain shifts. For patch-based negatives, it also improves the performance on the in-domain dataset.

### 3.1 ImageNet-100

|  | ImageNet | ImageNet-C | ImageNet-S | Stylized-ImageNet | ImageNet-R |
|---|---|---|---|---|---|
| MoCo-v2 [8] - $k = 16384$ | $77.88_{\pm 0.28}$ | $43.08_{\pm 0.27}$ | $28.24_{\pm 0.58}$ | $16.20_{\pm 0.55}$ | $32.92_{\pm 0.12}$ |
| + Texture-based - $\alpha = 2$ | $77.76_{\pm 0.17}$ | $43.58_{\pm 0.33}$ | $29.11_{\pm 0.39}$ | $16.59_{\pm 0.17}$ | $33.36_{\pm 0.15}$ |
| + Patch-based - $\alpha = 2$ | $\mathbf{79.35}_{\pm 0.12}$ | $\mathbf{45.13}_{\pm 0.35}$ | $31.76_{\pm 0.88}$ | $17.37_{\pm 0.19}$ | $34.78_{\pm 0.15}$ |
| + Patch-based - $\alpha = 3$ | $75.58_{\pm 0.52}$ | $44.45_{\pm 0.15}$ | $\mathbf{34.03}_{\pm 0.58}$ | $\mathbf{18.60}_{\pm 0.26}$ | $\mathbf{36.89}_{\pm 0.11}$ |
| MoCo-v2 [8] - $k = 8192$ | $77.73_{\pm 0.38}$ | $43.22_{\pm 0.39}$ | $28.45_{\pm 0.36}$ | $16.83_{\pm 0.12}$ | $33.19_{\pm 0.44}$ |
| + Patch-based - $\alpha = 2$ | $\mathbf{79.54}_{\pm 0.32}$ | $\mathbf{45.48}_{\pm 0.20}$ | $\mathbf{33.36}_{\pm 0.45}$ | $17.81_{\pm 0.32}$ | $\mathbf{36.31}_{\pm 0.37}$ |
| BYOL [16] | $78.76_{\pm 0.28}$ | $44.43_{\pm 0.35}$ | $35.84_{\pm 0.38}$ | $15.01_{\pm 0.19}$ | $39.53_{\pm 0.51}$ |
| + Patch-based - $\alpha = 0.05$ | $\mathbf{78.81}_{\pm 0.33}$ | $\mathbf{44.60}_{\pm 0.21}$ | $\mathbf{36.76}_{\pm 0.51}$ | $\mathbf{15.52}_{\pm 0.22}$ | $\mathbf{41.16}_{\pm 0.39}$ |
| InsDis [52] | 68.52 | 28.93 | 16.67 | 9.86 | 19.60 |
| CMC [47] | 79.34 | 39.28 | 24.04 | 13.88 | 32.68 |
| InfoMin [48] | 82.74 | 48.87 | 38.43 | 18.14 | 40.68 |
| Supervised | 86.26 | 49.17 | 34.95 | 21.20 | 39.76 |

Table 1: Top-1 accuracy on the ImageNet-100 dataset and its OOD variants. We consider the supervised baseline as well as several self-supervised baselines including MoCo-v2, BYOL, InsDis, CMC, and InfoMIN. For our main comparison using MoCo-v2 and BYOL, we also report the standard deviation of 3 runs. For MoCo models, $k$ represents the size of the memory bank .

We first follow the hyperparameters used in [31] to train MoCo-v2 on the ImageNet-100 dataset with a memory bank size $k = 16384$. We also use a similar hyperparameter configuration in which only memory bank size is halved. For patch-based augmentation parameters, we use patch size sampled from a uniform distribution $d \sim \mathcal{U}(16, 72)$. The parameter $\alpha$ is indicated behind each model name. We discuss the impact of $\alpha$ in detail in the next section. More ablations on the patch-based augmentations can be found in Appendix C.3. For ImageNet-C, we report the average accuracy across 5 levels of corruption severity.

Similar to [31], we repeat the experiments for 3 runs and report the mean and standard deviation in Table 1. As shown in the table, when following the previous memory bank size, using both patch-based and texture-based negatives improve the OOD generalizations. Specifically, patch-based augmentation increases the accuracy on ImageNet-S by $5.79\%$ and ImageNet-R by $5.97\%$ when $\alpha = 3$. When $\alpha = 2$, it also increases the in-domain accuracy by $1.47\%$ and accuracy on ImageNet-C by $2.05\%$. The similar trend shared by standard ImageNet and ImageNet-C with different $\alpha$ can be attributed to the resemblance of the images in the two dataset, especially those corrupted images with a lower level of severity. We show the performance of the model with $\alpha = 3$ is actually better on the

5

highest corruption level as shown in Appendix C.2. The improvement achieved using texture-based negatives is less, probably because the information contained in the texture image is restricted due to the limited access to the two fixed patches. When the memory bank is halved to be $8,192$, the baseline MoCo model has slightly worse performance, decreasing from $77.88$ to $77.73$. But with patch-based hard negative samples, the MoCo-v2 model instead achieved the best accuracy $79.54$ on the ImageNet-100 validation set, IamgeNet-C and IamgeNet-R. We discuss this more in a later section.



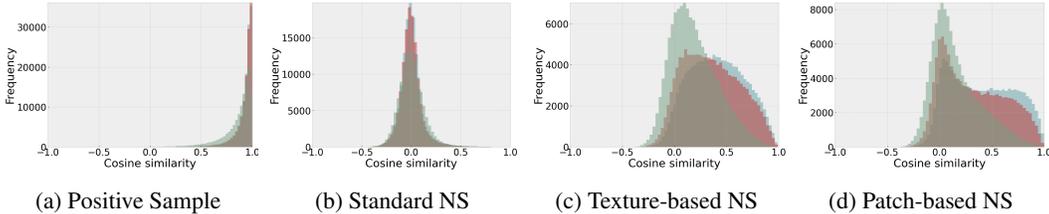(a) Positive Sample   (b) Standard NS   (c) Texture-based NS   (d) Patch-based NS

Figure 3: Histogram of cosine similarity between the representations of query sample and its paired positive sample, standard, texture-based, and patch-based negative samples (NS), using models trained without (blue) and with patch-based negative samples (red: $\alpha = 2$, green: $\alpha = 3$).

Similar to Figure 2, in Figure 3 we visualize the distribution of the cosine similarity between the query sample and both semantic and non-semantic samples calculated based on the models trained with and without patch-based negative samples. The shift of the distribution towards the origin in Figure 3d meets the expectation that our method reduces the similarity of input images and patch-based negative samples. Specifically, the average similarity decreases from $0.4040$ to $0.3252$ to $0.1593$ when $\alpha$ increases from 0, namely no patched-based negatives, to 2 to 3. Interestingly, we notice in Figure 3c that the similarity to texture-based negative samples also decreases, and the average similarity decreases from $0.4114$ to $0.3541$ to $0.1896$, although we did not explicitly penalize it. This demonstrates the resemblance of patch-based negative samples to texture-based negative samples. Given the better performance achieved with patch-based negative samples, for the rest of the experiments, we mainly focus on the patch-based methods. But we still conduct our analysis on the texture-based samples. We also find a marginal decrease in the positive similarity ($-0.0068$) and negative similarity ($-0.0008$) when $\alpha = 2$ and a substantial decrease in the positive similarity ($-0.0737$) and increase in the negative similarity ($0.0036$) when $\alpha = 3$. A similar figure for texture-based negative samples can be found in the Appendix in Figure 12.

## 3.2 ImageNet-1K

We follow the official hyperparameters [8] to train MoCo-v2 with our patch-based negative samples on the ImageNet-1K dataset. For the parameter $\alpha$ and patch size $d$, we follow the same configuration used on the ImageNet-100 dataset. We compare our results against the MoCo-v2 baseline [8] and the hard negative mixing algorithm, MoCHi [31]. Due to limited computational resources, we report the metrics evaluated with the official model without repeated runs. The results are shown in Table 2. More results can be found in Appendix C.6. Note that for ImageNet-C, we show the mCE metric [22], for which smaller is better. For the other datasets, we show the top-1 accuracy.

|  | ImageNet | ImageNet-C | ImageNet-S | Stylized-ImageNet | ImageNet-R |
|---|---|---|---|---|---|
| MoCo-v2 [8] | 67.60 | 87.7 | 17.47 | 5.55 | 27.81 |
| + MoCHi [31] | 67.56 | 88.7 | 16.32 | 5.94 | 25.71 |
| + Patch-base NS - $\alpha = 2$ | **67.92** | **87.6** | **18.58** | **6.34** | **28.95** |

Table 2: Top-1 accuracy on the ImageNet-1K dataset and its sketch, stylized, rendition variants, and mCE on the ImageNet-C dataset.

6

### 3.3 Extension to other non-semantic features

Non-semantic features are sometimes referred to as "shortcuts" in the contrastive learning literature [6, 7]. Models that leverage such features often exhibit unfavorable generalization to downstream tasks. For example, without color jittering, SimCLR [9] tends to utilize color histograms to reduce the training loss. In this section, we show that models trained with non-semantic negatives are coerced to avoid the shortcuts shared between query images and their non-semantic counterparts. In the example of the color shortcut, we note that the expected color distribution of our patch-based negatives is identical to that of the query images, and the actual distribution of samples is close. We conduct experiments with MoCo-v2 on the ImageNet-100 dataset while removing the color jittering from the augmentations. The accuracy of models with and without patch-based negatives are reported in Table 3. We found that patch-based negatives contribute significant effectiveness in preventing the models from learning such a color distribution shortcut.

| Model | Top-1 Accuracy |
|---|---|
| MoCo-v2 [8] | 70.44 |
| + Patch-based NS | 76.42 |

Table 3: Test accuracy of MoCo-v2 on the ImageNet-100 dataset after removing color jittering and adding patch-based negatives.

### 3.4 Memory bank size

Contrastive learning methods based on negative samples suffer from ineffective excavation of hard negatives [13, 31] and resort to large batch sizes [7] or memory bank [18]. In this section, we study whether our proposed negative samples can mitigate this problem on the STL-10 and ImageNet-100 datasets. We keep the hyperparameters intact and vary the memory bank size. We report the accuracy of the MoCo-v2 baseline with and without patch-based negative samples on STL-10 dataset in Table 4. We also compare with [43] which exploits hard negatives through reweighting. We found that with proper hyperparameters the MoCo-v2 baseline already beats the reweighting results with SimCLR.

| | |
|---|---|
| SimCLR [7] | 80.16 |
| + Debiased [10] | 84.90[§] |
| + Hard [43] | 87.42[§] |
| MoCo-v2 [8] | 88.00 |
| + Patch-based NS | 89.36 |

Table 4: Top-1 accuracy on the STL-10 dataset.
[§] denotes results visually extracted from Figure 2 in [43].
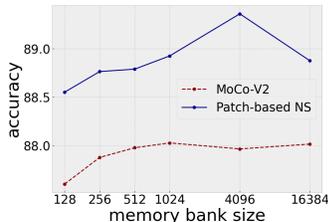


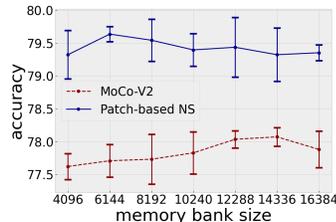Figure 4: Top-1 accuracy on the STL-10 dataset with different memory bank sizes.



Figure 5: Top-1 accuracy on the ImageNet-100 dataset with different memory bank sizes.

As shown in Figures 4 and 5, using patch-based non-semantic negatives consistently improves the MoCo baseline when the number of standard negatives varies. When slightly decreasing the memory bank size on the STL-10 dataset as shown in Figure 4 and ImageNet-100 in Table 1, the performance with patch-based negatives increases. This is probably because, according to Eq. 1 and analysis in Appendix B.1, a smaller memory bank size causes a larger contribution of non-semantic negatives to the loss, and consequently a larger regularization. To further demystify this observation, we conduct experiments with evenly sampled memory bank sizes between 4096 and 16384 and report the average and standard deviation across 3 runs in Figure 5. We confirm a consistent recession of baseline accuracy when decreasing memory back sizes [18, 43]. However, the steady improvement led by the non-semantic negatives effectively mitigates the problem - the decrease caused by a smaller memory bank is less substantial and using patch-based negatives always beats the baseline.

## 4 Discussion

### 4.1 Controlling the shape-texture trade-off with $\alpha$

There has been a growing interest in understanding the cause and impact of the trade-off between shape and texture bias of CNNs [14, 23, 35, 27]. CNNs trained on ImageNet are known to be

over-reliant on the texture features [14]. Contrastive learning with our non-semantic negatives serve as not only an effective method to reduce such reliance, but a natural tool to study such a trade-off.
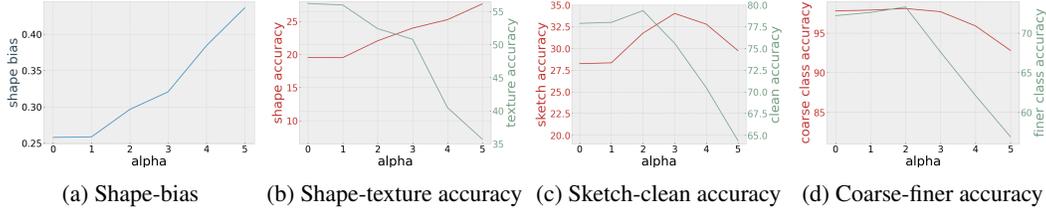


(a) Shape-bias  (b) Shape-texture accuracy  (c) Sketch-clean accuracy  (d) Coarse-finer accuracy

Figure 6: Larger $\alpha$ monotonically increases the model bias to shape features over texture features. Model performance is impacted by such a trade-off differently under different settings. In all scenarios, slightly calibrated shape bias improves model performance. The test settings represented in the red lines gain more from the increased shape bias than the settings represented in the green lines.

We train MoCo-v2 models with different $\alpha$ from 1 to 5 on the ImageNet-100 dataset. As shown in Figure 6a, we find that $\alpha$ effectively controls the trade-off on the model bias to the shape and texture features. $\alpha = 0$ is used to denote the baseline method. Specifically, a larger $\alpha$ in the loss function 1 leads to a larger penalty on the similarity between the representations of query samples and non-semantic samples, consequently a larger shape bias. We follow [14] to calculate shape bias on the stimuli images with conflicted shape and texture clues generated by style transfer. We show the corresponding accuracy on the shape and texture labels in Figure 6b.

As shown in Figure 6b, when $\alpha$ increases the texture accuracy on the stimuli dataset monotonically decreases while the shape accuracy monotonically increases. To further study how the trade-off between shape and texture bias impacts the model performance, we first compare the accuracy on the ImageNet validation dataset and ImageNet-Sketch dataset [49] when $\alpha$ varies in Figure 6c. We find that on both datasets, slightly increased shape bias over baseline ($\alpha = 2$) improves performances. Interestingly, the accuracy peak on the ImageNet-Sketch appears at $\alpha = 3$ while the peak appears at $\alpha = 2$ on the ImageNet validation dataset. In addition, for even larger $\alpha$ the accuracy on the ImageNet-Sketch dataset still outperforms the baseline while the standard accuracy gets hurt. This shows that different downstream tasks may benefit differently from differently shape-biased models.

We plot the histograms of similarities calculated by the models trained with different $\alpha$ in Appendix Figure 13. We find that large $\alpha$ makes the original pretext task challenging - the model cannot effectively pull together the representations of positive pairs. Specifically, when $\alpha$ increases from 1 to 5, the average similarity of the positive pairs decrease from $0.9267$ to $0.7541$. This demonstrates that it is hard for the model to learn representations that are completely independent of the texture features contained in the non-semantic images.

## 4.2 Rethinking the shape-texture trade-off through class-based analysis

The initial discussion on the shape-texture trade-off shows that humans rely more on the shape features while CNNs rely more on texture features and increasing shape bias can improve accuracy and robustness [14]. However, similar to [40, 23], we notice that increasing shape bias does not always improve the generalization and robustness of the models. To better understand this phenomenon, we provide two observations based on the analysis of the ImageNet-Texture dataset and our method to explain why a texture-biased model is helpful with classification on the ImageNet dataset.

First, we find that an increasing shape bias often leads to more errors among the fine-grained classes. The initial discussion of the shape-bias [14] only pays attention to the selected 16 coarse classes. We thus compare the finer and coarse class accuracy on the dog images of ImageNet dataset as in Figure 6d. For coarse class accuracy, the predictions are counted to be correct whenever the image is classified as a dog class, no matter which dog class is predicted, while for the finer class accuracy, only those predictions of target dog classes are counted to be correct. We notice that the finer class accuracy drops more significantly when shape bias increases as opposed to the coarse class accuracy. For example, when $\alpha = 3$, the finer class accuracy drops from $72.9$ to $67.6$ while the coarse class accuracy slightly decreases from $97.8$ to $97.7$. Therefore, for datasets with numerous fine-grained classes like ImageNet, a texture-biased model is more helpful for a higher accuracy, which confirms

n04589890, window screen, 0.58        n03207743, dishrag, 0.82        n02892201, plaque, 0.64

Figure 7: A ResNet-50 model trained on the ImageNet dataset achieves decent accuracy when only texture features are available on some classes. A normal image and its texture version are displayed for some of these classes. The caption indicates the class ID, name, and accuracy on texture images.

the previous conjecture [53] . Second, in Appendix Figure 9, we show a scatter plot of texture accuracy vs. standard accuracy of different model architectures and a histogram of accuracy on individual ImageNet-Texture classes. We identify several classes where using only texture features is sufficient to achieve a high classification accuracy. These classes are all missing in the previous study [14]. As shown in Figure 7 and Appendix Figure 10, texture serves as a more important clue than the shape for these classes.

## 5 Related Work

**Contrastive learning based self-supervised learning**    Recent contrastive learning based self-supervised learning methods including MoCo [18], SimCLR [7], InfoMin[48], SimSiam[9], BYOL [16], SwAV [6], Barlow Twins [55] have proven helpful in learning visual representations. These methods rely on different pretext tasks to increase the agreement among the different views of the same image. The augmentations used to generate these views are essential to the success of these contrastive learning methods [7, 6] by preventing shortcuts such as the use of simple color histogram [7]. There is an ongoing trend of developing novel augmentations [6, 48] or adaptively applying augmentations [53] and consistent improvement has been achieved with these studies. However, it is intractable to eliminate every shortcut and sometimes tricky to craft the correct positive pairs. Different from these methods, we show that augmentations that perturb the semantic features and craft negative samples can be more effective to impose additional regularization. For example, to prevent models from relying on local features, it is much easier to destroy global features and create negatives than to remove all the local features and create positives. Furthermore, by maximizing the difference between natural images and their non-semantic versions in the representation space, the models are coerced to avoid any potential shortcuts shared by them.

Methods like MoCo [18] and SimCLR [7] distinguish positive pairs from negative pairs that are picked from the rest of the dataset. However, most of the negatives prove to be unnecessary and insufficient [13, 31]. To excavate effective negative samples, these methods heavily depend on the large batch sizes [7] or memory bank [18]. Utilizing hard negative samples has long been recognized as an effective approach to boost model performance [17, 29, 54, 45]. In the contrastive learning studies, [10, 43] modify the contrastive learning loss to make it assign greater weights to the hard negative samples. [31] proposes to synthesize hard negative samples by taking linear combinations of the hardest negative samples. Our work is orthogonal to these ideas in the way that we propose to generate negative samples from given images themselves to reduce the reliance on the undesired features. In addition, two recent works [25, 32] study the application of adversarial examples as hard positive and negative samples in contrastive learning. [44] augments the images by manipulating their foregrounds and backgrounds to generate negative and positive samples. Compared with these studies, we mainly focus on the OOD evaluation of the models. In addition, our patch-based augmentation is also related to the self-Supervised learning methods that adopt the pretext task based on jigsaw [41, 39, 20], which we discuss in the Appendix B.2.

**Robustness and out-of-domain generalization of CNNs**    High test accuracy provides no guarantee that a network learns high-level semantic features instead of low-level superfluous features that exist in both training and test dataset [30]. An increasing number of studies have corroborated such concerns and found that CNNs can rely on local patches [4, 3], texture [14], high-frequency components [50] and even artificially added features [26] to achieve high test accuracy. These super-

ficial correlations become brittle under large domain shifts [22, 21]. This still remains an unsolved problem [46] and is rarely discussed in the contrastive learning setting.

Among all these undesired features, the shape-texture bias has been widely discussed in recent studies [14, 36]. Previous work has shown that CNNs trained on the ImageNet dataset are biased to texture features and such over-reliance can hurt the generalization performance of CNNs [14, 24]. Several studies have aimed at mitigating this problem [38, 35] or providing a better understanding [23, 27]. In this paper, we introduce a dataset called ImageNet-Texture, which can help future studies on these problems. Our method also effectively controls the trade-off between shape and texture bias. We provide new insights about this problem based on the analysis of our method and dataset.

## 6 Closing Remarks

**Conclusion** CNNs are prone to learn discriminative features that are vulnerable under domain shifts. In this paper, we first demonstrate the regularization power of contrastive learning to discard any undesired features by generating appropriate negative samples. We explore two approaches, the patch-based and texture-based augmentations, to craft negative samples with only local features preserved. We show that the representations learned by contrastive learning with such negative samples depend less on the local features, and consequently generalize better under OOD settings. We hope this paper can encourage people to rethink the role that negative samples play in contrastive learning, which hopefully leads to more efficient methods to generate negative samples.

**Limitations** The problem of dependence on superficial features exists in various domains beyond vision, such as language [37, 28]. Therefore, it is intriguing to consider generalizing such an idea to other modalities. In addition, as the mechanism to ensure that contrastive learning models trained on large datasets to discard the bias of the datasets is yet to be invented, severe social issues in fairness or privacy may be raised as a result [5, 33]. In this paper, we show how to calibrate the bias towards texture features using proposed negative samples. It is also worth considering whether contrastive learning can be used to address any of these negative effects triggered by bias in datasets.

## References

[1] M. Ashikhmin. Synthesizing natural textures. In *Proceedings of the 2001 symposium on Interactive 3D graphics*, pages 217–226, 2001.

[2] L. Beyer, O. J. Hénaff, A. Kolesnikov, X. Zhai, and A. van den Oord. Are we done with imagenet?, 2020.

[3] W. Brendel and M. Bethge. Approximating CNNs with bag-of-local-features models works surprisingly well on imagenet. In *International Conference on Learning Representations*, 2019.

[4] T. B. Brown, D. Mané, A. Roy, M. Abadi, and J. Gilmer. Adversarial patch. *arXiv preprint arXiv:1712.09665*, 2017.

[5] J. Buolamwini and T. Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR, 2018.

[6] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*, 2020.

[7] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.

[8] X. Chen, H. Fan, R. Girshick, and K. He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.

[9] X. Chen and K. He. Exploring simple siamese representation learning. *arXiv preprint arXiv:2011.10566*, 2020.

[10] C.-Y. Chuang, J. Robinson, Y.-C. Lin, A. Torralba, and S. Jegelka. Debiased contrastive learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 8765–8775. Curran Associates, Inc., 2020.

[11] J. S. De Bonet. Multiresolution sampling procedure for analysis and synthesis of texture images. In *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*, pages 361–368, 1997.

[12] A. A. Efros and T. K. Leung. Texture synthesis by non-parametric sampling. In *Proceedings of the seventh IEEE international conference on computer vision*, volume 2, pages 1033–1038. IEEE, 1999.

[13] J. Frankle, D. J. Schwab, A. S. Morcos, et al. Are all negatives created equal in contrastive instance discrimination? *arXiv preprint arXiv:2010.06682*, 2020.

[14] R. Geirhos, P. Rubisch, C. Michaelis, M. Bethge, F. A. Wichmann, and W. Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019.

[15] R. Geirhos, C. R. M. Temme, J. Rauber, H. H. Schütt, M. Bethge, and F. A. Wichmann. Generalisation in humans and deep neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 7549–7561, 2018.

[16] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, B. Piot, k. kavukcuoglu, R. Munos, and M. Valko. Bootstrap your own latent - a new approach to self-supervised learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 21271–21284. Curran Associates, Inc., 2020.

[17] B. Harwood, V. Kumar B G, G. Carneiro, I. Reid, and T. Drummond. Smart mining for deep metric learning. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

[18] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.

[19] D. J. Heeger and J. R. Bergen. Pyramid-based texture analysis/synthesis. In *Proceedings of the 22nd annual conference on Computer graphics and interactive techniques*, pages 229–238, 1995.

[20] O. Henaff. Data-efficient image recognition with contrastive predictive coding. In *International Conference on Machine Learning*, pages 4182–4192. PMLR, 2020.

[21] D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. Zhu, S. Parajuli, M. Guo, D. Song, J. Steinhardt, and J. Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. *arXiv preprint arXiv:2006.16241*, 2020.

[22] D. Hendrycks and T. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.

[23] K. Hermann, T. Chen, and S. Kornblith. The origins and prevalence of texture bias in convolutional neural networks. *Advances in Neural Information Processing Systems*, 33, 2020.

[24] K. L. Hermann and A. K. Lampinen. What shapes feature representations? exploring datasets, architectures, and training, 2020.

[25] C.-H. Ho and N. Nvasconcelos. Contrastive learning with adversarial examples. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 17081–17093. Curran Associates, Inc., 2020.

[26] H. Huang, X. Ma, S. M. Erfani, J. Bailey, and Y. Wang. Unlearnable examples: Making personal data unexploitable. In *International Conference on Learning Representations*, 2021.

[27] M. A. Islam, M. Kowal, P. Esser, S. Jia, B. Ommer, K. G. Derpanis, and N. Bruce. Shape or texture: Understanding discriminative features in {cnn}s. In *International Conference on Learning Representations*, 2021.

[28] R. Jia and P. Liang. Adversarial examples for evaluating reading comprehension systems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2021–2031, 2017.

[29] S. Jin, A. RoyChowdhury, H. Jiang, A. Singh, A. Prasad, D. Chakraborty, and E. Learned-Miller. Unsupervised hard example mining from videos for improved object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 307–324, 2018.

[30] J. Jo and Y. Bengio. Measuring the tendency of cnns to learn surface statistical regularities. *arXiv preprint arXiv:1711.11561*, 2017.

[31] Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, and D. Larlus. Hard negative mixing for contrastive learning. In *Neural Information Processing Systems (NeurIPS)*, 2020.

[32] M. Kim, J. Tack, and S. J. Hwang. Adversarial self-supervised contrastive learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 2983–2994. Curran Associates, Inc., 2020.

[33] S. Kiritchenko and S. M. Mohammad. Examining gender and race bias in two hundred sentiment analysis systems. *arXiv preprint arXiv:1805.04508*, 2018.

[34] C.-L. Li, K. Sohn, J. Yoon, and T. Pfister. Cutpaste: Self-supervised learning for anomaly detection and localization. *arXiv preprint arXiv:2104.04015*, 2021.

[35] Y. Li, Q. Yu, M. Tan, J. Mei, P. Tang, W. Shen, A. Yuille, and cihang xie. Shape-texture debiased neural network training. In *International Conference on Learning Representations*, 2021.

[36] G. Malhotra and J. Bowers. The contrasting roles of shape in human vision and convolutional neural networks. In *CogSci*, 2019.

[37] T. McCoy, E. Pavlick, and T. Linzen. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3428–3448, 2019.

[38] S. Mishra, A. Shah, A. Bansal, J. Choi, A. Shrivastava, A. Sharma, and D. Jacobs. Learning visual representations for transfer learning by suppressing texture. *arXiv preprint arXiv:2011.01901*, 2020.

[39] I. Misra and L. v. d. Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6707–6717, 2020.

[40] C. K. Mummadi, R. Subramaniam, R. Hutmacher, J. Vitay, V. Fischer, and J. H. Metzen. Does enhanced shape bias improve neural network robustness to common corruptions? In *ICLR*, 2021.

[41] M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European conference on computer vision*, pages 69–84. Springer, 2016.

[42] J. Portilla and E. P. Simoncelli. A parametric texture model based on joint statistics of complex wavelet coefficients. *International journal of computer vision*, 40(1):49–70, 2000.

[43] J. Robinson, C.-Y. Chuang, S. Sra, and S. Jegelka. Contrastive learning with hard negative samples. *International Conference on Learning Representations*, 2021.

[44] C. K. Ryali, D. J. Schwab, and A. S. Morcos. Leveraging background augmentations to encourage semantic focus in self-supervised contrastive learning. *arXiv preprint arXiv:2103.12719*, 2021.

[45] E. Simo-Serra, E. Trulls, L. Ferraz, I. Kokkinos, P. Fua, and F. Moreno-Noguer. Discriminative learning of deep convolutional feature point descriptors. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 118–126, 2015.

[46] R. Taori, A. Dave, V. Shankar, N. Carlini, B. Recht, and L. Schmidt. Measuring robustness to natural distribution shifts in image classification. *Advances in Neural Information Processing Systems*, 33, 2020.

[47] Y. Tian, D. Krishnan, and P. Isola. Contrastive multiview coding. In *ECCV*, 2020.

[48] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola. What makes for good views for contrastive learning. *arXiv preprint arXiv:2005.10243*, 2020.

[49] H. Wang, S. Ge, Z. Lipton, and E. P. Xing. Learning robust global representations by penalizing local predictive power. In *Advances in Neural Information Processing Systems*, pages 10506–10518, 2019.

[50] H. Wang, X. Wu, Z. Huang, and E. P. Xing. High-frequency component helps explain the generalization of convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8684–8694, 2020.

[51] L.-Y. Wei and M. Levoy. Fast texture synthesis using tree-structured vector quantization. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 479–488, 2000.

[52] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018.

[53] T. Xiao, X. Wang, A. A. Efros, and T. Darrell. What should not be contrastive in contrastive learning. In *International Conference on Learning Representations*, 2020.

[54] H. Xuan, A. Stylianou, X. Liu, and R. Pless. Hard negative examples are hard, but useful. In *European Conference on Computer Vision*, pages 126–142. Springer, 2020.

[55] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: Self-supervised learning via redundancy reduction. *arXiv preprint arXiv:2103.03230*, 2021.

# A ImageNet-Texture

See Figures 7 and 8 for examples of the ImageNet-Texture dataset and their counterparts in the original ImageNet dataset. The analysis of models pretrained on ImageNet and evaluated on ImageNet-Texture is presented in Section A.2.
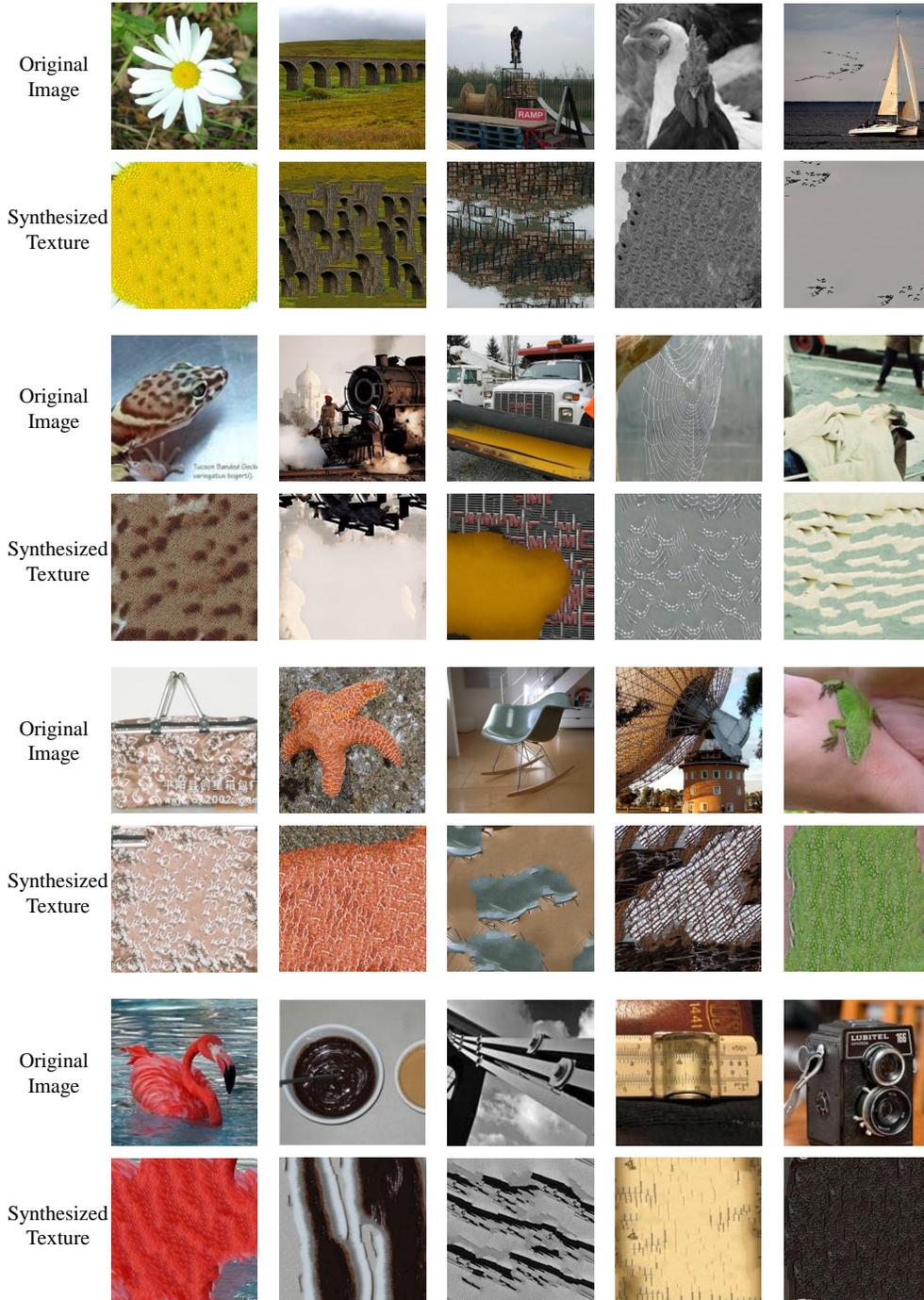
## A.1 Examples



Figure 7: Randomly picked images from our ImageNet-Texture dataset and their corresponding original images from the ImageNet dataset.
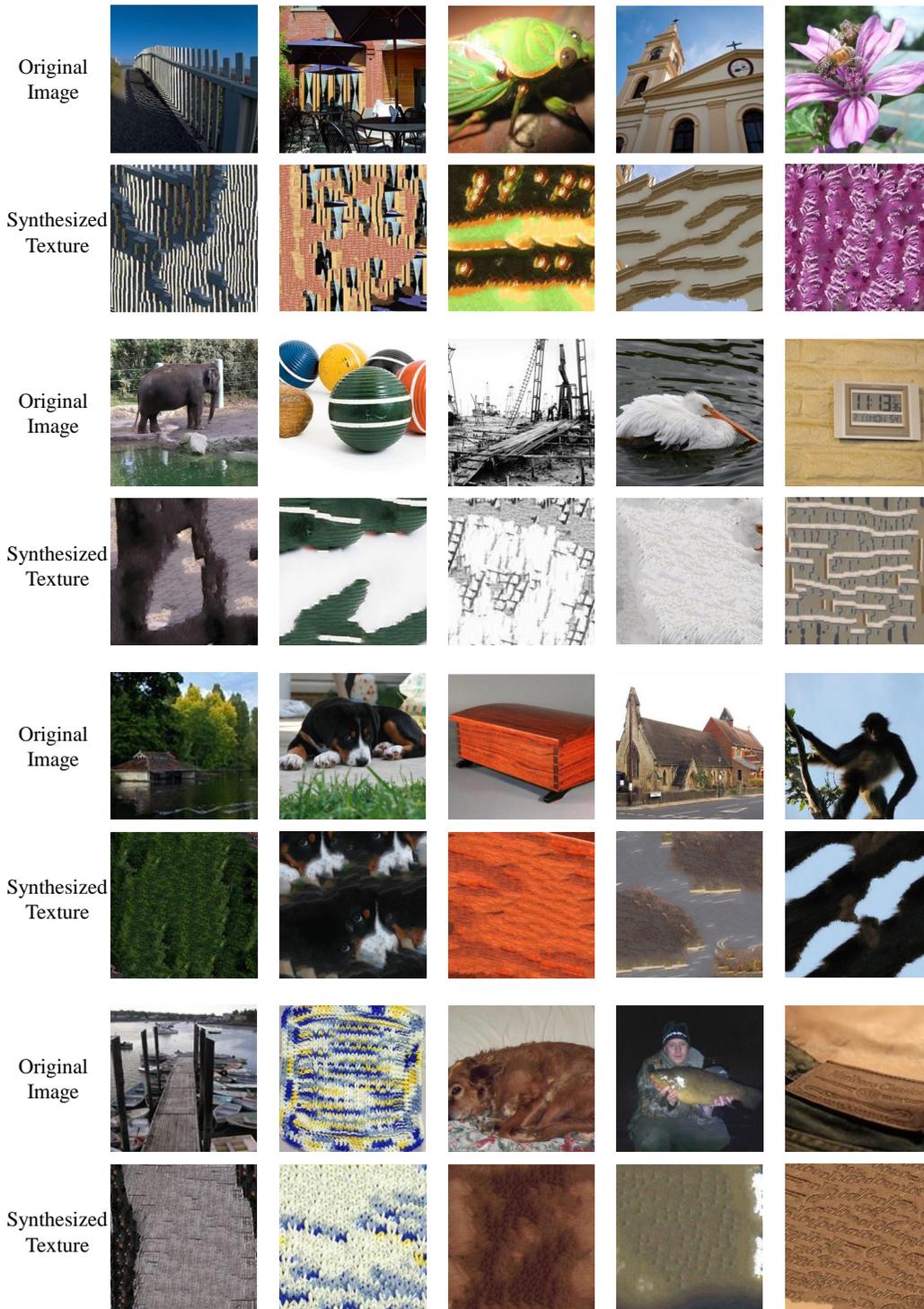
Figure 8: More randomly picked images from our ImageNet-Texture dataset and their corresponding original images from the ImageNet dataset.

## A.2 Analysis

The proposed ImageNet-Texture dataset contains images with preserved texture information and diminished shape information, which can be used to better understand how models behave when only texture features are available. We evaluate the pretrained models provided in Pytorch model zoo [2] on the ImageNet-Texture versus standard ImageNet and report the accuracy in Figure 9a. As shown in the figure, VGG networks generally have a higher accuracy on the ImageNet-Texture dataset. For the rest of the model architectures, we see a positive correlation between the standard accuracy and texture accuracy.

We also plot the histogram of the classes with different accuracy based on a pretrained ResNet-50 model in Figure 9b, which achieves a relatively high accuracy, $9.3\%$, among all the models. The class-based accuracy denotes the accuracy on the images of the same certain class. We find that on the 343 out of 100 classes, the model can only achieve an accuracy smaller than $2.5\%$, which shows that texture feature only is not sufficient to classify these classes. But we do notice a long tail of the distribution. Specifically, 15 classes have an accuracy larger than $50\%$. Sample images of these classes are demonstrated in Figure 10. Notably, texture plays an important role in distinguishing these classes. Shape is often less well-defined in these classes, for example in window screen and rapeseed.

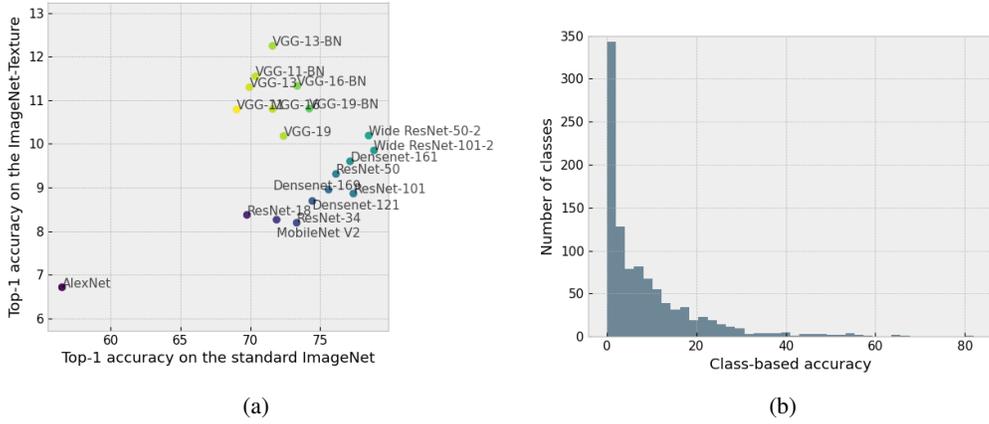

(a)                                       (b)

Figure 9: (a) Accuracy on the ImageNet-Texture versus standard ImageNet with different model architectures. (b) Histogram of accuracy on the individual ImageNet-Texture classes with ResNet-50.

## B  Discussion on the contrastive learning with non-semantic negatives

### B.1  Comparison of two ways to apply $\alpha$ in NCE loss

Since the non-semantic negative samples play different roles from the original negative samples, we introduce an additional parameter to control the penalty w.r.t. non-semantic negatives. There are two straightforward alternatives to implement $\alpha$, which we call $\mathcal{L}_{in}$ and $\mathcal{L}_{out}$:

$$\mathcal{L}_{in} = -\sum_{i \in I} \log \frac{\exp\left(z_i \cdot z_p / \tau\right)}{\exp\left(z_i \cdot z_p / \tau\right) + \exp\left(\alpha z_i \cdot z_{ns} / \tau\right) + \sum_{n \in \mathcal{N}} \exp\left(z_i \cdot z_n / \tau\right)}$$

$$\mathcal{L}_{out} = -\sum_{i \in I} \log \frac{\exp\left(z_i \cdot z_p / \tau\right)}{\exp\left(z_i \cdot z_p / \tau\right) + \alpha \exp\left(z_i \cdot z_{ns} / \tau\right) + \sum_{n \in \mathcal{N}} \exp\left(z_i \cdot z_n / \tau\right)}$$

We compare the relative importance of different pairs play in the gradient w.r.t. $z_i$:

$$\frac{\partial \mathcal{L}_{in}}{\partial z_i} = \frac{z_p / \tau \exp\left(z_i \cdot z_p / \tau\right) + \alpha z_{ns} / \tau \exp\left(\alpha z_i \cdot z_{ns} / \tau\right) + \sum_{n \in \mathcal{N}} z_n / \tau \exp\left(z_i \cdot z_n / \tau\right)}{\exp\left(z_i \cdot z_p / \tau\right) + \exp\left(\alpha z_i \cdot z_{ns} / \tau\right) + \sum_{n \in \mathcal{N}} \exp\left(z_i \cdot z_n / \tau\right)} - z_p / \tau$$

---

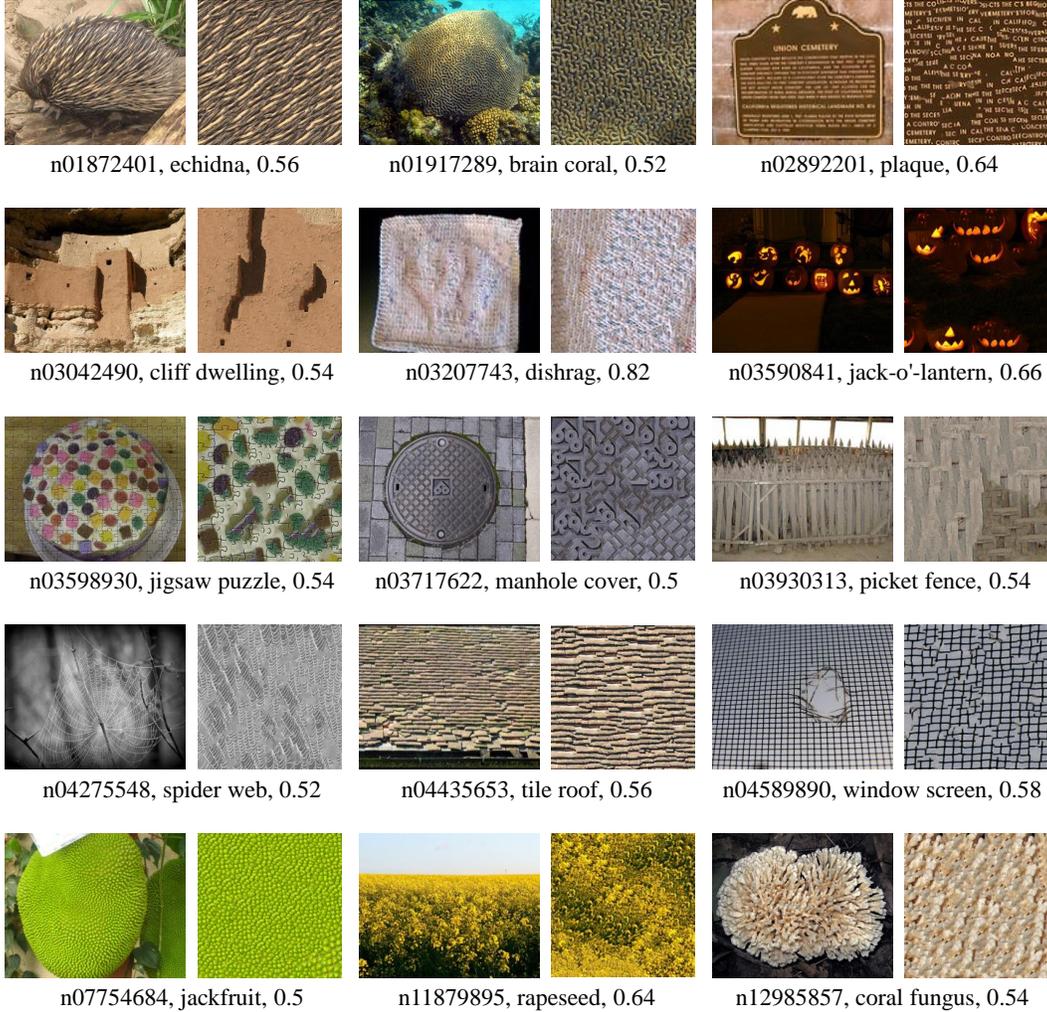[2] https://pytorch.org/vision/stable/models.html

15

Figure 10: On some classes, a ResNet-50 model trained on the standard ImageNet dataset can achieve $> 50\%$ accuracy when only texture features are available. For each class, one sample image and its texture version are displayed. The caption of each image pair indicates the ImageNet class ID, class name and ResNet-50 accuracy on the texture images.

$$\frac{\partial \mathcal{L}_{out}}{\partial z_i} = \frac{z_p/\tau \exp\left(z_i \cdot z_p/\tau\right) + \alpha z_{ns}/\tau \exp\left(z_i \cdot z_{ns}/\tau\right) + \sum_{n \in \mathcal{N}} z_n/\tau \exp\left(z_i \cdot z_n/\tau\right)}{\exp\left(z_i \cdot z_p/\tau\right) + \alpha \exp\left(z_i \cdot z_{ns}/\tau\right) + \sum_{n \in \mathcal{N}} \exp\left(z_i \cdot z_n/\tau\right)} - z_p/\tau$$

Since the denominator normalizes the 3 kinds of pairs equally, we only pay attention to the numerator. The difference between $\mathcal{L}_{out}$ and $\mathcal{L}_{in}$ is that $\mathcal{L}_{out}$ has $\alpha$ inside the exponential. Because of the exponential tail, it applies a exponentially larger weight to the negatives that are harder. Since non-semantic negatives are often harder as shown in Figure 2, $\mathcal{L}_{in}$ regularizes the non-semantic negatives more effectively than $\mathcal{L}_{out}$. Note that the implementation of $\mathcal{L}_{out}$ is equivalent to using $\alpha$ non-semantic negatives samples for each input image. In addition, the number of standard negative $|\mathcal{N}|$ in the loss also affects the relative importance of non-semantic. The smaller $|\mathcal{N}|$ is, the larger effect is played by the non-semantic negative.

## B.2 Comparison of patch-based negatives and jigsaw-based pretext tasks

Our patch-based augmentation is also closely related to some of the self-supervised learning methods which solve jigsaw as the pretext task. Specifically, [41] proposes to learn meaningful representations through predicting the order of shuffled 3x3 patches of a given image. [39] further extends this idea

to contrastive learning and learn representations that are invariant under such pre-text transformation. [20] exploit the connection between local patches to learn the representations by fusing it with contrastive predictive coding. The common idea behind these methods is to learn the representations based on the correct configuration of the patches in the natural images. However, our patch-based augmentation treats the entire image with unsorted patches as a wrong configuration of the a natural image and use it as the negative samples in the contrastive learning.

## C  Experimental details and additional results

### C.1  Implementation details

For ImageNet-1K dataset, we follow the official hyperparameters to MoCo-v2 [8]. For ImageNet-100 and STL-10 datasets, we follow [31] which uses a different memory bank size of 16384 and batch size 128 during the pretraining. For linear evaluation, we adopt a learning rate 10.0, batch size 128, and a learning schedule that decreases the learning rate by 0.1 at 30, 40, and 50 epochs. For BYOL pretraining, we follow [16] to utilize the LARS optimizer with cosine learning rate schedule with a global weight decay parameter of $10^{-6}$ and momentum of 0.9. We use a batch size of 1024 and an initial learning rate 4.8. For both MoCo and BYOL, the non-semantic negatives are input to the momentum branch. All of our models are trained on 4 GTX 1080 Ti gpus. Training single model takes around 1.5 and 12 days on the ImageNet-100 and ImageNet-1K datasets respectively.

### C.2  Results on ImageNet-C with different levels of severity

We show the accuracy of different MoCo-v2 models on the ImageNet-C-100 dataset with different corruption severity levels in Table 5. Specifically, as we have seen in the main body of the paper, a larger $\alpha$ often favors larger domain shift. This is further demonstrated by the fact that the model with $\alpha = 3$ outperforms the model with $\alpha = 2$ on the highest corruption level 5.

| Severity level | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| MoCo-v2 - $k = 16384$ | $65.58_{\pm 0.32}$ | $53.41_{\pm 0.27}$ | $42.31_{\pm 0.31}$ | $31.87_{\pm 0.31}$ | $22.22_{\pm 0.20}$ |
| + Texture-base NS | $66.00_{\pm 0.47}$ | $54.11_{\pm 0.43}$ | $42.96_{\pm 0.41}$ | $32.23_{\pm 0.26}$ | $22.60_{\pm 0.23}$ |
| + Patch-base NS - $\alpha = 2$ | $\mathbf{67.84}_{\pm 0.21}$ | $\mathbf{55.99}_{\pm 0.20}$ | $\mathbf{44.65}_{\pm 0.40}$ | $\mathbf{33.74}_{\pm 0.51}$ | $23.41_{\pm 0.50}$ |
| + Patch-base NS - $\alpha = 3$ | $65.82_{\pm 0.04}$ | $54.95_{\pm 0.13}$ | $44.10_{\pm 0.17}$ | $33.67_{\pm 0.25}$ | $\mathbf{23.71}_{\pm 0.28}$ |
| MoCo-v2 - $k = 8192$ | $65.44_{\pm 0.52}$ | $53.52_{\pm 0.49}$ | $42.71_{\pm 0.42}$ | $32.09_{\pm 0.33}$ | $22.36_{\pm 0.19}$ |
| + Patch-base NS - $\alpha = 2$ | $\mathbf{68.01}_{\pm 0.01}$ | $\mathbf{56.26}_{\pm 0.12}$ | $\mathbf{45.04}_{\pm 0.28}$ | $\mathbf{34.16}_{\pm 0.34}$ | $\mathbf{23.92}_{\pm 0.27}$ |

Table 5: Top-1 accuracy on ImageNet-C-100 dataset with different levels of severity.

### C.3  Additional ablations on the patch-based augmentations

#### C.3.1  Patch sizes

We show more results on the patch-based augmentations with different patch sizes by reporting the test accuracy on the ImageNet-100 validation set and corresponding ImageNet-Sketch dataset in Figure 11. Specifically, we sample patch sizes from uniform distributions $d \sim \mathcal{U}(x, y)$ of different interval $[x, y]$. We find that for the performance on the standard validation set, both the large or smaller patch sizes cause a less desired accuracy. But for the ImageNet-Sketch dataset, the larger patch sizes generally provide larger performance improvement. Our conjecture on this observation is that using larger patch sizes prevents the model to learn some of the local features that are shared between training and validation set while absent from the sketch dataset.

#### C.3.2  Data augmentations on individual patches

We test whether common data augmentation methods can be used to augment individual patches so as to further improve the performance. We report the accuracy on the ImageNet-100 validation set with the model trained with the patch-based negative samples in Table 6. The patches are potentially augmented with horizontal flip or vertical flip with 50% probability, or rotation in $0°, 90°, 180°$ or $270°$ with 25% probability respectively. We find that for a larger patch size, i.e. $d \sim \mathcal{U}(16, 72)$, such

|     | 8 | 16 | 32 |
| --- | --- | --- | --- |
| 56 | 78.68 | 78.94 | 78.85 |
| 72 | 78.92 | 79.35 | 79.34 |
| 108 | 78.82 | 78.98 | 78.9 |

(a) Standard accuracy

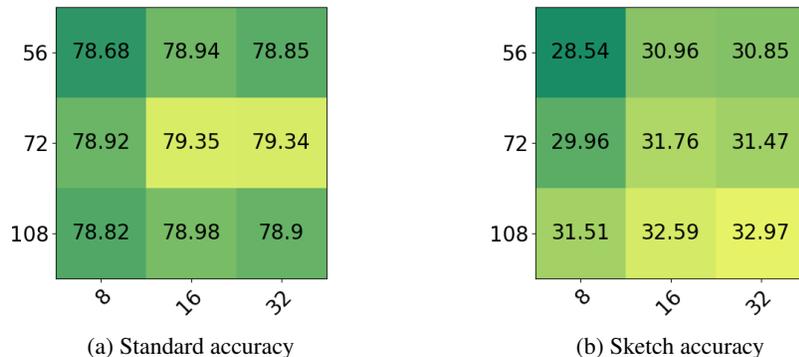|     | 8 | 16 | 32 |
| --- | --- | --- | --- |
| 56 | 28.54 | 30.96 | 30.85 |
| 72 | 29.96 | 31.76 | 31.47 |
| 108 | 31.51 | 32.59 | 32.97 |

(b) Sketch accuracy

Figure 11: Test accuracy of patch-based non-semantic augmentations with patch sizes sampled from different uniform distribution. The x-axis indicates the lower boundary of the sampling interval and the y-axis indicates the higher boundary.

augmentations always downplay the accuracy. But for a smaller patch size, such augmentations, especially with horizontal flip only, can improve the performance. For the larger patch size, we conjecture that it is important to ensure that the patches to be identical to the certain part of the input image. Therefore the non-semantic features are best preserved. However, for a smaller patch size, such augmentations have a less impact on the captured non-semantic information.

| Patch size | No aug. | + Horizontal flip | + Vertical flip | + Rotation |
| --- | --- | --- | --- | --- |
| $8 - 28$ | 78.58 | 78.72 | 78.66 | 78.64 |
| $16 - 72$ | 79.35 | 79.06 | 78.92 | 78.40 |

Table 6: Accuracy on ImageNet-100 validation set with different configurations of augmentations on the individual patches. Augmentations are cumulative across the columns (e.g. the "+ Vertical flip" model used horizontal and vertical flip).

## C.4 BYOL with different kinds of negative samples

The proper way to add negative samples to BYOL is still an open research problem. In this section, we provide additional experiments to demystify the role of non-semantic negative samples played in BYOL. Specifically, we train BYOL with a regular negative that is randomly picked from the same batch with the query sample excluded, or a regular negative as well as a patch-based non-semantic negative. We report the accuracy on the ImageNet-100 variants in the Table 7. We find that introducing a regular negative sample without further modification dramatically undermines the performance, which demonstrates that the improved performance is due to the negative mining strategy instead of adding an arbitrary negative sample.

|     | ImageNet | ImageNet-C | ImageNet-S | Stylized-ImageNet | ImageNet-R |
| --- | --- | --- | --- | --- | --- |
| BYOL | 78.76 | 44.43 | 35.84 | 15.01 | 39.53 |
| + Patch-based | 78.81 | 44.60 | 36.76 | 15.52 | 41.16 |
| + Reg. | 57.22 | 24.19 | 19.11 | 6.74 | 20.47 |
| + Patch-based and Reg. | 57.24 | 24.37 | 20.11 | 6.60 | 20.95 |

Table 7: Top-1 accuracy of BYOL trained with different negatives on the ImageNet-100 datasets.

## C.5 Hardness of texture-based negative samples

We show the distribution of similarities between the representations of input image versus different paired samples in Figure 12. The similarities are calculated by the models trained without (blue) and with (red) texture-based negative samples. We use $\alpha = 2$ when texture-based negative samples are

18

adopted. Comparing with Figure 3, the regularization of texture-based negative samples is generally weaker than the patch-based negative samples. For example, the average similarity w.r.t patch-based negative samples decrease from $0.4040$ to $0.3677$ when using texture-based negative samples for training, which is $0.3252$ when using patch-based negative samples.



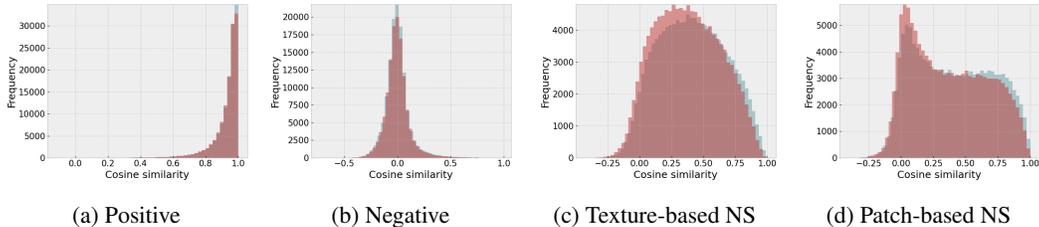| (a) Positive | (b) Negative | (c) Texture-based NS | (d) Patch-based NS |

Figure 12: The histogram of cosine similarity between the representations of different kinds of pairs using models trained without (blue) and with (red) texture-based negative samples.

## C.6 Hardness of patch-based negative samples using models trained with different $\alpha$

We show the distribution of similarities calculated by the models trained with patch-based negative samples and different $\alpha$ in Figure 13. When $\alpha$ increases from 0 to 5, we see a larger penalty on the similarities to both patch-based and texture-based negative samples. When $\alpha = 5$, the distributions of two negative samples are close to a normal distribution like the distribution of the standard negative samples. Specifically, the average similarity w.r.t. the patch-based negative samples further decreases to $0.4040$ to $0.1184$. In terms of the original pretext task, we find that the distribution of both standard negative and positive samples become wider and less concentrated to $0$ and $1$ respectively. This shows that a larger $\alpha$ makes the optimization of the original objective more difficult.

## C.7 More results on ImageNet-1K

| | ImageNet | ImageNet-C | ImageNet-S | Stylized | ImageNet-R |
|---|---|---|---|---|---|
| MoCo-v2 [8] | 67.60 | 87.7 | 17.47 | 5.55 | 27.81 |
| + MoCHi (128,1024,512) [31] | 66.62 | 90.0 | 16.00 | 5.32 | 25.29 |
| + MoCHi (512,1024,512) [31] | 67.56 | 88.7 | 16.32 | 5.94 | 25.71 |
| + Patch-base NS - $\alpha = 2, k = 65536$ | **67.92** | 87.6 | 18.58 | 6.34 | 28.95 |
| + Patch-base NS - $\alpha = 3, k = 65536$ | 60.80 | 92.1 | 18.79 | **6.80** | 28.11 |
| + Patch-base NS - $\alpha = 2, k = 32768$ | 67.83 | **87.2** | 18.70 | 6.06 | 28.50 |
| + Patch-base NS - $\alpha = 2, k = 16384$ | 67.34 | 87.5 | **19.49** | 6.49 | **29.45** |

Table 8: Top-1 accuracy on the ImageNet-1K dataset and its sketch, stylized, rendition variants, and mCE on the ImageNet-C dataset.

We show more results on the ImageNet-1K dataset using MoCo-v2 models trained with patch-based negative samples in Table 8. We also report the other released MoChi model [3], and the different hyperparameters are indicated following the model name in the table. For our method, we focus on the models when a larger penalty on the similarity w.r.t. the patch-based negative samples is added. Specifically, we report the results with a larger $\alpha = 3$ and smaller memory bank sizes $k = 32768$ or $k = 16384$. As shown in the table, increasing the penalty stably improves the generalization under OOD settings. For example, when $\alpha = 2$ and $k = 16384$, the accuracy on ImageNet-S and ImageNet-R increase from $17.47$ to $19.49$ and $27.81$ to $29.45$ respectively.

---

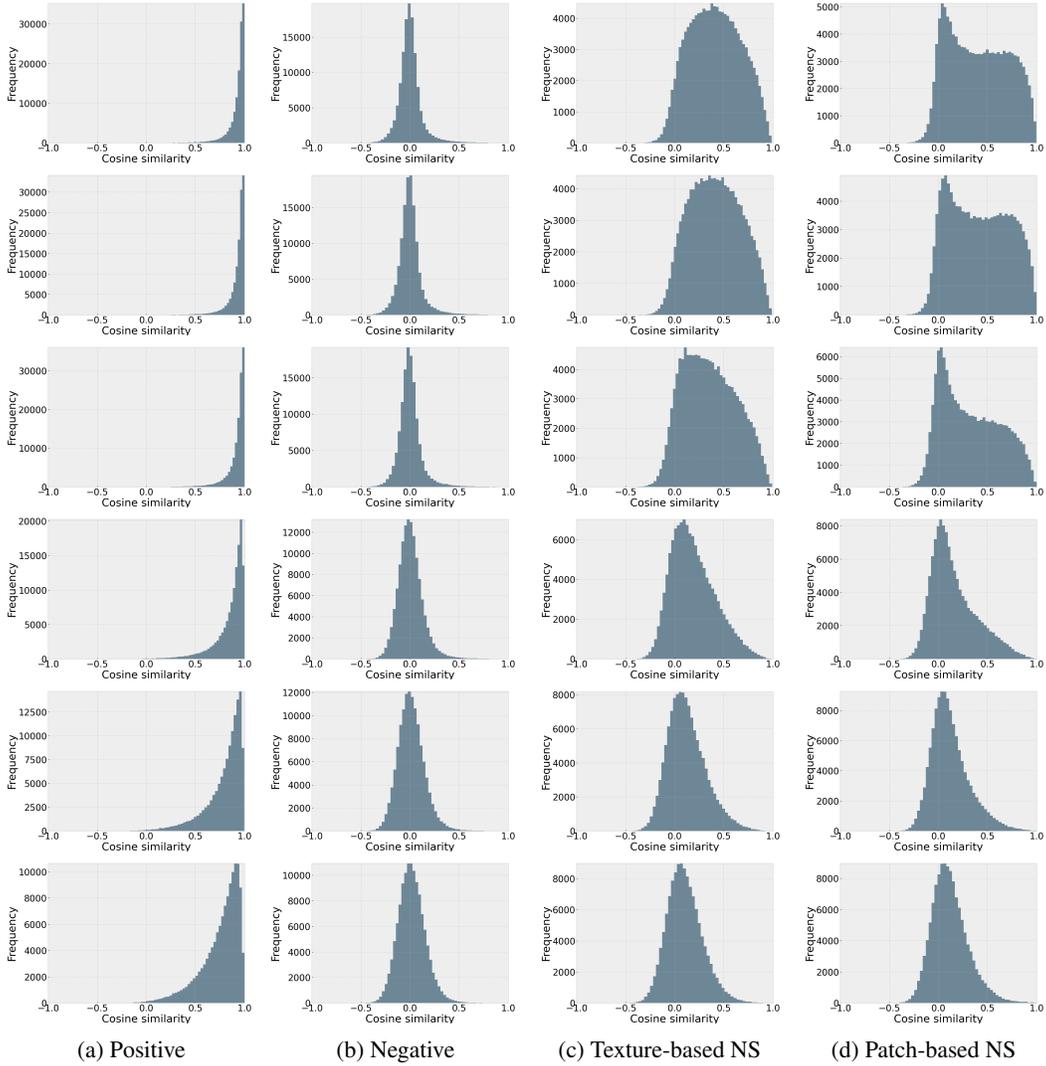[3] https://europe.naverlabs.com/research/computer-vision/mochi/

Figure 13: The histogram of cosine similarity between the representations of different kinds of pairs using models trained with patch-based negative samples and different $\alpha$ values from $0$ (top) to $5$ (bottom).