



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

When Does Contrastive Visual Representation Learning Work?

Citation for published version:

Cole, E, Yang, X, Wilber, K, Aodha, OM & Belongie, S 2022, When Does Contrastive Visual Representation Learning Work? in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Conference on Computer Vision and Pattern Recognition (CVPR), Institute of Electrical and Electronics Engineers (IEEE), pp. 14735-14744, IEEE/CVF Conference on Computer Vision and Pattern Recognition 2022, New Orleans, Louisiana, United States, 19/06/22. <https://doi.org/10.1109/CVPR52688.2022.01434>

Digital Object Identifier (DOI):

[10.1109/CVPR52688.2022.01434](https://doi.org/10.1109/CVPR52688.2022.01434)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



When Does Contrastive Visual Representation Learning Work?

Elijah Cole¹ Xuan Yang² Kimberly Wilber² Oisín Mac Aodha^{3,4} Serge Belongie⁵
¹Caltech ²Google ³University of Edinburgh ⁴Alan Turing Institute ⁵University of Copenhagen

Abstract

Recent self-supervised representation learning techniques have largely closed the gap between supervised and unsupervised learning on ImageNet classification. While the particulars of pretraining on ImageNet are now relatively well understood, the field still lacks widely accepted best practices for replicating this success on other datasets. As a first step in this direction, we study contrastive self-supervised learning on four diverse large-scale datasets. By looking through the lenses of data quantity, data domain, data quality, and task granularity, we provide new insights into the necessary conditions for successful self-supervised learning. Our key findings include observations such as: (i) the benefit of additional pretraining data beyond 500k images is modest, (ii) adding pretraining images from another domain does not lead to more general representations, (iii) corrupted pretraining images have a disparate impact on supervised and self-supervised pretraining, and (iv) contrastive learning lags far behind supervised learning on fine-grained visual classification tasks.

1. Introduction

Self-supervised learning (SSL) techniques can now produce visual representations which are competitive with representations generated by fully supervised networks for many downstream tasks [20]. This is an important milestone for computer vision, as removing the need for large amounts of labels at training time has the potential to scale up our ability to address challenges in domains where supervision is currently too difficult or costly to obtain. However, with some limited exceptions, the vast majority of current state-of-the-art approaches are developed and evaluated on standard datasets like ImageNet [43]. As a result, we do not have a good understanding of how well these methods work when they are applied to other datasets.

Under what conditions do self-supervised contrastive representation learning methods produce “good” visual representations? This is an important question for computer vision researchers because it adds to our understanding of SSL and highlights opportunities for new methods. This is

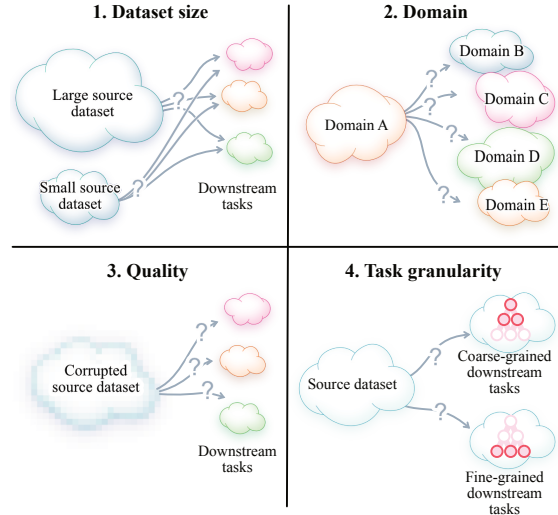


Figure 1. **What conditions are necessary for successful self-supervised pretraining on domains beyond ImageNet?** We investigate the impact of self-supervised and supervised training dataset size, the downstream domain, image quality, and the granularity of downstream classification tasks.

also an important question for domain experts with limited resources who might be interested in applying SSL to real-world problems. With these objectives in mind, we attempt to answer the following questions:

(i) What is the impact of data quantity? How many unlabeled images do we need for pretraining, and when is it worthwhile to get more? How much labeled data do we need for linear classifier training or end-to-end fine-tuning on a downstream task? In which regimes do self-supervised features rival those learned from full supervision?

(ii) What is the impact of the pretraining domain? How well do self-supervised representations trained on one domain transfer to another? Can we learn more general representations by combining datasets? Do different pretraining datasets lead to complementary representations?

(iii) What is the impact of data quality? How robust are self-supervised methods to training time image corruption such as reduced resolution, compression artifacts, or noise? Does pretraining on corrupted images lead to poor downstream performance on uncorrupted images?

(iv) What is the impact of task granularity? Does SSL

result in features that are only effective for “easy” classification tasks, or are they also useful for more challenging, “fine-grained” visual concepts?

We address the above questions through extensive quantitative evaluation across four diverse large-scale visual datasets (see Figure 1). We make several interesting observations and recommendations including:

- For an ImageNet-scale dataset, decreasing the amount of unlabeled training data by half (from 1M to 500k images) only degrades downstream classification performance by 1-2% (Figure 2). In many contexts this trade-off is reasonable, allowing for faster and cheaper pretraining. This also indicates that current self-supervised methods coupled with standard architectures may be unable to take advantage of very large pretraining sets.
- Self-supervised representations that are learned from images from the same domain as the test domain are much more effective than those learned from different domains (Table 1). Self-supervised training on our current datasets may not be sufficient to learn representations that readily generalize to many contexts.
- Neither (i) combining datasets before pretraining (Table 2) nor (ii) combining self-supervised features learned from different datasets (Table 3) leads to significant performance improvements. More work may be required before self-supervised techniques can learn highly generalizable representations from large and diverse datasets.
- Pretraining on corrupted images affects supervised and self-supervised learning very differently (Figure 4). For instance, self-supervised representations are surprisingly sensitive to image resolution.
- Current self-supervised methods learn representations that can easily disambiguate coarse-grained visual concepts like those in ImageNet. However, as the granularity of the concepts becomes finer, self-supervised performance lags further behind supervised baselines (Figure 5). The contrastive loss may lead to coarse-grained features which are insufficient for fine-grained tasks.

2. Related Work

SSL for visual representations. Early self-supervised representation learning methods typically centered around solving hand-designed “pretext tasks” like patch location prediction [18], rotation prediction [22], inpainting [40], cross-channel reconstruction [64], sorting sequences of video frames [35], solving jigsaw puzzles [38], or colorization [63]. However, more recent work has explored *contrastive learning-based* approaches where the pretext task is to distinguish matching and non-matching pairs of augmented input images [30, 39, 51]. The prototypical example is SimCLR [10, 11], which is trained to identify the matching image using a cross-entropy loss. Other variations on the contrastive SSL framework include using a

momentum encoder to provide large numbers of negative pairs (MoCo) [13, 27], adaptively scaling the margin in MoCo (EqCo) [67], and contrasting clustering assignments instead of augmented pairs (SwAV) [8]. Moving beyond the contrastive loss entirely, some papers recast the problem in a “learning-to-rank” framework (S2R2) [56], use simple feature prediction (SimSiam) [14], or predict the output of an exponential moving average network (BYOL) [26]. [6] investigates the role of negatives in contrastive learning, though we note that BYOL and SimSiam avoid using negatives explicitly. In this work, our focus is on self-supervised visual classification. We do not explore alternative settings such as supervised contrastive learning [33], contrastive learning in non-vision areas like language [42] or audio [44], or other methods that aim to reduce the annotation burden for representation learning such as large-scale weak supervision [37].

SSL beyond ImageNet. ImageNet classification has long been viewed as the gold standard benchmark task for SSL, and the gap between supervised and self-supervised performance on ImageNet has steadily closed over the last few years [8, 10, 26, 27]. There is now a growing expectation that SSL should reduce our dependence on manual supervision in challenging and diverse domains which may *not* resemble the traditional object classification setting represented by ImageNet. A number of papers have studied how well self-supervised representations pretrained on ImageNet perform on downstream tasks like fine-grained species classification [60], semantic segmentation [7], scene understanding [26], and instance segmentation [27].

More recently, researchers have begun to study the effectiveness of contrastive learning when *pretraining* on datasets other than ImageNet. In the case of remote sensing, the unique properties of the data have motivated the development of domain-specific contrastive learning techniques [4, 32]. In the medical domain, where images tend to be very dissimilar to ImageNet, it has been shown that contrastive pretraining on domain-specific images leads to significant gains compared to pretraining on ImageNet [11, 46]. [34] compared the representations learned from five different datasets, and showed that in most cases the best performing representations came from pretraining on similar datasets to the downstream task. In the case of fine-grained data, [54] found that contrastive pretraining on images of animals and plants did not lead to superior performance on downstream bird classification compared to pretraining on ImageNet. These apparently conflicting observations may be explained by the relationship between the pretraining and downstream data distributions, which we investigate in our experiments. [65] and [53] pretrained on several different datasets and showed that there was surprisingly little impact on downstream detection and segmentation performance, unless synthetic data was used for pretraining [65].

[50] pretrained on very large datasets (JFT-300M [47] and YFCC100M [49]), but did not observe an improvement over ImageNet pretraining in the standard regime.

We build on the above analysis by performing controlled, like-for-like, comparisons of SSL on several large datasets. This allows us to separate dataset-specific factors from general patterns in SSL performance, and deliver new insights into the necessary conditions for successful pretraining.

Analysis of SSL. A number of works have explored questions related to the conditions under which SSL is successful. [45] showed that self-supervised representations generalize better than supervised ones when the downstream concepts of interest are less semantically similar to the pretraining set. [20] showed that contrastive pretraining on ImageNet performs well on downstream tasks related to object recognition in natural images, while leaving more general study of pretraining in different domains to future work. While these works show that SSL on ImageNet can be effective, our experiments demonstrate that current SSL methods can perform much worse than supervised baselines on non-ImageNet domains, e.g. fine-grained classification.

Existing work has also investigated other aspects of SSL, e.g. [41] examined the invariances learned, [12] showed that easily learned features can inhibit the learning of more discriminative ones, [10, 53, 65] explored the impact of different image augmentations, [12, 53] compared representations from single vs. multi-object images, and [10, 25] varied the backbone model capacity. Most relevant to our work are studies that vary the amount of data in the pretraining dataset, e.g. [34, 53, 61, 65]. We extend this analysis by presenting a more detailed evaluation of the impact of the size of the unlabeled and labeled datasets, and investigate the role of data quality, data domain, and task granularity.

3. Methods

Datasets. We perform experiments on four complementary large-scale datasets: ImageNet [17], iNat21 [53], Places365 [66], and GLC20 [15]. Collectively, these datasets span many important visual properties, including: curated vs. “in-the-wild” images, fine- vs. coarse-grained categories, and object-centric images vs. scenes. Each dataset has at least one million images, which allows us to make fair comparisons against the traditional ImageNet setting. ImageNet (1.3M images, 1k classes) and Places365 (1.8M images, 365 classes) are standard computer vision datasets, so we will not describe them in detail. For ImageNet, we use the classic ILSVRC2012 subset of the full ImageNet-21k dataset. For Places365, we use the official variant “Places365-Standard (small images)” where all images have been resized to 256x256. iNat21 (2.7M images, 10k classes) contains images of plant and animal species and GLC20 (1M images, 16 classes) consists of remote sensing images. As both are recent datasets, we discuss

them in the supplementary material.

Fixed-size subsets. For some experiments we control for dataset size by creating subsampled versions of each dataset with sizes: 1M, 500k, 250k, 125k, and 50k images. We carry out this selection only once, and the images are chosen uniformly at random. We refer to these datasets using the name of the parent dataset followed by the number of images in parentheses, e.g. ImageNet (500k). Note that subsets of increasing size are *nested*, so e.g. ImageNet (500k) includes all of the images in ImageNet (250k). These subsets are also *static* across experiments, e.g. ImageNet (500k) always refers to the same set of 500k images. With the exception of Figures 2 and 3, we use the full dataset for any type of supervised training (i.e. linear evaluation, fine tuning, or supervised training from scratch). We always report results on the same test set for a given dataset, regardless of the training subset used.

Training details. All experiments in this paper are based on a ResNet-50 [28] backbone, which is standard in the contrastive learning literature [8, 10, 27]. We primarily perform experiments on SimCLR [10], a simple and popular contrastive learning method that contains all the building blocks for state-of-the-art self-supervised algorithms. We follow the standard protocol of first training with self-supervision alone and then evaluating the learned features using linear classifiers or end-to-end fine-tuning. Unless otherwise specified, we use hyperparameter settings based on [10] for all methods and datasets. While this may not lead to maximal performance, it is likely to be representative of how these methods are used in practice – due to the high computational cost of contrastive pretraining, extensive hyperparameter tuning is not feasible for most users. We also consider MoCo [27] and BYOL [26] in Figure 3. Full training details are provided in the supplementary material.

4. Experiments

We now describe our experiments in which we investigate the impact of data quantity, data domain, data quality, and task granularity on the success of contrastive learning.

4.1. Data quantity

First we consider the question of how much data is required to learn a “good” representation using SSL. There are two important notions of data quantity: (i) the number of *unlabeled images* used for pretraining and (ii) the number of *labeled images* used to subsequently train a classifier. Since labels are expensive, we would like to learn representations that generalize well with as few labeled images as possible. While unlabeled images are cheap to acquire, they still incur a cost because pretraining time is proportional to the size of the pretraining set. To understand when SSL is cost-effective, we need to understand how performance depends on these two notions of data quantity.

To study this question, we pretrain SimCLR using different numbers of unlabeled images. Each pretrained representation is then evaluated using different numbers of labeled images. In Figure 2 we present these results for iNat21 (left column), ImageNet (center column), and Places365 (right column). We also include results for supervised training from scratch (in black). We show linear evaluation results in the top row and corresponding fine-tuned results in the bottom row. Each curve in a figure corresponds to a different pretrained representation. The points along a curve correspond to different amounts of supervision used to train a linear classifier or fine-tune the network.

There is little benefit beyond 500k pretraining images. The gap between the 500k (blue) and 1M (orange) pretraining image curves is typically less than 1-2% in top-1 accuracy. This means that for a dataset with one million images, we can trade a small decrease in accuracy for a 50% decrease in pretraining time. If a 2-4% top-1 accuracy drop is acceptable, then the pretraining set size can be reduced by a factor of four (from 1M to 250k). However, the difference between 50k (pink) pretraining images and 250k (green) pretraining images is substantial for each dataset, often in excess of 10% top-1 accuracy. We conclude that SimCLR seems to saturate well before we get to ImageNet-sized pretraining sets. This is consistent with observations from the supervised learning literature, though more images are required to reach saturation [37].

Self-supervised pretraining can be a good initializer when there is limited supervision available. In the bottom row of Figure 2 we see that when only 10k or 50k labeled images are available, fine-tuning a SimCLR representation is significantly better than training from scratch. When supervision is plentiful, fine-tuned SimCLR representations achieve performance similar to supervised training from scratch. It is interesting to compare this to findings from the supervised setting which suggest that networks which are initially trained on distorted (i.e. augmented) images are unable to recover when subsequently trained with undistorted ones [3].

Self-supervised representations can approach fully supervised performance for some datasets, but only by using lots of labeled images. The ultimate goal of SSL is to match supervised performance without the need for large amounts of labeled data. Suppose we consider the right-most point on the black curves in Figure 2 as a proxy for “good” supervised performance. Then in both the linear and fine-tuned cases, the gap between SimCLR (pretrained on 1M images) and “good” supervised performance is quite large unless well over 100k labeled images are used. For instance, the gap between “good” supervised performance and a classifier trained using 50k labeled images on top of SimCLR (1M) is around 11% (11%) for Places365, 23% (21%) for ImageNet, and 58% (56%) for iNat21 in the lin-

ear (and fine-tuned) case. Although SSL works well when lots of supervision is available, further innovation is needed to improve the utility of self-supervised representations in the low-to-moderate supervision regime.

iNat21 is a valuable SSL benchmark. Figure 2 shows a surprisingly large gap ($\sim 30\%$) between supervised and self-supervised performance on iNat21 in the high supervision regime. In Figure 3 we see that other SSL methods exhibit similar limitations. The newer BYOL outperforms MoCo and SimCLR, but a considerable gap ($\sim 25\%$) remains. The high supervised performance shows that the task is possible, yet the self-supervised performance remains low. It seems that iNat21 reveals challenges for SSL that are not apparent in ImageNet, and we believe it is a valuable benchmark for future SSL research.

4.2. Data domain

In the previous section we observed that increasing the pretraining set size yields rapidly diminishing returns. In this section we consider a different design choice: *what kind of images* should we use for pretraining? Since most contrastive learning papers only pretrain on ImageNet, this question has not received much attention. We take an initial step towards an answer by studying the properties of SimCLR representations derived from four pretraining sets drawn from different domains.

We train SimCLR on iNat21 (1M), ImageNet (1M), Places365 (1M), and GLC20 (1M). By holding the pretraining set size constant, we aim to isolate the impact of the different visual domains. We present in-domain and cross-domain linear evaluation results for each representation in Table 1. In Table 2 we consider the effect of pretraining on *pooled datasets*, i.e. new image collections built by shuffling together existing datasets. Finally, in Table 3 we study different *fused representations*, which are formed by concatenating the outputs of different feature extractors.

Pretraining domain matters. In Table 1 we see that in-domain pretraining (diagonal entries) consistently beats cross-domain pretraining (off-diagonal entries). The gap can be surprisingly large, e.g. in-domain pretraining provides a 12% boost on iNat21 compared to the best cross-domain pretraining (ImageNet). One might have expected that a visually diverse dataset like ImageNet would lead to a better self-supervised representation than a more homogeneous dataset like GLC20 (even when evaluating on GLC20) but this is not what we observe.

The off-diagonal entries of Table 1 show that training SimCLR on ImageNet leads to the best cross-domain performance, while GLC20 leads to the worst cross-domain performance. Since the pretraining protocols and dataset sizes are held constant, we suggest that the characteristics of the image sets themselves are responsible for the differences we observe. The strong cross-domain performance of

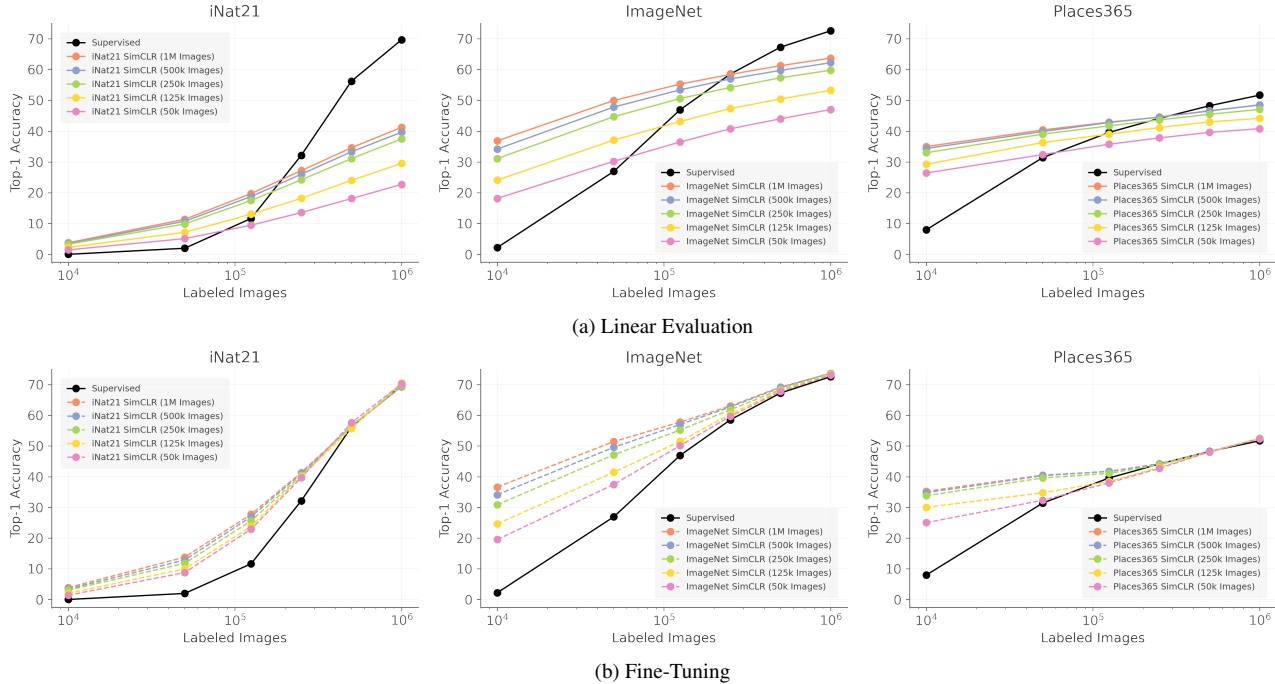


Figure 2. **How much data does SimCLR need?** Linear evaluation results (top row) and fine-tuning results (bottom row) as a function of the number of *unlabeled images* used for pretraining and the number of *labeled images* used for downstream supervised training. The “Supervised” curve (black) corresponds to training from scratch on different numbers of labeled images. It is the same for the top and bottom plots in each column. Most SSL papers focus on the “high data” regime, using $\sim 10^6$ images (e.g. all of ImageNet) for both pretraining and classifier supervision, but there are significant opportunities for improvement in the “low-data” regime. Even with 10^6 labeled images for linear classifier training, SimCLR performs far worse than supervised learning on iNat21, suggesting that iNat21 could be a more useful SSL benchmark than ImageNet in future.

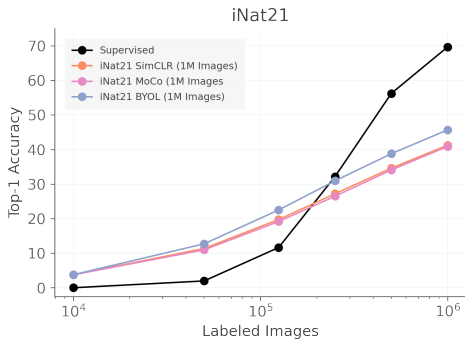


Figure 3. **How does SimCLR compare to other self-supervised methods?** Linear evaluation results on iNat21 for SimCLR, MoCo, and BYOL. All methods are pretrained on 1M images for 1000 epochs and follow the same linear evaluation protocol. The more recent BYOL performs better than the others, but a large gap remains to supervised performance.

SimCLR pretrained on ImageNet may be due to *semantic similarity* – perhaps it is better to pretrain on a dataset that is semantically similar to the downstream task, even in a self-supervised context. This makes sense because there are classes in ImageNet that are similar to classes in iNat21 (animals) and Places365 (scenes). This also explains the weak performance of GLC20, since remote sensing imagery is

Pretraining	iNat21	ImageNet	Places365	GLC20
iNat21 (1M) SimCLR	0.493	0.519	0.416	0.707
ImageNet (1M) SimCLR	0.373	0.644	0.486	0.716
Places365 (1M) SimCLR	0.292	0.491	0.501	0.693
GLC20 (1M) SimCLR	0.187	0.372	0.329	0.769
Supervised (All Images)	0.791	0.741	0.539	0.826

Table 1. **Does pretraining domain matter?** Linear evaluation results for representations derived from different million-image datasets. We train the linear classifiers using the full training sets. The results in the “Supervised” row correspond to supervised training from scratch on the full training set. We report MAP for GLC20 and top-1 accuracy for other datasets. In all cases, in-domain pretraining outperforms cross-domain pretraining. In each column we highlight the **best** and second-best results.

not similar to the other datasets.

Adding cross-domain pretraining data does not necessarily lead to more general representations. We have seen that pretraining on different domains leads to representations with significantly differing capabilities. This leads to a natural question: *what happens if we combine our datasets and then learn a representation?*

Table 2 gives linear evaluation results for SimCLR pretrained on different “pooled” datasets. In each row, n images from dataset A and m images from dataset B are shuffled together to produce a pretraining set of size $n + m$. For instance, the pretraining dataset in the first row of Table 2

Pretraining				Evaluation		
<i>iNat21</i>	<i>ImageNet</i>	<i>Places365</i>	<i>GLC20</i>	<i>iNat21</i>	<i>ImageNet</i>	<i>Places365</i>
250k	250k	-	-	0.444	0.597	0.467
-	250k	250k	-	0.334	0.596	0.490
250k	-	250k	-	0.428	0.531	0.483
250k	250k	250k	250k	0.410	0.574	0.482
In-Domain (250k)				0.451	0.608	0.485
In-Domain (500k)				0.477	0.629	0.499
In-Domain (1M)				0.493	0.644	0.501

Table 2. **The effect of dataset pooling.** Linear evaluation results for self-supervised representations derived from *pooled datasets*, where two or more datasets are shuffled together. We train the linear classifiers using the full training sets. The “In-Domain” results correspond to pretraining on subsets of the dataset named at the top of the column. Pooling datasets increases pretraining set size and diversity, but we find that performance *decreases* relative to comparable in-domain pretraining. The “In-Domain (1M)” row corresponds to the diagonal entries of Table 1.

consists of 250k iNat21 images and 250k ImageNet images shuffled together.

If we compare the “In-Domain (500k)” row against the (equally sized) pooled datasets in the first three rows of Table 2, we see that the in-domain pretraining on 500k images is always better. Similarly, the “In-Domain (1M)” row beats the 1M-image pooled dataset (consisting of 250k images from the four datasets). The more diverse pooled pretraining sets always lead to worse performance compared to the more homogeneous pretraining sets of the same size.

Table 2 also allows us to say whether it is worthwhile to *add* pretraining data from a different domain (as opposed to swapping out some in-domain data for some data from a different domain, as we have been discussing so far). The “In-Domain (250k)” row is better than the 1M-image pooled dataset and almost all of the 500k-image pooled datasets. It seems that adding pretraining data from a different domain typically *hurts* performance. In contrast, Figure 2 shows that increasing the amount of *in-domain* pretraining data consistently improves performance.

We hypothesize that the reason for this lackluster performance is that diverse images are easier to tell apart, which makes the contrastive pretext task easier. If the contrastive task is too easy, the quality of the representation suffers [6, 12]. While more investigation is needed, the fact that increasing pretraining data diversity can hurt performance suggests a “diversity-difficulty trade-off” that should be considered when creating pretraining sets for SSL.

Self-supervised representations can be largely redundant. From Table 1 it is clear that pretraining on different datasets leads to representations that differ significantly. For instance, iNat21 SimCLR beats ImageNet SimCLR on iNat21 (+12.4%) and ImageNet SimCLR beats iNat21 SimCLR on ImageNet (+12.7%). Do these representations learn complementary information, or do they just capture the same information to different degrees?

ImageNet	iNat21	Dim.	ImageNet	iNat21
SimCLR	-	2048	0.647	0.380
-	SimCLR	2048	0.520	0.506
Sup.	-	2048	0.711	0.434
-	Sup.	2048	0.490	<u>0.769</u>
Sup.	Sup.	4096	0.712	0.772
SimCLR	SimCLR	4096	0.641	0.520
SimCLR & Sup.	-	4096	0.720	0.472
-	SimCLR & Sup.	4096	0.527	0.772
SimCLR	Sup.	4096	0.605	<u>0.769</u>
Sup.	SimCLR	4096	<u>0.717</u>	0.553

Table 3. **The effect of representation fusion.** Linear evaluation results for different combinations of supervised and self-supervised representations on ImageNet and iNat21. We train the linear classifiers using the full training sets. For comparability, the in-domain supervised results in this table (ImageNet Sup. evaluated on ImageNet and iNat21 Sup. evaluated on iNat21) are for linear classifiers trained on representations learned from full supervision. “Dim.” is the representation dimensionality. In each column we highlight the **best** and second-best results.

To probe this question we concatenate features from different pretrained networks and carry out linear evaluation on these “fused” representations. In Table 3 we present linear evaluation results for fused representations on ImageNet and iNat21. Combining ImageNet SimCLR and iNat21 SimCLR is worse than ImageNet SimCLR alone on ImageNet (-0.6%), but better than iNat21 SimCLR alone on iNat21 (+1.4%). These effects are small relative to the > 12% difference between ImageNet SimCLR and iNat21 SimCLR. This suggests that the two self-supervised representations are largely redundant.

There is a larger effect when combining supervised and self-supervised representations. For iNat21, adding ImageNet Sup. (i.e. supervised ImageNet features) on top of iNat21 SimCLR improves performance significantly (+4.7%). However, adding iNat21 Sup. on top of ImageNet SimCLR actually decreases performance (-4.2%). These results are consistent with the hypothesis that dataset semantics are important even for SSL. Since ImageNet is semantically broader than iNat21 (ImageNet has animal classes, but also many other things), features learned from ImageNet (supervised or self-supervised) should be more helpful for iNat21 than vice-versa.

4.3. Data quality

We have seen that the characteristics of the pretraining data can have a significant impact on the quality of self-supervised representations. In this section we dig deeper into this question by studying the impact of pretraining on artificially degraded images. This serves two purposes. First, this is a practical question since there are many settings where image quality issues are pervasive e.g. medical imaging [48] or camera trap data [5]. Second, it can help us understand the robustness properties of SSL.

To create a corrupted dataset we apply a particular image

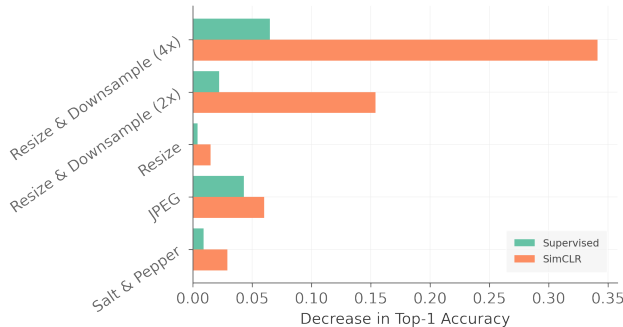


Figure 4. **What is the effect of pretraining image corruption?** Decrease in linear evaluation accuracy on ImageNet due to pretraining on corrupted versions of the ImageNet training set. The zero point corresponds to pretraining (supervised or SimCLR) on uncorrupted images followed by linear evaluation. “Supervised” and “SimCLR” have different zero points. All linear classifiers are trained using the full uncorrupted ImageNet training set.

corruption to each image in the dataset. This is a one-time offline preprocessing step, so corruptions that have a random component are realized only once per image. Given a corrupted dataset we then pretrain as normal. During linear evaluation, we use the original clean images for training and testing, i.e. the corrupted images are only used for pretraining.

In Figure 4 we present linear evaluation results on ImageNet for a simple but diverse set of corruptions. The zero point corresponds to pretraining on uncorrupted images, and we measure how much performance drops when pretraining on corrupted images. The “Salt and Pepper” corruption is salt and pepper noise applied independently to each pixel, in each channel, with probability 0.01. The “JPEG” corruption is JPEG compression with a very low quality level of 10. For “Resize”, we resize each image so that the short side is 256 pixels while preserving the aspect ratio. This reduces the resolution of the crops used for training. For our downsampling corruptions, we follow the resize operation with downsampling by 2x or 4x and then upsampling by the same factor. This holds constant the image size and the fraction of the image occupied by each object, but reduces resolution. Implementation details and examples can be found in the supplementary.

Image resolution is critical for SSL. “Downsample (2x)” and “Downsample (4x)” are by far the most damaging corruptions for SimCLR, reducing accuracy by around 15% and 34%, respectively. Since SimCLR already involves extreme cropping, we might expect more robustness to changes in image resolution. This finding could be partially explained by the difficulty of generalizing to higher-resolution images during linear classifier training [52]. However, supervised pretraining faces the same challenge but the effect of downsampling is much less dramatic. This

suggests that the performance drop is due to deficiencies in the features learned by SimCLR.

SSL is relatively robust to high-frequency noise. “JPEG” and “Salt & Pepper” both add high-frequency noise to the image. For SimCLR, these corruptions have a much milder impact than the downsampling corruptions. One possible explanation is that downsampling destroys texture information, which is known to be a particularly important signal for convolutional neural networks [21, 31]. For supervised pretraining the ranking of corruptions is very different, with “JPEG” landing between 2x and 4x downsampling.

4.4. Task granularity

We have seen that the properties of pretraining datasets are important for determining the utility of self-supervised representations. But are there downstream tasks for which self-supervised representations are particularly well or poorly suited? We consider *fine-grained classification* and show that classification performance depends on *task granularity*, i.e. how fine or coarse the labels are. While there are formal methods for measuring dataset granularity [16], we claim by intuition that iNat21 is more fine-grained than ImageNet, which is more fine-grained than Places365.

In Figure 5 we use label hierarchies (which are available for ImageNet, iNat21, and Places365) to explicitly study how performance depends on label granularity. We treat “distance from the root of the hierarchy” as a proxy for granularity, so labels further from the root are considered to be more fine-grained. We perform (i) linear classifier training (for SimCLR) and (ii) end-to-end training from scratch (for “Supervised”) using the labels at the finest level of the taxonomy and re-compute accuracy values as we progressively coarsen the predictions and labels. We do not re-train at each level of granularity. A complete description of this process can be found in the supplementary materials.

The performance gap between SSL and supervised learning grows as task granularity becomes finer. We start with the iNat21 results in Figure 5. The supervised and SimCLR pretrained models perform similarly at the coarsest levels of the label hierarchy (“Kingdom”). Both models perform worse as task granularity increases, but the SimCLR model degrades much more rapidly (“Species”). This suggests that SimCLR may fail to capture fine-grained semantic information as effectively as supervised pretraining. We also observe a growing supervised/self-supervised gap for ImageNet and Places365. The magnitude of this gap seems to track dataset granularity, since iNat21 (most fine-grained) has the largest gap and Places365 (least fine-grained) has the smallest gap. The fact that supervised learning achieves high performance on iNat21 while SSL lags behind suggests that iNat21 could be a valuable benchmark dataset for the next phase of SSL research.

Are the augmentations destructive? State-of-the-art con-

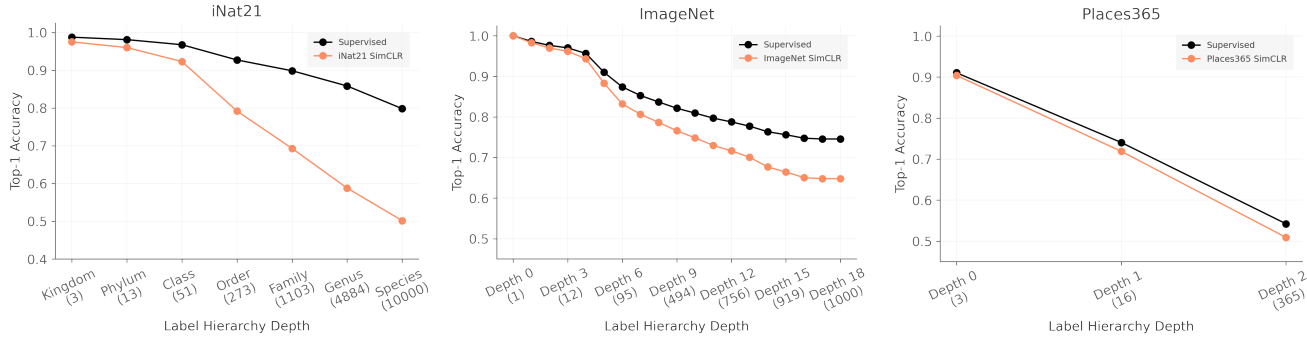


Figure 5. **How does performance depend on label granularity?** Linear evaluation at different levels of label granularity for iNat21, ImageNet, and Places365. Each plot compares supervised learning from scratch against a linear classifier trained on top of in-domain SimCLR. Both are trained using the full training sets. We plot top-1 accuracy against label granularity, which is more fine-grained as we move from left to right. The numbers on the x -axis are the class counts at a given level of the label hierarchy. We do not re-train at coarser granularity levels, we just change the evaluation label set. The definitions of the hierarchy levels are given in the supplementary material.

trastive learning techniques are designed for ImageNet, so the default augmentation policy may be poorly tuned for other datasets [60]. For instance, if color is a key fine-grained feature for species classification then the “color jitter” augmentation used by SimCLR may destroy important information for iNat21 classification. Could this explain the rapid drop in performance exhibited by iNat21 SimCLR for fine-grained classes? Notice that there is a similar, though less extreme, fine-grained performance drop for ImageNet SimCLR in Figure 5. Since the ImageNet-tuned augmentations are presumably not destructive for ImageNet, it does not seem likely that this fully explain our observations.

Does contrastive learning have a coarse-grained bias?

We hypothesize that the contrastive loss tends to cluster images based on overall visual similarity. The intuition is that fine-grained features are often subtle, and subtle features are unlikely to be very useful for distinguishing between pairs of images in the contrastive pretext task. If our hypothesis is correct then the boundaries between different clusters would not be well-aligned with the boundaries between fine-grained classes. This effect could be overlooked when evaluating on coarse-grained classes, but would become apparent on a more fine-grained task. Additional analysis is required to fully understand this “granularity gap” in SSL, which we leave to future work.

5. Conclusion

We have presented a comprehensive set of experiments to address several aspects of the question: *when does contrastive visual representation learning work?* In Section 4.1 we found that we need fewer than 500k pretraining images before encountering severe diminishing returns. However, even the best self-supervised representations are still much worse than peak supervised performance without hundreds of thousands of labeled images for classifier training. In Section 4.2 we found that self-supervised pretraining on 1M images from different domains results in representations

with very different capabilities, and that simple methods for combining different datasets do not lead to large gains. In Section 4.3 we showed that image resolution is critical for contrastive learning and, more broadly, that some image corruptions can degrade a self-supervised representation to the point of unusability while others have almost no impact. Finally, in Section 4.4 we found that supervised pretraining retains a substantial edge when it comes to fine-grained classification. These experiments highlight several areas where further research is needed to improve current SSL algorithms, most of which were not evident from traditional evaluation protocols, i.e. top-1 accuracy on ImageNet.

Limitations. We mainly perform experiments using one self-supervised method. We focus on SimCLR because it reflects the essence of state-of-the-art contrastive learning methods without introducing additional architectural complexities. While our MoCo and BYOL experiments are not much different from SimCLR, it is important to validate our results on other self-supervised methods. It would also be interesting to explore alternative backbone architectures [9, 19], though after controlling for training settings, ResNet-50 remains competitive with newer architectures [58, 59]. We study only classification tasks, so additional work is also required to understand how these results translate to segmentation [57] or detection [29, 68]. Finally, we only consider datasets up to roughly ImageNet scale. We believe this is the most practical setting for most use cases, but it is possible that some patterns may be different for significantly larger datasets and models [23, 24].

Acknowledgements. We thank Mason McGill for detailed feedback, and Grant Van Horn, Christine Kaeser-Chen, Yin Cui, Sergey Ioffe, Pietro Perona, and the rest of the Perona Lab for insightful discussions. This work was supported by the Caltech Resnick Sustainability Institute, an NSF Graduate Research Fellowship (grant number DGE1745301), and the Pioneer Centre for AI (DNRF grant number P1).

References

- [1] Pillow. <https://python-pillow.org/>. 17
- [2] Wordnet interface. <https://www.nltk.org/howto/wordnet.html>. 16
- [3] Alessandro Achille, Matteo Rovere, and Stefano Soatto. Critical learning periods in deep neural networks. In *ICLR*, 2019. 4
- [4] Kumar Ayush, Burak Uzkent, Chenlin Meng, Kumar Tanmay, Marshall Burke, David Lobell, and Stefano Ermon. Geography-aware self-supervised learning. In *ICCV*, 2021. 2
- [5] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *ECCV*, 2018. 6
- [6] Tiffany Tianhui Cai, Jonathan Frankle, David J Schwab, and Ari S Morcos. Are all negatives created equal in contrastive instance discrimination? *arXiv:2010.06682*, 2020. 2, 6
- [7] Yue Cao, Zhenda Xie, Bin Liu, Yutong Lin, Zheng Zhang, and Han Hu. Parametric instance classification for unsupervised visual feature learning. In *NeurIPS*, 2020. 2
- [8] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *NeurIPS*, 2020. 2, 3
- [9] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. 8
- [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 2, 3, 16, 17
- [11] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. In *NeurIPS*, 2020. 2
- [12] Ting Chen, Calvin Luo, and Lala Li. Intriguing properties of contrastive losses. In *NeurIPS*, 2021. 3, 6
- [13] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv:2003.04297*, 2020. 2, 16
- [14] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, 2021. 2
- [15] Elijah Cole, Benjamin Deneu, Titouan Lorieul, Maximilien Servajean, Christophe Botella, Dan Morris, Nebojsa Jojic, Pierre Bonnet, and Alexis Joly. The geolifeclaf 2020 dataset. *arXiv:2004.04192*, 2020. 3, 16
- [16] Yin Cui, Zeqi Gu, Dhruv Mahajan, Laurens Van Der Maaten, Serge Belongie, and Ser-Nam Lim. Measuring dataset granularity. *arXiv:1912.10154*, 2019. 7
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 3
- [18] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, 2015. 2
- [19] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 8
- [20] Linus Ericsson, Henry Gouk, and Timothy M Hospedales. How well do self-supervised models transfer? In *CVPR*, 2021. 1, 3
- [21] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. In *ICLR*, 2019. 7
- [22] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *ICLR*, 2018. 2
- [23] Priya Goyal, Mathilde Caron, Benjamin Lefaudeux, Min Xu, Pengchao Wang, Vivek Pai, Mannat Singh, Vitaliy Liptchinsky, Ishan Misra, Armand Joulin, et al. Self-supervised pretraining of visual features in the wild. *arXiv preprint arXiv:2103.01988*, 2021. 8
- [24] Priya Goyal, Quentin Duval, Isaac Seessel, Mathilde Caron, Mannat Singh, Ishan Misra, Levent Sagun, Armand Joulin, and Piotr Bojanowski. Vision models are more robust and fair when pretrained on uncurated images without supervision. *arXiv preprint arXiv:2202.08360*, 2022. 8
- [25] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan Misra. Scaling and benchmarking self-supervised visual representation learning. In *ICCV*, 2019. 3
- [26] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. In *NeurIPS*, 2020. 2, 3, 16
- [27] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 2, 3, 16
- [28] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3
- [29] Olivier J Hénaff, Skanda Koppula, Jean-Baptiste Alayrac, Aaron van den Oord, Oriol Vinyals, and João Carreira. Efficient visual pretraining with contrastive detection. In *ICCV*, 2021. 8
- [30] Olivier J Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. *arXiv:1905.09272*, 2019. 2
- [31] Katherine L Hermann, Ting Chen, and Simon Kornblith. The origins and prevalence of texture bias in convolutional neural networks. In *NeurIPS*, 2020. 7
- [32] Jian Kang, Ruben Fernandez-Beltran, Puhong Duan, Sicong Liu, and Antonio J Plaza. Deep unsupervised embedding for remotely sensed images based on spatially augmented momentum contrast. *Transactions on Geoscience and Remote Sensing*, 2020. 2
- [33] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and

- Dilip Krishnan. Supervised contrastive learning. In *NeurIPS*, 2020. 2
- [34] Klemen Kotar, Gabriel Ilharco, Ludwig Schmidt, Kiana Ehsani, and Roozbeh Mottaghi. Contrasting contrastive self-supervised representation learning pipelines. In *ICCV*, 2021. 2, 3
- [35] Hsin-Ying Lee, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Unsupervised representation learning by sorting sequences. In *ICCV*, 2017. 2
- [36] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv:1608.03983*, 2016. 16
- [37] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, 2018. 2, 4
- [38] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, 2016. 2
- [39] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv:1807.03748*, 2019. 2
- [40] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016. 2
- [41] Senthil Purushwalkam and Abhinav Gupta. Demystifying contrastive self-supervised learning: Invariances, augmentations and dataset biases. In *NeurIPS*, 2020. 3
- [42] Nils Rethmeier and Isabelle Augenstein. Long-tail zero and few-shot learning via contrastive pretraining on and for small data. *arXiv:2010.01061*, 2020. 2
- [43] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 1
- [44] Aaqib Saeed, David Grangier, and Neil Zeghidour. Contrastive learning of general-purpose audio representations. In *ICASSP*, 2021. 2
- [45] Mert Bulent Sariyildiz, Yannis Kalantidis, Diane Larlus, and Karteek Alahari. Concept generalization in visual representation learning. In *ICCV*, 2021. 3
- [46] Hari Sowrirajan, Jingbo Yang, Andrew Y Ng, and Pranav Rajpurkar. Moco pretraining improves representation and transferability of chest x-ray models. In *Medical Imaging with Deep Learning*, 2021. 2
- [47] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *ICCV*, 2017. 3
- [48] Siyi Tang, Amirata Ghorbani, Rikiya Yamashita, Sameer Rehman, Jared A Dunnmon, James Zou, and Daniel L Rubin. Data valuation for medical imaging using shapley value and application to a large-scale chest x-ray dataset. *Scientific reports*, 2021. 6
- [49] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, 2016. 3
- [50] Yonglong Tian, Olivier J Henaff, and Aaron van den Oord. Divide and contrast: Self-supervised learning from uncurated data. In *ICCV*, 2021. 3
- [51] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *ECCV*, 2020. 2
- [52] Hugo Touvron, Andrea Vedaldi, Matthijs Douze, and Hervé Jégou. Fixing the train-test resolution discrepancy. In *NeurIPS*, 2019. 7
- [53] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Revisiting contrastive methods for unsupervised learning of visual representations. In *NeurIPS*, 2021. 2, 3
- [54] Grant Van Horn, Elijah Cole, Sara Beery, Kimberly Wilber, Serge Belongie, and Oisín Mac Aodha. Benchmarking representation learning for natural world image collections. In *CVPR*, 2021. 2, 12, 15, 16
- [55] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The iNaturalist Species Classification and Detection Dataset. In *CVPR*, 2018. 15
- [56] Ali Varamesh, Ali Diba, Tinne Tuytelaars, and Luc Van Gool. Self-supervised ranking for representation learning. *arXiv:2010.07258*, 2020. 2
- [57] Wenguan Wang, Tianfei Zhou, Fisher Yu, Jifeng Dai, Ender Konukoglu, and Luc Van Gool. Exploring cross-image pixel contrast for semantic segmentation. In *ICCV*, 2021. 8
- [58] Ross Wightman, Hugo Touvron, and Hervé Jégou. Resnet strikes back: An improved training procedure in timm. *arXiv:2110.00476*, 2021. 8
- [59] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick. Early convolutions help transformers see better. In *NeurIPS*, 2021. 8
- [60] Tete Xiao, Xiaolong Wang, Alexei A Efros, and Trevor Darrell. What should not be contrastive in contrastive learning. In *ICLR*, 2020. 2, 8
- [61] Xingyi Yang, Xuehai He, Yuxiao Liang, Yue Yang, Shanghang Zhang, and Pengtao Xie. Transfer learning or self-supervised learning? a tale of two pretraining paradigms. *arXiv:2007.04234*, 2020. 3
- [62] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv:1708.03888*, 2017. 16
- [63] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *ECCV*, 2016. 2
- [64] Richard Zhang, Phillip Isola, and Alexei A Efros. Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In *CVPR*, 2017. 2
- [65] Nanxuan Zhao, Zhirong Wu, Rynson WH Lau, and Stephen Lin. What makes instance discrimination good for transfer learning? In *ICLR*, 2021. 2, 3
- [66] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. 3
- [67] Benjin Zhu, Junqiang Huang, Zeming Li, Xiangyu Zhang, and Jian Sun. Eqco: Equivalent rules for self-supervised contrastive learning. *arXiv:2010.01929*, 2020. 2

- [68] Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin D Cubuk, and Quoc V Le. Rethinking pre-training and self-training. In *NeurIPS*, 2020. 8

A. Additional Results

A.1. How does task granularity affect different self-supervised learning methods?

In Figure 5 we saw that there is a large gap between supervised and self-supervised (SimCLR) performance on iNat21. Figure A1 extends Figure 5 by adding results for MoCo and BYOL. Across all granularity levels, MoCo is slightly worse than SimCLR and BYOL is significantly better. For all three self-supervised methods, performance drops rapidly as the evaluation is made more fine-grained. While BYOL is much better than SimCLR, it still lags behind fully supervised performance by 20% top-1 accuracy.

A.2. Do larger models scale better in terms of pre-training set size?

In Figure 2 we observe that doubling the pretraining set size from 500k images to 1M images leads to small benefits (1-2%) across three large-scale datasets. However, all of those results are based on a ResNet-50. Does the story change for larger or smaller models? In Figure A2 we study this question using ResNet-34, ResNet-50, and ResNet-101. When we double the size of the pretraining set from 125k to 250k, ResNet-50 and ResNet-101 make significantly larger gains than ResNet-34. However, doubling the size of the pretraining set from 500k to 1M produces gains of <2% for all models. While ResNet-101 gains more than ResNet-50 with each increase in pretraining set size, the gap between them is very small by the time we reach 1M images. This is the same conclusion we reached in Figure 2.

A.3. Does semantic similarity explain patterns in self-supervised performance?

In Section 4.2 we saw that (i) in-domain SimCLR pretraining always beats cross-domain SimCLR pretraining and (ii) ImageNet is the best dataset for cross-domain pretraining. One hypothesis which could explain these patterns is that *semantic similarity between the pretraining dataset and the downstream task leads to better performance*. This would require that modern self-supervised methods capture high-level semantic information. In this section we consider evidence for this hypothesis.

ImageNet SimCLR performs well on iNat21 classes that are similar to ImageNet classes. ImageNet includes around 200 mammal categories, 60 bird categories, and 30 categories of insects and reptiles. A breakdown of the categories in iNat21 can be found in [54]. In Figure A3 we analyze per-categories accuracy averaged over six *taxonomic classes* of animals (Arachnida, Insecta, Amphibia, Mammalia, Reptilia) and two taxonomic classes of plants (Liliopsida and Magnoliopsida). Surprisingly, ImageNet SimCLR outperforms iNat21 SimCLR on mammals (Mammalia) and nearly matches the performance of iNat21 Sim-

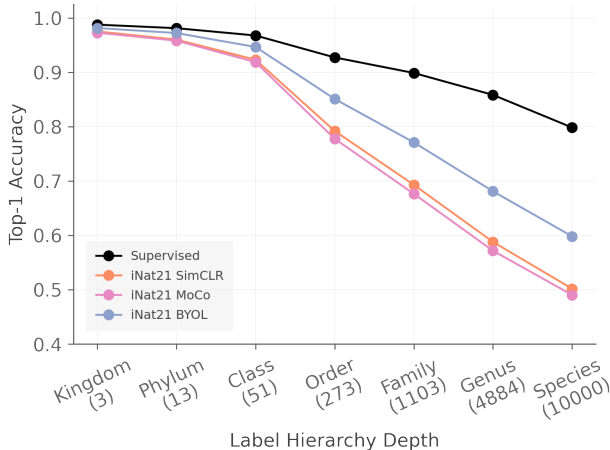


Figure A1. **How does performance depend on label granularity?** Linear evaluation at different levels of label granularity for iNat21. We compare end-to-end training from scratch against linear classifiers trained on top of in-domain self-supervised representations (SimCLR, MoCo, and BYOL). All classifiers (linear and end-to-end) are trained using the full iNat21 training set. This plot is identical to Figure 5 except that we have added curves for MoCo and BYOL.

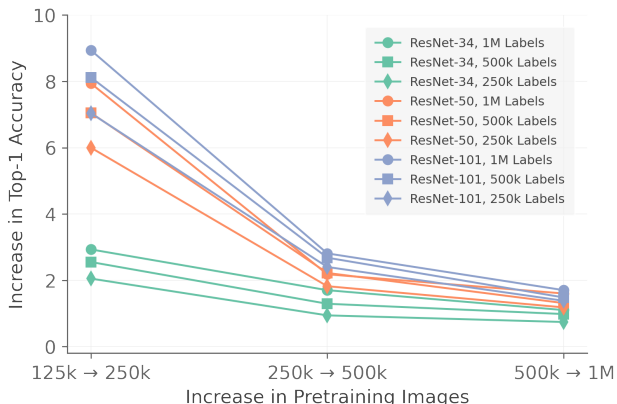


Figure A2. **Increasing pretraining set size leads to rapid diminishing returns across different model sizes.** Linear evaluation results on iNat21 for SimCLR. We show the increase in top-1 accuracy on iNat21 that results from doubling pretraining set size. Each color is a different architecture. For a given color, each line uses a different amount of labeled data for linear classifier training.

CLR on birds (Aves). We also evaluate Places365 SimCLR pretraining, which does not have any categories corresponding to animals or plants. We do not see any taxonomic classes for which Places365 SimCLR performs close to iNat21 SimCLR.

Most of the ImageNet classes for which iNat21 SimCLR beats ImageNet SimCLR are animals or plants. We find similar effects in the context of ImageNet classification.

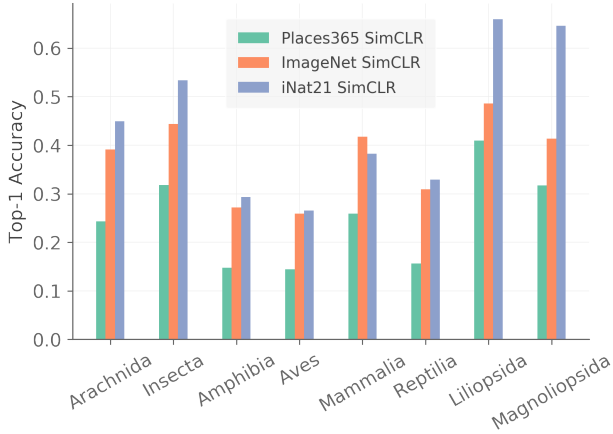


Figure A3. **Semantic similarity may predict transfer performance.** Linear evaluation results on iNat21 for different pretrained representations. We compare representations pretrained on Places365, ImageNet, and iNat21 (full datasets, not subsampled) in terms of top-1 linear classification accuracy on iNat21. The result for each taxonomic class (Arachnida, Insecta, Amphibia, Aves, Mammalia, Reptilia, Liliopsida, Magnoliopsida) is the average of the per-species accuracy over all species in that taxonomic class.

When we compare the per-category accuracy for ImageNet SimCLR with the per-category accuracy for iNat21 SimCLR, we find that ImageNet SimCLR leads to a higher accuracy for all but 80 categories. Of those 80 categories, 68 (i.e. 85%) are animals or plants.

In-domain pretraining helps some classes and hurts others. To develop a deeper understanding of the effect of pretraining domain, we compute the per-class accuracy improvement that results from using in-domain SimCLR instead of ImageNet SimCLR. We present these results for iNat21 and Places365 in Figure A4. For these results we pretrain on the full datasets, not the million-image subsets. We see that in-domain pretraining leads to an improvement for $\sim 60\%$ of classes, while the rest stay the same or degrade. In Table A1 we list the most harmed and most improved classes. Interestingly, all of the most improved classes for iNat21 are plants. Around 40% of the images in iNat21 are plants, but of course the self-supervised method does not have access to the labels. We also notice that many of the most harmed classes for iNat21 are similar to classes we might find in ImageNet, e.g. birds, mammals, and reptiles. This is consistent with the hypothesis that the success of SimCLR is partially governed by the semantic similarity between the pretraining set and the downstream task, even though no labels are used for representations learning. The patterns seem less clear for Places365.



Figure A4. **In-domain contrastive learning improves accuracy on most (but not all) classes.** Increase in per-class linear evaluation results for different pretrained representations compared to an ImageNet SimCLR baseline. In Table 1 we saw that in-domain pretraining was better than cross-domain pretraining. Here we break down those results in terms of the per-class accuracy increase for in-domain SimCLR with respect to ImageNet SimCLR (represented by the dashed line). For both Places365 (green line) and iNat21 (orange line), in-domain SimCLR pretraining benefits around 60% of classes while around 40% of classes are either the same or worse off. Note that the curves for Places365 and iNat21 are sorted independently so the species ordering is different for each. See Table A1 for lists of the most harmed and most improved classes for both datasets.

A.4. Is SimCLR overly tuned for ImageNet?

One possible explanation for the strong cross-domain performance of ImageNet SimCLR we observe in Table 1 is that the training procedures, augmentations, and hyperparameters for SimCLR were designed with ImageNet in mind. This might lead SimCLR to produce better representations when trained on ImageNet than it does when trained on other datasets. However, we see that in-domain SimCLR is better than ImageNet SimCLR for iNat21, Places365, and GLC20. If SimCLR is somehow “overfit” to ImageNet, that effect seems to be overwhelmed by the effect of domain similarity.

A.5. What is the effect of native image resolution?

ImageNet and iNat21 have larger images than Places365 and GLC20. While images are always resized to 224x224 before they are passed in to the network, that happens after random crops are chosen. This means that we are training on more detailed 224x224 images for ImageNet and iNat21 compared to Places365 and GLC20. This could affect cross-domain performance comparisons such as those in e.g. Table 1. To understand the impact of this difference, we compare pretraining on resized images to pretraining on the original images for ImageNet and iNat21. We provide linear evaluation results in Table A2. It seems that resizing

iNat21		Places365	
Most Improved	Most Harmed	Most Improved	Most Harmed
summer-cypress	Ferruginous Hawk	/a/airport_terminal	/s/slum
Greater Tickseed	Western Banded Gecko	/r/roof_garden	/h/home_office
Annual Blue-eyed Grass	Desert Cottontail	/r/restaurant	/s/swamp
Jamaica Snakeweed	Arizona Alligator Lizard	/g/gazebo/interior	/a/arena/performance
California Jacob's ladder	Petticoat Mottlegill	/b/bedroom	/b/beach
tomato	Elk	/b/booth/indoor	/c/canal/urban
leatherleaf fern	Ruddy Ground-Dove	/r/rice_paddy	/o/orchard
northern bugleweed	Long-tailed Weasel	/c/castle	/o/ocean
mock azalea	Little Blue Dragonlet	/m/museum/indoor	/g/garage/indoor
Mexican False Calico	Signal Crayfish	/l/locker_room	/u/underwater/ocean_deep

Table A1. **In-domain pretraining helps some classes and harms others.** Lists of the ten most improved and the ten most harmed classes when we change from ImageNet SimCLR pretraining to in-domain SimCLR pretraining. See Figure A4 for the corresponding curves showing the distribution of accuracy improvement over all classes.

Pretraining	iNat21	ImageNet	Places365	GLC20
iNat21	0.506	0.520	0.413	0.865
iNat21 (Resize)	0.505	0.500	0.412	0.865
Change	-0.001	-0.020	-0.001	0.000
ImageNet	0.380	0.647	0.488	0.710
ImageNet (Resize)	0.394	0.632	0.471	0.712
Change	+0.014	-0.015	-0.017	+0.002

Table A2. **Analysis of the effect of native image size.** Linear evaluation results for representations pretrained on resized versions of ImageNet and iNat21. ImageNet and iNat21 have images that vary in size, many of which are much larger than the 256x256 images in Places365 and GLC20. Here we analyze the effect of pretraining on resized variants of ImageNet and iNat, which have been preprocessed so that all images have a short side of 256. We use the “Resize” corruption described in Appendix C. Note that the downsampling results in Figure 4 start from resized datasets – in this table we are analyzing the effect of the initial resizing.

can introduce a 1-2% difference in top-1 accuracy, which can be significant on datasets like ImageNet where the performance improvements of new methods are also on the order of 1-2%.

A.6. Is class difficulty preserved between different representations?

To analyze the differences between self-supervised representations a bit further, we ask whether the same classes are “difficult” or “easy” under different representations. In Figure A5 we illustrate how per-class accuracy changes for iNat21 (top row) and Places365 (bottom row) when switching between ImageNet SimCLR and in-domain SimCLR. The panels in the left column define the hardest and easiest examples based on ImageNet SimCLR, while the panels in the right column define the hardest and easiest examples based on in-domain SimCLR. We observe that class difficulty is not preserved between ImageNet SimCLR and iNat21 SimCLR (top row), but it is largely preserved between ImageNet SimCLR and Places365 SimCLR (bottom row). We also note that the overall patterns are the same whether we track the easiest/hardest examples for ImageNet SimCLR and move to the in-domain representation (left column) or track the easiest/hardest examples for in-domain

SimCLR and move to ImageNet SimCLR (right column).

A.7. What is the effect of within-dataset diversity?

In Table 2 we saw that adding pretraining images from a different dataset provides little to no benefit whereas adding pretraining images from the same dataset consistently helps. The surprising conclusion is that a larger, more diverse pretraining dataset can be worse than a smaller, homogeneous pretraining dataset. In this section we present a preliminary study of a milder form of data diversity by changing the number of classes in our pretraining data while holding the number of images constant. We construct three equally sized subsets of ImageNet: one with 200 classes (500 images per class), one with 500 classes (200 images per class), and one with 1k classes (100 images / class). We present linear evaluation results in Table A3. The class information is only used to construct the datasets, which are then used for self-supervised pretraining. Linear classifiers trained on top of these representations use full training sets as usual.

If we assume that class count is a valid proxy for visual diversity, then Table A3 indicates that increasing diversity improves performance on Places365 (+3.9% top-1) but degrades performance on iNat21 (-2.4% top-1). All else being equal, we might intuitively expect a more diverse pretraining set to be beneficial. This seems to be the case for Places365. However, the result for iNat21 shows that this is not necessarily the case. It is possible that more homogeneous pretraining data leads to more fine-grained self-supervised features, which would account for the decrease in performance with increasing diversity for iNat21. Since Places365 is not very fine-grained, it would not benefit from this effect. However, this is a small-scale experiment on one dataset so it should be interpreted with caution.

If these results stand up to under further scrutiny, then we would need to reconcile this finding with our results in Table 2, which show that increased diversity (achieved by replacing some in-domain data with some data from another domain) degrades performance even for Places365. The simplest explanation is that the increased diversity here is much milder - we are simply changing how images are dis-

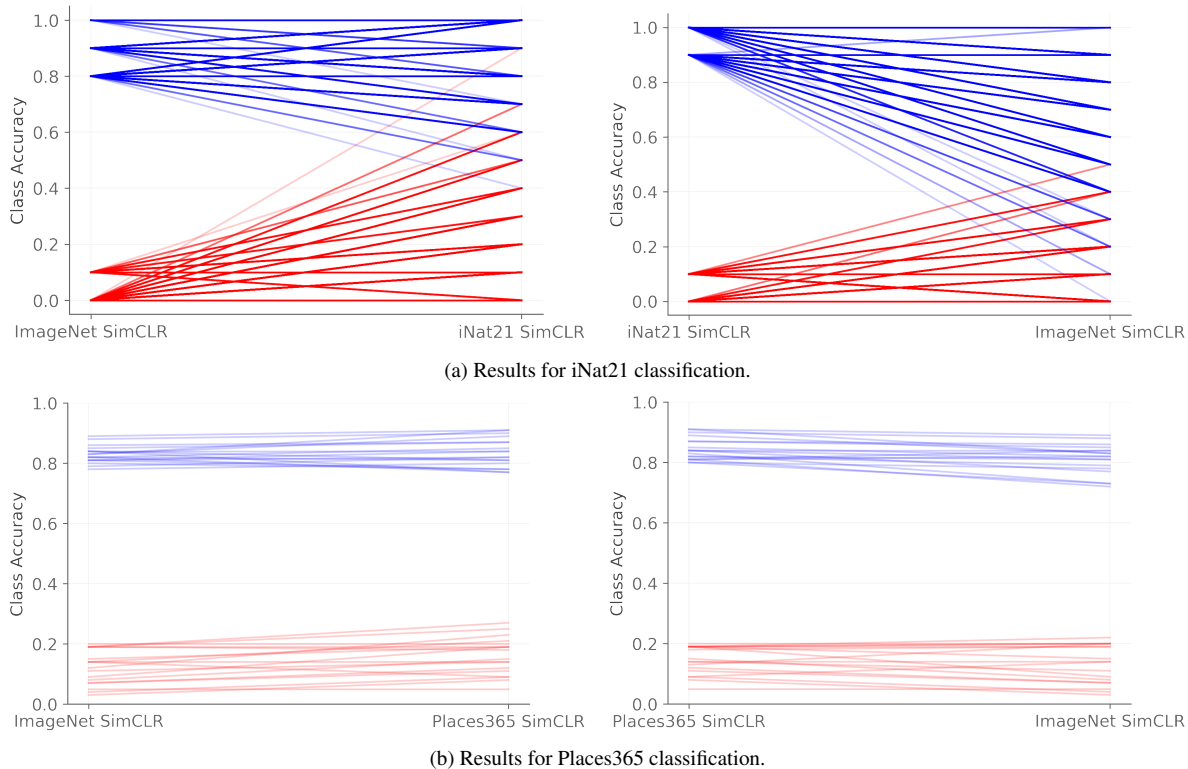


Figure A5. **Difficulty depends on the representation.** Visualization of the change in per-class linear evaluation results when the underlying self-supervised representation is changed. We show the hardest 5% of classes (red lines) and the easiest 5% of classes (blue lines) for the representation named in the bottom left corner of each panel. Left and right plots simply reverse which representation is being used to define the easy and hard classes. Each line represents one class, and shows how the accuracy for that class increases or decreases when we replace the representation named in the bottom-left corner of each panel with the representation named in the bottom-right corner of each panel. Note that the iNat21 validation set has 10 images per class, so all class accuracy values for the top plots lie in $\{0, 0.1, \dots, 0.9, 1.0\}$.

Classes	Images / Class	ImageNet	iNat21	Places365
200	500	0.509	0.314	0.390
500	200	0.531	0.305	0.415
1000	100	0.522	0.290	0.429

Table A3. **What is the effect of image diversity within a dataset?** Linear evaluation results for self-supervised representations based on 100k ImageNet images distributed over different numbers of classes.

tributed over classes, not adding images from other datasets entirely. Our results indicate that this mild diversity is beneficial for pretraining, but too much diversity may render the contrastive pretraining task too easy, resulting in weaker features.

B. Qualitative Examples

Images from different domains. In our paper we consider four datasets: ImageNet, iNat21, Places365, and GLC20. We illustrate their qualitative differences by showing some randomly chosen images from each dataset in Figure A6. By comparing the first row (ImageNet) with the second

(iNat21) and third (Places365) rows, we can see that there are ImageNet images that are semantically similar to images from iNat21 (e.g. the animals in the first and third images) and Places365 (e.g. the bridge scene in the fourth image). The images from GLC20 (bottom row) are quite distinct from the images from the other three datasets.

Corrupted images. In Figure A7 we show examples of the image corruptions we use in Figure 4. While all of these corruptions may seem subjectively mild, Figure 4 shows that they can have a considerable impact on the quality of the learned representations.

C. Implementation Details

C.1. Datasets

iNat21. The 2021 iNaturalist Challenge dataset (iNat21) is a fine-grained species classification dataset [54], with 2.7M training images covering 10k species. Unlike prior iNaturalist datasets [55], iNat21 has an approximately balanced training set. The 100k official validation images are evenly sampled from each species, and we use it as our test set.

GLC20. GeoLifeCLEF 2020 [15] is a collection of remote sensing imagery, designed to facilitate research on location-based species prediction, while also serving as a land cover (LC) classification dataset. Each image is associated with a vector describing the distribution of land cover classes in a 256m^2 patch centered at that location. For the purposes of this work, we binarize this vector (1 for any land cover class whose proportion is nonzero, 0 otherwise) and treat the task as multi-label classification. We only use the half of the dataset from the US, which means we have 1M training images covering 16 land cover classes. Throughout the paper, we refer to this subset of the GeoLifeCLEF 2020 dataset as GLC20. We use the official validation set as a test set, which has around 27k images that were held out in spatial blocks to mitigate the effects of spatial autocorrelation. Note that the labels for this dataset are noisy, so we are mainly interested in GLC20 as a pretraining set.

C.2. Training hyperparameters

SimCLR pretraining. Unless otherwise specified, we use the same settings as the ImageNet experiments in [10]. One exception is that we omit Gaussian blur from the augmentation set since [10] found that it provides a relatively small benefit, around 1% top-1 accuracy on ImageNet. Full details of the augmentations are given in Section C.4. We train with a batch size of 4096 for 1000 epochs and use 16 TPUs for training. We use the LARS optimizer [62] with a learning rate of 4.8 (following $0.075 \times \text{batch size}/256$), decayed on a cosine schedule [36] with a 10-epoch linear warmup and no restarts. For small datasets (size 50k or smaller), we use a lower learning rate of 0.4 (following $0.025 \times \text{batch size}/256$) decayed on a cosine schedule. Our projection head has two layers and an output dimension of 128. A temperature parameter of $\tau = 0.1$ is set for the contrastive loss. Batch normalization statistics are synchronized across TPU cores with a decay parameter of 0.9.

MoCo pretraining. We use the same settings as the ImageNet experiments in [27], with the improvements noted in [13]. As in our SimCLR experiments, we train with a batch size of 1024 using 16 TPUs. For comparability, we use the same augmentation strategy as we do for SimCLR and train for 1000 epochs. Like [10] but unlike [13, 27], we do not standardize images by subtracting per-channel means and dividing by per-channel standard deviations.

BYOL pretraining. We use the same settings as the ImageNet experiments in [26]. As in our SimCLR experiments, we train with a batch size of 4096 using 16 TPUs. For comparability, we use the same augmentation strategy as we do for SimCLR (which happens to be the default for [26]) and train for 1000 epochs. Like [10] but unlike [26], we do not standardize images by subtracting per-channel means and dividing by per-channel standard deviations.

Linear supervised training. Linear classifiers are trained

for 90 epochs using SGD with Nesterov momentum. We use a momentum of 0.9, a batch size of 1024, and a learning rate of 0.4, following the scaling rule $0.1 \times \text{batch size}/256$. The learning rate follows a cosine decay schedule without linear warmup or restarts [36]. Unless otherwise specified, we do not use weight decay / L2 regularization or data augmentation. We take a square center crop with edge length equal to 87.5% of the short side of the image and resize to 224×224 . We use four Tesla V100 GPUs for training.

End-to-end fine-tuning. We use the same settings as linear supervised training with the following exceptions. We train using a smaller batch size of 512 and a lower learning rate of 0.1, following the learning rate scaling rule $0.05 \times \text{batch size}/256$. To mitigate overfitting we use L2 regularization (10^{-4}) for the classifier head and data augmentation (random cropping and horizontal flips). These augmentations use the same implementation as the cropping and flipping used for SimCLR pretraining.

End-to-end supervised training from scratch. We use the same hyperparameters as end-to-end fine-tuning with the following exceptions. We train for 90 epochs using a traditional piece-wise constant learning rate schedule where the initial learning rate of 0.1 is decayed by a factor of 10 at epochs 30 and 60. We also use L2 regularization of 10^{-4} throughout the network.

C.3. Taxonomies

Three of our datasets are equipped with label taxonomies: ImageNet, iNat21, and Places365. We describe these taxonomies below.

ImageNet. We use the WordNet [2] label hierarchy for ImageNet. The finest labels are the standard ImageNet-1k class labels. To coarsen these labels, we start at the deepest level of the hierarchy and merge all leaf nodes at that level with their parents. This produces a new hierarchy, whose leaf nodes will now be used as categories. Each category set is named “Depth k ” where k is the depth of the leaf node that is further from the root. We repeat this process until the leaf nodes merge with the root.

iNat21. Since the categories in iNat21 are animal and plant species, the “tree of life” serves as a natural taxonomy. The taxonomic levels are *Species* (finest, 10k categories), *Genus* (4884 categories), *Family* (1103 categories), *Order* (273 categories), *Class* (51 categories), *Phylum* (13 categories), and *Kingdom* (coarsest, 3 categories). For additional details see [54].

Places365. Places 365 is equipped with a 3-tier hierarchy. The finest labels are the standard category labels for the dataset (“Depth 2”). These categories fall into 16 scene types which constitute the “Depth 1” level of the hierarchy. Examples include water, mountain, transportation, sports, industry, etc. Then the “Depth 0” level consists of a coarser grouping of these scene types into three categories: indoor,

outdoor (natural), and outdoor (man-made).

C.4. Augmentations

In this paper we use three augmentation operations: random horizontal flipping, random cropping, and random color jittering. When training SimCLR, we use all three augmentations. When fine-tuning we only use random horizontal flipping and random cropping as in [10]. We do the same when training from scratch. We do not use any data augmentation when training linear classifiers. For each of these operations we use the implementation from [10] with default settings. We give brief descriptions of each augmentation operation below.

Random horizontal flipping. With probability $1/2$, flip the image over the vertical center line.

Random cropping. Randomly select a rectangular subset of the image covering between 8% and 100% of the whole image, with aspect ratio between $3/4$ and $4/3$.

Random color jitter. Randomly perturb the brightness, contrast, saturation, and hue of the image according to a strength parameter s . See [10] for the exact implementation. We set the strength parameter to $s = 1.0$.

C.5. Corruptions

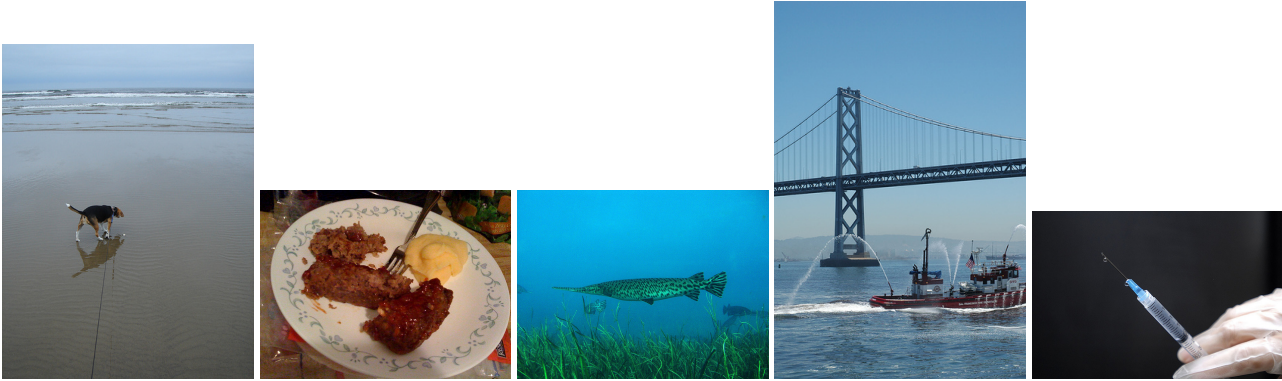
In Section 4.3 of the main paper we investigate the impact of pretraining on artificially degraded images. Here we provide implementation details for each of the image corruption operations.

Resize. We resize the image so that the shorter side is 256 pixels long, but we preserve the aspect ratio. As described below, this corruption allows us to make comparisons which control for image size. Images are resized using the standard PIL [1] function `PIL.Image.resize` with the default nearest-neighbor interpolation.

Resize and downsample. We first apply the “Resize” corruption and then downsample by 2x or 4x before up-sampling by the same factor. The initial resizing is important because some of our datasets have larger images than others and larger images are less affected by down-sampling by a constant factor than their smaller counterparts. Downsampling and up-sampling is accomplished using `PIL.Image.resize` with default settings, just like the “Resize” corruption.

JPEG compression. We use the standard PIL function `PIL.Image.save` to perform JPEG compression. We set the `quality` parameter to 10, which is low enough to cause significant visual artifacts.

Salt and pepper noise. Each pixel in each channel is independently corrupted with probability $1/100$, and corrupted pixels are set to 0 (“pepper”) or 1 (“salt”) with equal probability.



(a) ImageNet



(b) iNat21



(c) Places365



(d) GLC20

Figure A6. **Examples from the datasets used.** We show five randomly selected images from each dataset: ImageNet (top row), iNat21 (second row), Places365 (third row), and GLC20 (bottom row). Note that all images in GLC20 and Places365 are 256×256 pixels, while ImageNet and iNat21 have higher-resolution images and varying aspect ratios. “Places365-Standard” does have varying image resolutions, but we use “Places365-Standard (small images)” which is an official variant that has been resized to 256×256 .



Figure A7. **Examples of corrupted images.** We show the effect of different image corruptions on one randomly chosen image from each dataset: ImageNet (top row), iNat21 (second row), Places365 (third row), and GLC20 (bottom row).