On Robustness and Transferability of Convolutional Neural Networks

Josip Djolonga^{*} Jessica Yung^{*} Michael Tschannen^{*} Rob Romijnders Lucas Beyer Alexander Kolesnikov Joan Puigcerver Matthias Minderer Alexander D'Amour Dan Moldovan Sylvain Gelly Neil Houlsby Xiaohua Zhai Mario Lucic Google Research, Brain Team

Abstract

Modern deep convolutional networks (CNNs) are often criticized for not generalizing under distributional shifts. However, several recent breakthroughs in transfer learning suggest that these networks can cope with severe distribution shifts and successfully adapt to new tasks from a few training examples. In this work we study the interplay between outof-distribution and transfer performance of modern image classification CNNs for the first time and investigate the impact of the pre-training data size, the model scale, and the data preprocessing pipeline. We find that increasing both the training set and model sizes significantly improve the distributional shift robustness. Furthermore, we show that, perhaps surprisingly, simple changes in the preprocessing such as modifying the image resolution can significantly mitigate robustness issues in some cases. Finally, we outline the shortcomings of existing robustness evaluation datasets and introduce a synthetic dataset SI-SCORE we use for a systematic analysis across factors of variation common in visual data such as object size and position.

1. Introduction

Deep convolutional networks have attained impressive results across a plethora of visual classification benchmarks [36, 60] where the training and testing distributions match. In the real world, however, the conditions in which the models are deployed can often differ significantly from the conditions in which the model was trained. It is thus imperative to understand the impact *dataset shifts* [50] have on the performance of these models. This problem has gained a lot of traction and several systematic investigations have shown unexpectedly high sensitivity of image classifiers to various dimensions, including photometric perturbations [27], natural perturbations obtained from video data [54], as well as model-specific adversarial perturbations [23].

The problem of dataset shift, or *out-of-distribution (OOD)* generalization, is closely related to a learning paradigm



Figure 1: We explore the fundamental interplay between in-distribution performance, out-of-distribution (OOD) performance, and transfer learning performance (red arrows in the graph on the right), with respect to the major design choices listed on the left. The relationship between in-distribution and OOD performance is highly under-explored along these axes, whereas that between OOD and transfer performance has not been studied before to the best of our knowledge.

known as *transfer learning* [56, §13]. In transfer learning we are interested in constructing models that can improve their performance on some target task by leveraging data from different related problems. In contrast, under dataset shift one assumes that there are two environments, namely training and testing [56], with the constraint that the model cannot be adapted using data from the target environment. As a consequence, the two environments typically have to be more similar and their differences more structured than in the transfer setting (c.f. Section 2).

In the context of transfer learning, detailed scaling laws characterizing the interplay between the in-distribution and transfer performance as a function of pre-training data set size, model size, architectural choices such as normalization, and transfer strategy have been established recently [37, 72, 36]. Model and dataset scale were identified as key factors for transfer performance. The similarities between transfer learning and OOD generalization suggests that these axes are also relevant for OOD generalization and raises the question of what the corresponding scaling laws are. While some axes have been partially explored by prior work [27, 70], the big picture is largely unknown. Even more importantly, is in-distribution performance enough to characterize OOD performance, or can transfer performance give a more fine-grained characterization of OOD performance

^{*}Shared first authorship. Please send e-mail correspondence to {josipd,lucic}@google.com.

of a population of models than in-distribution performance? To the best of our knowledge, this question has not been systematically explored before in the literature.

Contributions We systematically investigate the interplay between the in-distribution accuracy of image classification models on the training distribution, their generalization to OOD data (without adaptation), and their transfer learning performance with adaptation in the low-data regime (see Fig. 1 for an illustration). Specifically:

- (i) We present the first meta-analysis of existing OOD metrics and transfer learning benchmarks across a wide variety of models, ranging from self-supervised to fully supervised models with up to 900M parameters. We show that increasing the model and data scale disproportionately improves transfer and OOD performance, while only marginally improving the performance on the IMAGENET validation set.
- (ii) Focusing on OOD robustness, we analyze the effects of the training set size, model scale, and the training regime and testing resolution, and find that the effect of scale overshadows all other dimensions.
- (iii) We introduce a novel dataset for fine-grained OOD analysis to quantify the robustness to object size, object location, and object orientation (rotation angle). We believe that this is a first systematic study to show that the models become less sensitive (and hence more robust) to each of these factors of variation as the dataset size and model size increase.

2. Background

Robustness of image classification models Understanding and correcting for dataset shifts are classical problems in statistics and machine learning, and have as such received substantial attention, see e.g. the monograph [50]. Formally, let us denote the observed variable by X and the variable we want to predict by Y. A dataset shift occurs when we train on samples from $P_{\text{train}}(X, Y)$, but are at test time evaluated under a different distribution $P_{\text{test}}(X,Y)$. Storkey [56] discusses and precisely defines different possibilities for how P_{train} and P_{test} can differ. We are mostly interested in covariate shifts, i.e., when the conditionals $P_{\text{train}}(Y|X) = P_{\text{test}}(Y|X)$ agree, but the marginals $P_{\text{train}}(X)$ and $P_{\text{test}}(X)$ differ. Most robustness datasets proposed in the literature targeting IMAGENET models are such instances—the images X come from a source $P_{\text{test}}(X)$ different from the original collection process $P_{\text{train}}(X)$, but the label semantics do not change. As a robustness score one typically uses the expected accuracy, i.e., $P_{\text{test}}(Y = f(X))$, where f(X) is the class predicted by the model.

Dataset shift types IMAGENET-V2 is a recollected version of the IMAGENET validation set [52]. The authors attempted to replicate the data collection process, but found

that all models drop significantly in accuracy. Recent work attributes this drop to statistical bias in the data collection [17]. IMAGENET-C and IMAGENET-P [27] are obtained by corrupting the IMAGENET validation set with classical corruptions, such as blur, different types of noise and compression, and further cropping the images to 224×224 . These datasets define a total of 15 noise, blur, weather, and digital corruption types, each appearing at 5 severity levels or intensities. OBJECTNET [3] presents a new test set of images collected directly using crowd-sourcing. OBJECTNET is particular as the objects are captured at unusual poses in cluttered, natural scenes, which can severely degrade recognition performance. Given this clutter, and arguably better suitability as a detection than recognition task [5], Y|Xmight be hard to define and the dataset goes beyond a covariate shift. In contrast, the IMAGENET-A dataset [30] consists of real-world, unmodified, and naturally occurring examples that are misclassified by ResNet models. Hence in addition to the covariate shift due to the data source, this dataset is not model-agnostic and exhibits a strong selection bias [56].

Attempting to focus on naturally occurring invariances, [54] annotated two video datasets: IMAGENET-VID-ROBUST and YOUTUBE-BB-ROBUST, derived from the IMAGENET-VID [11] and YOUTUBE-BB [51] datasets respectively. In [54] the authors propose the pm-k metric given an anchor frame and up to k neighboring frames, a prediction is marked as correct only if the classifier correctly classifies all 2k + 1 frames around and including the anchor. We present the details of each dataset in Appendix A.

Transferability of image classification models In transfer learning [48], a model might leverage the data it has seen on a related distribution, $P_{\rm pre-train}$, to perform better on a new task $P_{\rm train}$. Note that in contrast to the covariate shift setting, the disparity between $P_{\rm pre-train}$ and the new task is typically larger, but one is further given samples from $P_{\rm train}$. While there exist many approaches on how to transfer to the new task, the most common approach in modern deep learning, which we use, is to (i) train a model on $P_{\rm pre-train}$ (using perhaps an auxiliary, self-supervised task [15, 22]), and then (ii) train a model on $P_{\rm train}$ by initializing the model weights from the model trained in the first step.

Recently, a suite of datasets has been collected to benchmark modern image classification transfer techniques [72]. The Visual Task Adaptation Benchmark (VTAB) defines 19 datasets with 1000 labeled samples each, categorized into three groups: *natural* (most similar to IMAGENET) consists of standard natural classification tasks (e.g., CIFAR); *specialized* contains medical and satellite images; and *structured* (least similar to IMAGENET) consists mostly of synthetic tasks that require understanding of the geometric layout of scenes. We compute an overall transfer score as the mean across all 19 datasets, as well as scores for each subgroup of tasks. We provide details for all of the tasks in Appendix A.



Figure 2: The relationship between transfer learning, IMAGENET, and robustness performance. (Left) Average score on all transfer benchmarks versus IMAGENET performance. (Center) Average score on all robustness benchmarks versus average transfer performance. (Right) Correlation between different groups of transfer datasets (natural, specialized, structured), and robustness metrics.

3. A meta-analysis of robustness and transferability metrics

While many robustness metrics have been proposed to capture different sources of brittleness, it is not well understood how these metrics relate to each other. We investigate the practical question of how useful the various metrics are in guiding design choices. Further, we empirically analyze the relationship between robustness and transferability metrics, which is lacking in the literature, despite their close relationship. To analyze these questions, we evaluated 39 different models over 23 robustness metrics and the 19 transfer tasks.

Metrics For robustness, we measure the model accuracy on the IMAGENET, IMAGENET-V2 (the *matched frequency* variant) and OBJECTNET datasets. We also consider video datasets, IMAGENET-VID and YOUTUBE-BB; we use both the accuracy metric and the pm-10 metric (suffix -W). On IMAGENET-C we report the AlexNet-accuracy-weighted [39] accuracy over all corruption times (called *mean corruption error* in [27]). To evaluate the transferability of the models, we use the VTAB-1K benchmark introduced in Section 2. We evaluate average transfer performance across all 19 datasets, with 1000 examples each, as well as per-group performance. To transfer a model we performed a sweep over two learning rates and schedules. We report the median testing accuracy over three fine-tuning runs with parameters selected using a 800-200 example train-validation split.

Models To perform this meta-analysis we consider several model families. We evaluate ResNet-50 [24] and six EfficientNet (B0 through B5) models [60] including variants using AutoAugment [10] and AdvProp [69], which have been trained on IMAGENET. We include self-supervised SimCLR [6] (variants: linear classifier on fixed representation (lin), fine-tuned on 10% (ft-10), and 100% (ft-100) of the IMAGENET data), and self-supervised-semi-supervised (S4L) [71] models that have been fine-tuned to 10% and 100% of the IMAGENET data. We also consider a set of models that use other data sources. Specifically, three NoisyStudent [70] variants which use IMAGENET and unlabelled data from the JFT dataset, BiT (BigTransfer) [36] models that have been first trained on IMAGENET, IMAGENET-21K,

or JFT and then transferred to IMAGENET by fine-tuning, and the Video-Induced Visual Invariance (VIVI) model [66], which uses IMAGENET and unlabelled videos from the YT8M dataset [1]. Finally, we consider the BigBiGAN [14] model which has been first trained as a class-conditional generative model and then fine-tuned as an IMAGENET classifier. All details can be found in Appendix E.

How informative are robustness metrics for discriminating between models? The goal of a metric is to discriminate between different models and thus guide design choices. We therefore quantify the usefulness of each metric in terms of how much it improves the discriminability between the various models beyond the information provided by IMA-GENET accuracy. Specifically, we train logistic regression classifiers to discriminate between the 12 model groups outlined above. We compared the performance of a classifier using only IMAGENET accuracy as input feature, to a classifier using IMAGENET and up to two of the other metrics, see Fig. 4 and Appendix A. We found that most of the tested metrics provide little increase in model discriminability over IMAGENET accuracy. We further, similarly to [61], found that all metrics are highly rank-correlated with each other, which we present in Appendix A. Of course, these results are conditioned on the size and composition of our dataset, and may differ for a different set of models. However, based on our collection of 39 models in 12 groups, the most informative metrics are those based on different datasets and/or video, rather than IMAGENET-derived datasets.

How related are OOD robustness and transfer metrics? Next, we turn to transfer learning. It has been observed that better IMAGENET models transfer better [37, 72]. Since robustness metrics correlated strongly with IMAGENET accuracy, we might expect a similar relationship. To get an overall view, we compute the mean of all robustness metrics, and compare it to transfer performance. Figure 2 (center) shows this average robustness plotted against transfer performance, while Figure 2 (left) shows transfer versus IMAGENET accuracy. Indeed, we observe a large correlation coefficient $\rho = 0.73$ between robustness and transfer metrics; however, the correlation is not stronger than



Figure 3: (**Top**) Reduction (in %) in classification error relative to the classification error of the model trained for 112k steps on 1M examples (bottom left corner) as a function of training iterations and training set size. The results are for a ResNet-101x3 trained on IMAGENET-21K subsets. (**Bottom**) Relative reduction (in %) in classification error going from a ResNet-50 to a ResNet-101x3 as a function of training steps and training set size (IMAGENET-21K subsets). The reduction generally increases with the training set size and longer training. Hence, the right *scaling laws* not only lead to in-distribution improvements, but also to simultaneous improvements across a heterogeneous set of OOD benchmarks. We investigate why these larger models achieve stronger performance across all benchmarks in Section **5**.

between transfer and IMAGENET. Further, we compute the correlation of the residual robustness score (mean robustness minus IMAGENET accuracy) against transfer score, and find only a weak relationship of $\rho = 0.12$. This indicates that robustness metrics, on aggregate, do not provide additional signal that predicts model transferability beyond that of the base IMAGENET performance. We do, however, see some interesting differences in the relative performances of different model groups. Certain model groups, while attaining reasonable IMAGENET/robustness scores, transfer less well to VTAB. Therefore, *there are factors unrelated to robust inference that do influence transferability*. One example is batch normalization which is outperformed by group normalization with weight standardization in transfer [36]. Next, we break down the correlation by



Figure 4: Informativeness of robustness metrics. Values indicate the difference in accuracy of a logistic classifier trained to discriminate between model types based on IMAGENET accuracy plus one additional metric, compared to a classifier trained only on IMAGENET accuracy (higher is better, top 10 metrics shown). Bars show mean \pm s.d. of 1000 bootstrap samples from the 39 models.

robustness metrics and transfer datasets in Fig. 2 (right). We see that each metric correlates similarly with the task groups. However, for the groups that require more distant transfer (Specialized, Structured), no metric predicts transferability well. Perhaps surprisingly, raw IMAGENET accuracy is the best predictor of transfer to *structured* tasks, indicating that robustness metrics do not relate to challenging transfer tasks, at least not more than raw IMAGENET accuracy.

Summary Metrics based on ImageNet have very little additional discriminative power over ImageNet accuracy, while those not based on ImageNet have more, but their additional discriminative power is still low—popular robustness metrics provide marginal complementary information. Transferability is also related to IMAGENET accuracy, and hence robustness. We observe that while there is correlation, transfer highlights failures that are somewhat independent of robustness. Further, no particular robustness metric appears to correlate better with any particular group of transfer tasks than IMAGENET does. Inspired by these results, we next investigate strategies known to be effective for IMAGENET and transfer learning on the OOD robustness benchmarks.

4. Scaling laws for OOD performance

Increasing the scale of pre-training data, model architecture, and training steps have recently led to diminishing improvements in terms of IMAGENET accuracy. By contrast, it has been recently established that scaling along these axes can lead to substantial improvements in transfer learning performance [36, 60]. In the context of robustness, this type of scaling has been explored less. While there are some results hinting that scale can improve robustness [27, 52, 70, 64], no principled study decoupling the different scale axes has been performed. Given the strong correlation between transfer performance and robustness, this motivates the systematic investigation of the effects of the pre-training data size, model architecture size, training steps, and input resolution. While paramount to the out-of-distribution performance, as we find, these pretraining design choices have not yet received a great deal of attention from the community.

Setup We consider the standard IMAGENET training setup [24] as a baseline, and scale up the training accordingly. To study the impact of dataset size, we consider the IMAGENET-21K [11] and JFT [57] datasets for the experiments, as pre-training on either of them has shown great performance in transfer learning [36]. We scale from the IM-AGENET training set size (1.28M images) to the IMAGENET-21K training set size (13M images, about 10 times larger than IMAGENET). To explore the effect of the model size, we use a ResNet-50 as well as the deeper and $3 \times$ wider ResNet-101x3 model. We further investigate the impact of the training schedule as larger datasets are known to benefit from longer training for transfer learning [36]. To disentangle the impact of dataset size and training schedules, we train the models for every pair of dataset size and schedule.

We fine-tune the trained models to IMAGENET using the BiT HyperRule [36], and assess their OOD generalization performance in the next section. Throughout, we report the reduction in classification error relative to the model which was trained on the smallest number of examples and for the fewest iterations, and which hence achieves the lowest accuracy. Other details are presented in Appendix B.

Pre-training dataset size impact The results for the ResNet-101x3 model are presented in Fig. 3. When pretrained on IMAGENET-21K, the OOD classification error significantly decreases with increasing pre-training dataset size and duration: We observe relative error reductions of 20-30% when going from 112k steps on 1M data points to 1.12M steps on 13M data points. The reductions are least pronounced for YOUTUBE-BB(-W). Note that training for 1.12M steps leads to a lower accuracy than training for only 457k steps unless the full IMAGENET-21K dataset is used. For models trained on JFT we observe a similar behavior except that training for 1.12M steps often leads to a higher accuracy than training for 457k steps even when only 1M or 5M data points are used (c.f. Appendix B). These results suggest that, if the models have enough capacity, increasing the amount of pre-training data, without any additional changes, leads to substantial gains in all datasets simultaneously which is in line with recent results in transfer learning [36].

Model size impact Figure 3 shows the relative reduction in classification error when using a ResNet-101x3 instead of a ResNet-50 as a function of the number of training steps and the dataset size. It can be seen that increasing the model



Figure 5: Comparison of different types of evaluation preprocessing and resolutions. (Default, blue): Accuracy obtained for the preprocessing and resolution proposed by the authors of the respective models. (Best, orange): The accuracy when selecting the best resolution from {64, 128, 224, 288, 320, 384, 512, 768}. (FixRes, green): Applying FixRes for the same set of resolutions and selecting the best resolution. Increasing the evaluation resolution and additionally using FixRes helps across a large range of models and pretraining datasets on IMAGENET-A and OBJECTNET.

size can lead to substantial reductions of 5–20%. For a fixed training duration, using more data always helps. However, on IMAGENET-21K, training too long can lead to increases in the classification error when the model size is increased, unless the full IMAGENET-21K is used. This is likely due to overfitting. This effect is much less pronounced when JFT is used for training. JFT results are presented in Appendix B. Again, reductions in classification error are least pronounced for YOUTUBE-BB/YOUTUBE-BB-W.

Testing resolution and OOD robustness During training, images are typically cropped randomly, with many crop sizes and aspect ratios, to prevent overfitting. In contrast, during testing, the images are usually rescaled such that the shorter side has a pre-specified length, and a fixed-size center crop is taken and then fed to the classifier. This leads to a mismatch in object sizes between training and testing. Increasing the resolution at which images are tested leads to an improvement in accuracy across different architectures [63, 64]. Furthermore, additional benefits can be obtained by applying FixRes — fine-tuning the network on the training set with the test-time preprocessing (i.e. omitting random cropping with aspect ratio changes), and at a higher resolution. We explore the effect of this discrepancy on the robustness of different architectures. As some of the robustness datasets were collected differently from IMAGENET, discrepancies in the cropping are likely. We investigate both adjusting testtime resolution and applying FixRes. For FixRes, we use a simple setup with a single schedule and learning rate for all models (except using a $10 \times$ smaller learning rate for the BiT models), and without heavy color augmentation as in [63] or label smoothing as in [64]. We did not extensively tune hyperparameters, but chose a setup that works reasonably well across architectures and training datasets. Note that changing the resolution can be seen as scaling the computational resources available to the model, as both training and inference costs grow with the resolution.

Following the protocol of the FixRes paper [63], we evaluate each model for all resolutions in {64, 128, 224, 288, 320, 384, 512, 768} to *illustrate the potential of adapting the testing resolution* (in practice we do not have access to an OOD validation set so we cannot select the optimal solution in advance). For conciseness, we show the accuracy for IMAGENET-A and OBJECTNET at the testing resolution proposed by the authors of the respective architecture along with the highest accuracy across testing resolutions (Figure 5). The results for other datasets and resolutions are deferred to Appendix C.

We start by discussing observations that apply to most models, excluding the BiT models which will be discussed below. While FixRes only leads to marginal benefits on IMAGENET, it can lead to substantial improvements on the robustness metrics. Choosing the optimal testing resolution leads to a significant increase in accuracy on IMAGENET-A and OBJECTNET in most cases, and applying FixRes often leads to additional substantial gains. For OBJECTNET, fine-tuning with testing preprocessing (i.e. fine-tuning with central cropping instead of random cropping as used during training) can help even without increasing resolution.

Increasing the resolution and/or applying FixRes often slightly helps on IMAGENET-V2. For IMAGENET-C, the optimal testing resolution often corresponds to the resolution used for training, and applying FixRes rarely changes this picture. This is not surprising as the IMAGENET-C images are cropped to 224 pixels by default, and increasing the resolution does not add any new information to the image. For the video-derived robustness datasets IMAGENET-VID-ROBUST and YOUTUBE-BB-ROBUST, evaluating at a larger testing resolution and/or applying FixRes at a higher resolution can substantially improve the accuracy on the anchor frame and the robustness accuracy for small EfficientNet and ResNet models, but does not help the larger ones. For the BiT models, the resolution suggested by the authors is almost always optimal, except on OBJECTNET and IMAGENET-A, where changing the preprocessing considerably helps. FixRes arguably does not lead to improvements as it was already applied in BiT as a part of the BiT HyperRule.

Summary These empirical results point to the following conclusion: for models with enough capacity, increasing the amount of pre-training data, with no additional changes, leads to substantial gains in all considered OOD generalization tasks *simultaneously*. Secondly, resolution adjustments as outlined above can address the considerable distribution shift caused by resolution mismatch.

F.O.V.	DATASET CONFIGURATION	IMAGES
Size	Objects upright in the center, sizes from 1% to	92884
	100% of the image area in 1% increments.	
LOCATION	Objects upright. Sizes are 20% of the image	479184
	area. We do a grid search of locations, dividing	
	the x-coordinate and y-coordinate dimensions	
	into 20 equal parts each, for a total of 441 coor-	
	dinate locations.	
ROTATION	Objects in the center, sizes equal to 20%, 50%,	39540
	80% or 100% of the image size. Rotation an-	
	gles ranging from 1 to 341 degrees counter-	
	clockwise in 20-degree increments.	

Table 1: Synthetic dataset details. The first column shows the relevant factor of variation (F.O.V.). When there are multiple values for multiple factors of variation, we generate the full cross product of images.

5. SI-SCORE: A fine-grained analysis of robustness to common factors of variation

The results in Section 4 do not reveal the underlying reasons for the success of larger models trained on more data on all robustness metrics. Intuitively, one would expect that these models are more invariant to specific factors of variation, such as object location, size, and rotation. However, a systematic assessment hinges on testing data which can be varied according to these axes in a controlled way. At the same time, the combinatorial nature of the problem precludes any large-scale systematic data collection scheme.

In this work we present a scalable alternative and construct a novel synthetic dataset for fine-grained evaluation: SI-SCORE (Synthetic Interventions on Scenes for Robustness Evaluation). In a nutshell, we paste a large collection of objects onto uncluttered backgrounds (Figure 6, Figure 14a), and can thus conduct controlled studies by systematically varying the object class, size, location, and orientation.¹

Synthetic dataset details The foregrounds were extracted from OpenImages [40] using the provided segmentation masks. We include only object classes that map to ImageNet classes. We also removed all objects that are tagged as occluded or truncated, and manually removed highly incomplete or inaccurately labeled objects. The backgrounds were images from nature taken from *pexels.com* (the license therein allows one to reuse photos with modifications). We manually filtered the backgrounds to remove ones with prominent objects, such as images focused on a single animal or person. In total, we converged to 614 object instances across 62 classes, and a set of 867 backgrounds.

We constructed three subsets for evaluation, one corresponding to each factor of variation we wanted to investigate, as shown in Table 1. In particular, for each object instance,

¹The synthetic dataset and code used to generate the dataset are opensourced on GitHub and are being hosted by the Common Visual Data Foundation.



Figure 6: (Left) Sample images from our synthetic dataset. We consider 614 foreground objects from 62 classes and 867 backgrounds and vary the object location, rotation angle, and object size for a total of 611 608 images. (**Right**) In the first column, for each location on the grid, we compute the average accuracy. Then, we normalize each location by the 95th percentile across all locations, which quantifies the gap between the locations where the model performs well, and the ones where it under-performs (first column, dark blue versus white). Then, we consider models trained with more data, compute the same normalized score, and plot the *difference* with respect to the first column. We observe that, as dataset size increases, sensitivity to object location decreases – the outer regions improve in relative accuracy more than the inner ones (e.g. dark blue vs white on the second and third columns). The effect is more pronounced for the larger model. The full set of results is presented in Figure 17 in Appendix D.

we sample two backgrounds, and for each of these objectbackground combinations, we take a cross product over all the factors of variation. For the datasets with multiple values for more than one factor of variation, we take a cross product of all the values for each factor of variation in the set (object size, rotation, location). For example, for the rotation angle dataset, there are four object sizes and 18 rotation angles, so we do a cross product and have 72 factor of variation combinations. For the object size and rotation datasets, we only consider images where objects are at least 95% in the image. For the location dataset, such filtering removes almost all images where objects are near the edges of the image, so we do not do such filtering. Note that since we use the central coordinates of objects as their location, at least 25% of each object is in the image even if we do not do any filtering. The results in the following sections are similar when filtering out objects that are less than 50% or 75% in the image.

Learned invariances as a function of scale We study one factor of variation at a time. For example, when studying the impact of changing the location of the object center, we measure the average performance for each location over a uniform grid. Building on our investigation in the previous section, we test whether increasing model size and dataset size improves robustness to these three factors of variation by evaluating the ResNet-50 and ResNet-101x3 models. We observe that the models indeed become more invariant to object location (Figure 6), rotation (Figure 7, left), and size (Figure 7, right) as the pre-training set size increases. Specifically, as we pre-train on more data, the average prediction accuracy across various object locations, sizes, and rotation angles becomes more uniform. Furthermore, the larger ResNet-101x3 model is indeed more robust. Analogous results on the JFT dataset are presented in Appendix D.

6. Related work

There has been a growing literature exploring the robustness of image classification networks. Early investigations in face and natural image recognition found that performance degrades by introducing blur, Gaussian noise, occlusion, and compression artifacts, but less by color distortions [12, 35]. Subsequent studies have investigated brittleness to similar corruptions [53, 76], as well as to impulse noise [31], photometric perturbations [62], and small shifts and other transformations [2, 17, 74]. CNNs have also been shown to over-rely upon texture rather than shape to make predictions, in contrast to human behavior [20]. Robustness to adversarial attacks [23] is a related, but distinct problem, where performance under worst-case perturbations are studied. In this paper we did not study such adversarial robustness, but have focused on average-case robustness to natural perturbations.

2Several techniques have been shown to improve model robustness on these datasets. Using better data augmentation can improve performance on data with synthetic noise [29, 43]. Auxiliary self-supervision [7, 71] can improve robustness to label noise and common corruptions [28]. Transductive fine-tuning using self-supervision on the test data improves performance under distribution shift [58]. Training with adversarial perturbations improves many robustness benchmarks if one uses separate Batch-Norm parameters for clean and adversarial data [69]. Finally, additional pre-training using very large auxiliary datasets has recently shown significant improvements in robustness. Noisy Student [70] reports good performance on several robustness datasets, while Big Transfer (BiT) [36] reports strong performance on the OBJECTNET dataset [3].

Deep networks are often trained by pre-training the network on a different problem and then fine-tuning on the



Figure 7: (Left) In the first row of both plots we show the ratio of the accuracy and the best accuracy (across all rotations). For the second row (model trained on 2.6M instances) and other rows, we compute the same normalized score and visualize the difference with the first row. Larger positive differences with the first row imply a more uniform behavior across object rotations. We observe that, as the dataset size increases, the average prediction accuracy across various rotation angles becomes more uniform. The effect is more pronounced for the larger model. (Right) Similarly, the average accuracy across various object sizes becomes more uniform for both models. As expected, the improvement is most pronounced for small object sizes covering 10–20% of the pixels. The full set of results is presented in Figures 15 and 16 in Appendix D.

target task. This pre-training is often referred to as representation learning; representations can be trained using supervised [32, 36], weakly-supervised [44], or unsupervised data [13, 14, 66, 70]. Recent benchmarks have been proposed to evaluate transfer to several datasets, to assess generalization to tasks with different characteristics, or those disjoint from the pre-training data [65, 72]. While state-ofthe-art performance on many competitive datasets is attained via transfer learning [70, 36], the implications for final robustness metrics remain unclear.

Creating synthetic datasets by inserting objects onto backgrounds has been used for training [75, 16, 21] and evaluating models [36], but previous works do not systematically vary object size, location or orientation, or analyze translation and rotation robustness only at the image level [18].

Given the lack of a consensus on what "natural" perturbations *are*, there are no established general laws on how models behave under various data shifts. Concurrently, [61] investigated whether higher accuracy on synthetic datasets translates to superior performance on natural OOD datasets. They also identify model size and training data set size as the only technique providing a benefit. In [26] the authors list several of the hypotheses that appear in the literature, and collect new datasets that provide (both positive and negative) evidence for their soundness.

7. Limitations and future work

We analyzed OOD generalization and transferability of image classifiers, and demonstrated that model and data scale together with a simple training recipe lead to large improvements. However, these models do exhibit substantial performance gaps when tested on OOD data, and further research is required. Secondly, this approach hinges on the availability of curated datasets and significant computing capabilities which is not always practical. Hence, we believe that transfer learning, i.e. train once, apply many times, is the most promising paradigm for OOD robustness in the short term. One limitation of this study is that we consider image classification models fine-tuned to the IMAGENET label space which were developed with the goal of optimizing the accuracy on the IMAGENET test set. While existing work shows that we do not overfit to IMAGENET, it is possible that these models have correlated failure modes on datasets which share the biases with IMAGENET [52]. This highlights the need for datasets which enable fine-grained analysis for all important factors of variation and we hope that our dataset will be useful for researchers.

The introduced synthetic data can be used to investigate other qualitative differences between models. For example, when comparing ResNet-50s trained on ImageNet, a ResNet using GroupNorm does better on smaller objects than one with BatchNorm, whereas the model with BatchNorm does better on larger objects (Figure 14b in the appendix). While a thorough investigation is beyond the scope of this work, we hope that SI-SCORE will be useful for such future studies.

Instead of requiring the model to work under various dataset shifts, one can ask an alternative question: assuming that the model will be deployed in an environment significantly different from the training one, can we at least quantify the model uncertainty for each prediction? This important property remains elusive for moderate-scale neural networks [55], but could potentially be improved by large-scale pretraining which we leave for future work.

References

- Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. arXiv:1609.08675, 2016. 3
- [2] Aharon Azulay and Yair Weiss. Why do deep convolutional networks generalize so poorly to small image transformations? *Journal of Machine Learning Research*, 20, 2019.
 7
- [3] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. In Advances in Neural Information Processing Systems, 2019. 2, 7, 13
- [4] Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler, Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, et al. Deepmind lab. arXiv preprint arXiv:1612.03801, 2016. 13
- [5] Ali Borji. Objectnet dataset: Reanalysis and correction. In arXiv 2004.02042, 2020. 2
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. arXiv:2002.05709, 2020. 3, 22
- [7] Ting Chen, Xiaohua Zhai, Marvin Ritter, Mario Lucic, and Neil Houlsby. Self-supervised GANs via auxiliary rotation loss. In *Conference on Computer Vision and Pattern Recognition*, 2019. 7
- [8] Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 2017. 13
- [9] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014. 13
- [10] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *Conference on Computer Vision and Pattern Recognition*, 2019. 3
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, 2009. 2, 5
- [12] Samuel Dodge and Lina Karam. Understanding how image quality affects deep neural networks. In *International Conference on Quality of Multimedia Experience*, 2016. 7
- [13] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *International Conference on Computer Vision*, 2015. 8
- [14] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. In Advances in Neural Information Processing Systems, 2019. 3, 8, 22
- [15] Alexey Dosovitskiy, Philipp Fischer, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with exemplar convolutional neural networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(9), 2015. 2

- [16] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, paste and learn: Surprisingly easy synthesis for instance detection. In *International Conference on Computer Vision*, 2017. 8
- [17] Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Jacob Steinhardt, and Aleksander Madry. Identifying statistical bias in dataset replication. arXiv: 2005.09619, 2020. 2, 7
- [18] Logan Engstrom, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. A rotation and a translation suffice: Fooling CNNs with simple transformations. *CoRR*, abs/1712.02779, 2017. 8
- [19] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research*, 2013. 13
- [20] 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 *International Conference on Learning Representations*, 2019. 7
- [21] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D. Cubuk, Quoc V. Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation, 2020. 8
- [22] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. arXiv:1803.07728, 2018. 2
- [23] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. arXiv:1412.6572, 2014. 1, 7
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference* on Computer Vision and Pattern Recognition, 2016. 3, 5
- [25] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2019. 13
- [26] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. arXiv preprint arXiv:2006.16241, 2020. 8
- [27] Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common corruptions and surface variations. arXiv: 1807.01697, 2018. 1, 2, 3, 4, 13
- [28] Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using self-supervised learning can improve model robustness and uncertainty. In Advances in Neural Information Processing Systems, 2019. 7
- [29] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. arXiv:1912.02781, 2019. 7
- [30] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. *arXiv*: 1907.07174, 2019. 2, 13

- [31] Hossein Hosseini, Baicen Xiao, and Radha Poovendran. Google's cloud vision api is not robust to noise. In *International Conference on Machine Learning and Applications*, 2017. 7
- [32] Minyoung Huh, Pulkit Agrawal, and Alexei A Efros. What makes imagenet good for transfer learning? arXiv:1608.08614, 2016. 8
- [33] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *IEEE Conference on Computer Vision* and Pattern Recognition, 2017. 13
- [34] Kaggle and EyePacs. Kaggle diabetic retinopathy detection, July 2015. 13
- [35] Samil Karahan, Merve Kilinc Yildirim, Kadir Kirtaç, Ferhat Sükrü Rende, Gultekin Butun, and Hazim Kemal Ekenel. How image degradations affect deep CNN-based face recognition? In *International Conference of the Biometrics Special Interest Group*, 2016. 7
- [36] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (BiT): General visual representation learning. *European Conference on Computer Vision*, 2020. 1, 3, 4, 5, 7, 8, 13, 14, 22
- [37] Simon Kornblith, Jonathon Shlens, and Quoc V. Le. Do better imagenet models transfer better? In *Conference on Computer Vision and Pattern Recognition*, 2019. 1, 3
- [38] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009. 13
- [39] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In Advances in Neural Information Processing Systems, 2012.
 3, 13
- [40] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. arXiv: 1811.00982, 2020. 6
- [41] Yann LeCun, Fu Jie Huang, and Leon Bottou. Learning methods for generic object recognition with invariance to pose and lighting. In *IEEE Conference on Computer Vision* and Pattern Recognition, 2004. 13
- [42] Fei-Fei Li, Rob Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE Transactions on Pattern Analysis* and Machine Intelligence, 2006. 12
- [43] Raphael Gontijo Lopes, Dong Yin, Ben Poole, Justin Gilmer, and Ekin D Cubuk. Improving robustness without sacrificing accuracy with patch gaussian augmentation. arXiv:1906.02611, 2019. 7
- [44] Dhruv Mahajan, Ross B. Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten. Exploring the limits of weakly supervised pretraining. In *European Conference on Computer Vision*, 2018. 8
- [45] Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentanglement testing sprites dataset. https://github.com/deepmind/dsprites-dataset/, 2017. 13

- [46] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011. 13
- [47] M-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In *Indian Conference* on Computer Vision, Graphics and Image Processing, Dec 2008. 13
- [48] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10), 2009. 2
- [49] O. M. Parkhi, A. Vedaldi, A. Zisserman, and C. V. Jawahar. Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 13
- [50] Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset shift in machine learning*. The MIT Press, 2009. 1, 2
- [51] Esteban Real, Jonathon Shlens, Stefano Mazzocchi, Xin Pan, and Vincent Vanhoucke. Youtube-boundingboxes: A large high-precision human-annotated data set for object detection in video. In *Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [52] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? arXiv: 1902.10811, 2019. 2, 4, 8, 13
- [53] Prasun Roy, Subhankar Ghosh, Saumik Bhattacharya, and Umapada Pal. Effects of degradations on deep neural network architectures. arXiv:1807.10108, 2018. 7
- [54] Vaishaal Shankar, Achal Dave, Rebecca Roelofs, Deva Ramanan, Benjamin Recht, and Ludwig Schmidt. A systematic framework for natural perturbations from videos. *arXiv*:1906.02168, 2019. 1, 2, 13
- [55] Jasper Snoek, Yaniv Ovadia, Emily Fertig, Balaji Lakshminarayanan, Sebastian Nowozin, D. Sculley, Joshua V. Dillon, Jie Ren, and Zachary Nado. Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift. In Advances in Neural Information Processing Systems, 2019. 8
- [56] Amos Storkey. When training and test sets are different: characterizing learning transfer. *Dataset shift in machine learning*, 2009. 1, 2
- [57] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *International Conference on Computer Vision*, 2017. 5, 14
- [58] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei A Efros, and Moritz Hardt. Test-time training for out-ofdistribution generalization. arXiv:1909.13231, 2019. 7
- [59] C. Szegedy, Wei Liu, Yangqing Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Conference on Computer Vision and Pattern Recognition*, 2015. 14
- [60] Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. arXiv:1905.11946, 2019. 1, 3, 4, 22

- [61] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring robustness to natural distribution shifts in image classification. Advances in Neural Information Processing Systems, 33, 2020. 3, 8
- [62] Dogancan Temel, Jinsol Lee, and Ghassan AlRegib. Cure-or: Challenging unreal and real environments for object recognition. In *International Conference on Machine Learning and Applications*, 2018. 7
- [63] Hugo Touvron, Andrea Vedaldi, Matthijs Douze, and Hervé Jégou. Fixing the train-test resolution discrepancy. In Advances in Neural Information Processing Systems, 2019. 5, 6
- [64] Hugo Touvron, Andrea Vedaldi, Matthijs Douze, and Hervé Jégou. Fixing the train-test resolution discrepancy: Fixefficientnet. arXiv:2003.08237, 2020. 4, 5, 6
- [65] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, and Hugo Larochelle. Meta-dataset: A dataset of datasets for learning to learn from few examples. *arXiv:1903.03096*, 2019. 8
- [66] Michael Tschannen, Josip Djolonga, Marvin Ritter, Aravindh Mahendran, Neil Houlsby, Sylvain Gelly, and Mario Lucic. Self-supervised learning of video-induced visual invariances. In *Conference on Computer Vision and Pattern Recognition*, 2020. 3, 8, 22
- [67] Bastiaan S Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation equivariant cnns for digital pathology. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2018. 13
- [68] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2010. 13
- [69] Cihang Xie, Mingxing Tan, Boqing Gong, Jiang Wang, Alan Yuille, and Quoc V Le. Adversarial examples improve image recognition. arXiv:1911.09665, 2019. 3, 7, 22
- [70] Qizhe Xie, Eduard Hovy, Minh-Thang Luong, and Quoc V Le. Self-training with noisy student improves imagenet classification. arXiv:1911.04252, 2019. 1, 3, 4, 7, 8, 22
- [71] Xiaohua Zhai, Avital Oliver, Alexander Kolesnikov, and Lucas Beyer. S4I: Self-supervised semi-supervised learning. In *International Conference on Computer Vision*, 2019. 3, 7, 22
- [72] Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruyssen, Carlos Riquelme, Mario Lucic, Josip Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, Lucas Beyer, Olivier Bachem, Michael Tschannen, Marcin Michalski, Olivier Bousquet, Sylvain Gelly, and Neil Houlsby. A large-scale study of representation learning with the visual task adaptation benchmark. *arXiv: 1910.04867*, 2019. 1, 2, 3, 8, 12
- [73] Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018. 14
- [74] Richard Zhang. Making convolutional networks shiftinvariant again. In *International Conference on Machine Learning*, 2019. 7

- [75] Nanxuan Zhao, Zhirong Wu, Rynson W. H. Lau, and Stephen Lin. Distilling localization for self-supervised representation learning. arXiv: 2004.06638, 2020. 8
- [76] Yiren Zhou, Sibo Song, and Ngai-Man Cheung. On classification of distorted images with deep convolutional neural networks. In *International Conference on Acoustics, Speech* and Signal Processing, 2017. 7

A. Analysis of existing robustness and transfer metrics

Here, we provide additional details related to the analyses and benchmarks presented in Section 3.

A.1. Robustness metric correlation



Figure 8: Spearman's rank correlation between accuracies on the eight robustness datasets. Samples were taken from 39 models across various model families presented in Table 2.

A.2. Dimensionality of the space of robustness metrics

To estimate how many different dimensions are measured by the robustness metrics beyond what is already explained by IMAGENET accuracy, we proceeded as follows. For each of the robustness metrics shown in Figure 8 and 10, a linear regression was fit to predict that metric's value for the 39 models, using IMAGENET accuracy as the sole predictor variable. Then, the residuals were computed for each metric by subtracting the linear regression prediction. The plot shows the fraction of variance explained for the first 4 principal components of the space of residuals of the robustness metrics. As a null hypothesis, we assumed that there is no correlation structure in the metric residuals. To construct corresponding null datasets, we randomly permuted the values for each metric independently, which destroys the correlation structure between metrics. Figure 9a shows that only the first principal component is significantly above the value expected under the null hypothesis.

A.3. Informativeness of robustness metrics

To estimate how useful different combinations of robustness metrics are for discriminating between model types, we trained logistic regression classifiers to discriminate between the 12 model groups outlined in the main paper. We consider IMAGENET accuracy as a baseline metric and therefore compare the performance of a classifier using only IMAGENET accuracy as input feature, to a classifier using IMAGENET either one (Figure 10, left) or two (Figure 10, right) additional metrics as input features. Figure 10 shows difference in accuracy to the baseline (IMAGENET) classifier. These results can serve practitioners with a limited budget as a rough guideline for which metric combinations are the most informative. In our experiments, the most informative combination of metrics in addition to IMAGENET accuracy was OBJECTNET and YOUTUBE-BB, although other combinations performed similarly within the statistical uncertainty.

A.4. Visual Task Adaptation Benchmark Details

The Visual Task Adaptation Benchmark (VTAB) [72] contains 19 tasks. Either the full dataset or 1000-example training sets may be used. We use the version with 1000-example training sets (VTAB-1k).

The tasks are divided into three groups: *natural* consists of standard natural image classification problems; *specialized* consists of domain-specific images captured with specialist equipment (e.g. medical images); *structured* consists of classification tasks that require geometric understanding of a scene. The *natural* group contains the following datasets: Caltech101 [42],



DATASET	INSTANCES	Cls.
IMAGENET [39] IMAGENET-A [30] IMAGENET-C [27] OBJECTNET [3] IMAGENET-V2 [52] IMAGENET-VID [54] YTBB-ROBUST [54]	$50\ 000750015 \times 4 \times 50\ 00018\ 57410\ 00022\ 17951\ 826$	$ \begin{array}{r} 1000\\ 200\\ 1000\\ 113\\ 1000\\ 293\\ 229 \end{array} $

(b) The name and reference, number of instances, and the number of classes overlapping with ImageNet for each dataset.

(a) The space of robustness metrics.

Figure 9: (Left) The space of robustness metrics spans approximately one statistically significant dimension after accounting for IMAGENET accuracy. Errorbars show 95% confidence intervals based on 1000 bootstrap samples (for the true data) or 1000 random permutations (for the null distribution). See Section A.2 for details. (Right) Details for the datasets used in this study. The datasets were used only for evaluation.



Figure 10: Informativeness of robustness metrics (related to Figure 4). (Left) Similar to Figure 4, but showing all 23 robustness metrics. Difference in accuracy of a logistic classifier trained to discriminate between model types based on IMAGENET accuracy plus one additional metric, compared to a classifier trained only on IMAGENET accuracy (higher is better, top 10 metrics shown). Bars show mean±s.d. of 1000 bootstrap samples from the 39 models. (Right) Increase in classifier accuracy over IMAGENET accuracy when including up to two robustness metrics as explanatory variables. The diagonal shows the single-feature values from (left).

CIFAR-100 [38], DTD [9], Flowers102 [47], Pets [49], Sun397 [68], SVHN [46]. The *specialized* group contains remote sensing datasets EuroSAT [25] and Resisc45 [8], and medical image datasets Patch Camelyon [67] and Diabetic Retinopathy [34]. The *structured* group contains the following tasks: counting and distance prediction on CLEVR [33], pixel-location and orientation prediction on dSprites [45], camera elevation and object orientation on SmallNORB [41], object distance on DMLab [4] and vehicle distance on KITTI [19].

B. Scale and OOD generalization

Training Details The models are firstly pre-trained on IMAGENET-21K and JFT, and are then fine-tuned on IMAGENET to match the label space for evaluation. We follow the pre-training and BiT-HyperRule fine-tuning setup proposed in [36].

Specifically, for pre-training, we use SGD with momentum with initial learning rate of 0.1, and momentum 0.9. We use

linear learning rate warm-up for 5000 optimization steps and multiply the learning rate by $\frac{\text{batch size}}{256}$. We use a weight decay of 0.0001. We use the random image cropping technique from [59], and random horizontal mirroring followed by resizing the image to 224×224 pixels. We use a global batch size of 1024 and train on a Cloud TPUv3-128. We pre-train models for the cross product of the following combinations:

- Dataset Size: {1.28M (1× ImageNet train set), 2.6M (2× ImageNet train set), 5.2M (4× ImageNet train set), 9M (7× ImageNet train set), 13M (10× ImageNet train set)}.
- Train Schedule (steps): {113K (90 ImageNet epochs), 229K (180 ImageNet epochs), 457K (360 ImageNet epochs), 791K (630 ImageNet epochs), 1.1M (900 ImageNet epochs)}.

For fine-tuning, we use the BiT-Hyperrule as described in [36]: batch size 512, learning rate 0.003, no weight decay, the classification head initialized to zeros, Mixup [73] with $\alpha = 0.1$, fine-tuning for 20 000 steps with 384×384 image resolution.

Additional Results Here we highlight the results equivalent to Figure 3, with the only difference that we consider subsets of the JFT [57] dataset, instead of IMAGENET-21K (Figure 11). We present the results on the synthetic dataset in Appendix D.

		ImageNet-A	ImageNet-C	ImageNet-V2	ObjectNet	ImageNet-Vid	YouTube-BB	ImageNet-Vid-W	YouTube-BB-W
teps	1120K	2.9 16.9 23.8	5.3 22.9 <mark>30.8</mark>	7.5 23.0 <mark>30.3</mark>	3.4 15.6 22.0	8.5 21.1 <mark>29.4</mark>	4.6 11.2 9.8	3.1 17.6 24.4	5.3 12.4 13.3
ing s	457K	2.4 13.6 17.4	4.1 19.8 23.0	7.5 21.3 20.2	2.1 15.2 16.7	8.1 22.2 24.4	5.0 8.7 11.8	6.4 16.8 20.4	6.0 11.8 12.5
Train	112K	0.0 5.2 6.7	0.0 8.8 10.9	0.0 11.1 12.7	0.0 6.9 8.0	0.0 14.0 16.5	0.0 6.8 6.1	0.0 9.4 10.3	0.0 7.4 7.7
		1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size
		ImageNet-A	ImageNet-C	ImageNet-V2	ObjectNet	ImageNet-Vid	YouTube-BB	ImageNet-Vid-W	YouTube-BB-W
teps	1120K	4.8 14.8 21.0	15.9 21.5 <mark>27.2</mark>	16.4 18.1 <mark>23.8</mark>	9.6 12.5 17.4	17.4 16.9 21.6	9.7 6.3 6.5	12.1 16.8 20.9	9.5 8.4 10.2
ing s	457K	4.5 12.8 16.8	14.3 <mark>20.7</mark> 23.2	16.3 19.7 19.4	8.2 14.3 16.2	15.6 20.0 21.0	11.3 7.7 9.3	15.8 15.4 <mark>20.5</mark>	9.6 9.5 9.6
Train	112K	2.7 6.9 8.3	12.8 16.3 18.5	14.4 20.5 21.7	9.1 12.7 13.0	9.4 15.8 18.6	5.8 6.9 5.6	7.9 11.4 14.9	4.7 7.8 7.1
		1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size	1M 5M 13M Dataset size

Figure 11: (**Top**) Reduction (in %) in classification error relative to the classification error of the model trained for 112k steps on 1M examples (bottom left corner) as a function of training steps and training set size. The results are for ResNet-50 trained on JFT subsets. (**Bottom**) Relative reduction (in %) in classification error going from ResNet-50 to ResNet-101x3 as a function of training steps and training set size and longer training.

C. Effect of the testing resolution

Cropping details Before applying the respective model, we first resize every image such that the shorter side has length $\lfloor 1.15 \cdot r \rfloor$ while preserving the aspect ratio and take a central crop of size $r \times r$. For the widely used 224×224 testing resolution, this leads to standard single-crop testing preprocessing, where the images are first resized such that the shorter side has length 256.

Training details for FixRes For fine-tuning to the target resolution (FixRes) we use SGD with momentum with initial learning rate of 0.004 (except for the BiT models for which we use 0.0004), and momentum 0.9, accounting for varying batch size by multiplying the learning rate with $\frac{\text{batch size}}{256}$. We train for $15\,000 \cdot \frac{\text{batch size}}{2048}$, decaying the learning rate by a factor of 10 after 1/3 and 2/3 of the iterations. The batch size is chosen based on the model size to avoid memory overflow; we use 2048 in most cases. We train on a Cloud TPUv3-64. We emphasize that we did not extensively tune the training parameters for FixRes, but chose a setting that works well across models and data sets.

Additional results In Figure 12 we provide an extended version of Figure 5 that shows the effect of FixRes for all datasets and models. In Figure 13 we plot the performance of all models and their FixRes variants as a function of the resolution.



⁽d) Two VIVI variants (R50x1 and R50x3), both co-trained with ImageNet.

Figure 12: Comparison of different types of evaluation preprocessing and resolutions. Default: Accuracy obtained for the preprocessing and resolution proposed by the authors of the respective models. Best: The accuracy when selecting the best resolution from $\{64, 128, 224, 288, 320, 384, 512, 768\}$. FixRes: Applying FixRes for the same set of resolutions and selecting the best resolution. Increasing the evaluation resolution and additionally using FixRes helps across a large range of models and pretraining datasets.



Figure 13: Comparison of different types of evaluation preprocessing and resolutions, without modifying the model and after applying FixRes. For brevity the same shorthands are used in the model names as in Figure 12.

D. Additional results on SI-SCORE, the synthetic dataset



Figure 14: (Left) Additional sample images from our synthetic dataset. (Right) From SI-SCORE, we find that an ImageNet-trained ResNet-50 has higher classification accuracy on smaller objects if it uses GroupNorm and higher accuracy on larger objects if it uses BatchNorm. Investigating this phenomena in detail is outside the scope of this paper - here we simply highlight the potential of investigating models using datasets such as SI-SCORE.

Relative performance improvement (ResNet-50, ImageNet-21K)

												_
	1.3 M	39.7	58.7	69.3	75.1	79.1	80.8	86.4	92.4	93.3	100.0	- 16
size	2.6 M	4.6	5.8	3.4	2.5	1.7	3.2	2.5	-1.9	3.7	0.0	- 8
set	5.2 M	10.1	6.7	5.3	2.7	3.5	4.3	0.4	-2.4	1.4	0.0	- 0
Data	9.0 M	10.0	10.3	7.1	6.3	5.4	5.4	2.7	0.3	-0.7	0.0	8
	13.0 M	14.1	12.3	9.9	9.1	6.3	9.3	5.8	4.1	0.9	0.0	16
		10	20	30	40	50 Area	60 1 (%)	70	80	90	100	-

	R	elative	perfor	mance	e impro	vemei	nt (Res	Net-10)1-x3,	Imagel	Net-21	<)
	1.3 M	47.3	68.8	78.5	83.1	87.9	89.2	92.8	95.9	95.7	100.0	- 16
set size	2.6 M	7.7	3.1	0.4	1.3	-2.5	-0.8	-1.5	-0.2	-2.5	0.0	- 8
	5.2 M	9.6	3.4	1.4	0.2	-1.7	-2.1	-2.3	-2.3	-3.3	0.0	- 0
Data	9.0 M	11.1	6.6	5.0	3.9	0.8	0.3	3.5	3.0	1.9	-1.6	8
	13.0 M	12.2	5.9	1.7	-1.3	-2.4	-1.8	-3.7	-2.3	-3.3	0.0	16
		10	20	30	40	50 Area	60 1 (%)	70	80	90	100	_

Relative performance improvement (ResNet-50, JFT)

	1.3 M	36.9	60.5	66.3	75.0	78.1	82.7	84.9	95.8	96.6	100.0	- 16
size	2.6 M	6.8	0.8	6.9	3.2	3.2	3.0	2.2	-2.8	-0.6	0.0	- 8
set	5.2 M	13.2	6.4	12.9	8.1	7.9	4.3	5.3	0.1	-1.4	-0.2	- 0
Data	9.0 M	11.9	5.0	8.5	3.9	6.0	4.9	5.0	-2.7	-3.3	0.0	8
	13.0 M	13.5	7.9	11.4	7.3	6.0	5.8	8.2	-1.4	0.1	0.0	16
		10	20	30	40	50 Area	60 1 (%)	70	80	90	100	_

	Relative performance improvement (ResNet-101-x3, JFT)														
0	1.3 M	48.6	68.2	76.1	82.1	84.7	86.1	87.9	95.1	96.1	99.9	- 16			
set size	2.6 M	6.8	3.3	2.0	-0.4	0.8	1.6	3.1	-0.9	-2.1	0.1	- 8			
	5.2 M	8.5	4.1	2.3	-1.2	-2.2	-1.1	0.3	-1.7	-2.3	0.1	- 0			
Date	9.0 M	13.1	6.9	5.3	1.6	1.3	0.5	3.0	1.4	-1.2	0.1	8			
	13.0 M	16.8	12.0	10.1	5.2	5.2	4.7	7.0	0.9	0.5	0.1	16			
		10	20	30	40	50 Area	60 1 (%)	70	80	90	100				

Figure 15: In the first row of both plots we show the ratio of the accuracy and the best accuracy (across all areas). For the second row (model trained on 2.6M instances), and other rows, we compute the same normalized score and visualize the difference with the first row. Larger differences imply a more uniform behavior across relative object areas. We observe that, as the dataset size increases, the average prediction accuracy across various object areas becomes more uniform. The effect is more pronounced for the larger model. As expected, the improvement is most pronounced for small object sizes covering 10-20% of the pixels.



Figure 16: In the first row of both plots we show the ratio of the accuracy and the best accuracy (across all rotations). For the second row (model trained on 2.6M instances), and other rows, we compute the same normalized score and visualize the difference with the first row. Larger differences imply a more uniform behavior across object rotations. We observe that, as the dataset size increases, the average prediction accuracy across various rotation angles becomes more uniform. The effect is more pronounced for the larger model.



Figure 17: In the first column, for each location on the grid, we compute the average accuracy. Then, we normalize each location by the 95th percentile across all locations, which quantifies the gap between the locations where the model performs well, and the ones where it under-performs (first column, dark blue vs white). Then, we consider models trained with more data, compute the same normalized score, and plot the *difference* with respect to the first column. We observe that, as dataset size increases, sensitivity to object location decreases – the outer regions improve in relative accuracy more than the inner ones (e.g. dark blue vs white in the second to fifth columns). The effect is more pronounced for the larger model.



Figure 18: In the main paper, we presented results on the location dataset when not filtering out images where the objects were partially occluded, since that would exclude many locations from the dataset. For completeness, we present results filtering out objects that are less than 50% or 75% in the image in this figure and Figure 19.

In the first column, for each location on the grid, we compute the average accuracy. Then, we normalize each location by the 95^{th} percentile across all locations, which quantifies the gap between the locations where the model performs well, and the ones where it under-performs (first column, dark blue vs white). Then, we consider models trained with more data, compute the same normalized score, and plot the *difference* with respect to the first column. We observe that, as dataset size increases, sensitivity to object location decreases – the outer regions improve in relative accuracy slightly more than the inner ones (e.g. dark blue vs white in the second to fifth columns). The effect is more pronounced for the larger model. We filter out all test images for which the foreground object is not at least 50% within the image.



Figure 19: In the main paper, we presented results on the location dataset when not filtering out images where the objects were partially occluded, since that would exclude many locations from the dataset. For completeness, we present results filtering out objects that are less than 50% or 75% in the image in this figure and Figure 18.

In the first column, for each location on the grid, we compute the average accuracy. Then, we normalize each location by the 95^{th} percentile across all locations, which quantifies the gap between the locations where the model performs well, and the ones where it under-performs (first column, dark blue vs white). Then, we consider models trained with more data, compute the same normalized score, and plot the *difference* with respect to the first column. We observe that, as dataset size increases, sensitivity to object location decreases – the outer regions improve in relative accuracy more than the inner ones (e.g. dark blue vs white on the second and third columns). The effect is harder to see since most pixels near the edges have been filtered out — here we filter out all test images for which the foreground object is not at least 75% within the image.

E. Overview of model abbreviations

MODEL NAME	Түре	TRAINING DATA	ARCHITECTURE	Depth	Сн.
r50-imagenet-100	SUPERVISED	IMAGENET	RESNET	50	1
r50-imagenet-10	SUPERVISED	IMAGENET, 10%	RESNET	50	1
BIT-IMAGENET-R50-X1	SUPERVISED [36]	IMAGENET	ResNet	50	1
BIT-IMAGENET-R50-X3	SUPERVISED [36]	IMAGENET	ResNet	50	3
BIT-IMAGENET-R101-X1	SUPERVISED [36]	IMAGENET	RESNET	101	1
BIT-IMAGENET-R101-X3	SUPERVISED [36]	IMAGENET	RESNET	101	3
BIT-IMAGENET21K-R50-X1	SUPERVISED [36]	IMAGENET21k	RESNET	50	1
BIT-IMAGENET21K-R50-X3	SUPERVISED [36]	IMAGENET21k	RESNET	50	3
BIT-IMAGENET21K-R101-X1	SUPERVISED [36]	IMAGENET21k	RESNET	101	1
BIT-IMAGENET21K-R101-X3	SUPERVISED [36]	IMAGENET21k	RESNET	101	3
BIT-JFT-R50-X1	SUPERVISED [36]	JFT	RESNET	50	1
BIT-JFT-R50-X3	SUPERVISED [36]	JFT	RESNET	50	3
BIT-JFT-R101-X1	SUPERVISED [36]	JFT	ResNet	101	1
BIT-JFT-R101-X3	SUPERVISED [36]	JFT	RESNET	101	3
BIT-JFT-R152-X4	SUPERVISED [36]	JFT	RESNET	50	4
r50-imagenet-10-exemplar	Self-sup. & cotraining [71]	IMAGENET, 10%	RESNET	50	1
R50-IMAGENET-10-ROTATION	Self-sup. & cotraining [71]	IMAGENET, 10%	RESNET	50	1
r50-imagenet-100-exemplar	Self-sup. & cotraining [71]	IMAGENET	ResNet	50	1
R50-IMAGENET-100-ROTATION	Self-sup. & cotraining [71]	IMAGENET	RESNET	50	1
SIMCLR-1X-SELF-SUPERVISED	Self-supervised [6], fine tuning	IMAGENET	RESNET	50	1
SIMCLR-2X-SELF-SUPERVISED	Self-supervised [6], fine tuning	IMAGENET	ResNet	50	2
SIMCLR-4X-SELF-SUPERVISED	SELF-SUPERVISED [6], FINE TUNING	IMAGENET	RESNET	50	4
simclr-1x-fine-tuned-10	Self-supervised [6], fine tuning	IMAGENET, 10%	ResNet	50	1
SIMCLR-2X-FINE-TUNED-10	Self-supervised [6], fine tuning	IMAGENET, 10%	RESNET	50	2
SIMCLR-4X-FINE-TUNED-10	SELF-SUPERVISED [6], FINE TUNING	IMAGENET, 10%	RESNET	50	3
simclr-1x-fine-tuned-100	Self-supervised [6], fine tuning	IMAGENET	ResNet	50	1
SIMCLR-2X-FINE-TUNED-100	SELF-SUPERVISED [6], FINE TUNING	IMAGENET	RESNET	50	2
SIMCLR-4X-FINE-TUNED-100	SELF-SUPERVISED [6], FINE TUNING	IMAGENET	RESNET	50	4
efficientnet-std-b0	SUPERVISED [60]	IMAGENET	EFFICIENTNET	18	1
efficientnet-std-b4	SUPERVISED [60]	IMAGENET	EFFICIENTNET	37	1
EFFICIENTNET-ADV-PROP-B0	SUPERVISED & ADVERSARIAL [69]	IMAGENET	EFFICIENTNET	18	1
EFFICIENTNET-ADV-PROP-B4	SUPERVISED & ADVERSARIAL [69]	IMAGENET	EfficientNet	37	1
EFFICIENTNET-ADV-PROP-B7	SUPERVISED & ADVERSARIAL [69]	IMAGENET	EfficientNet	64	2
efficientnet-noisy-student-b0	SUPERVISED & DISTILLATION [70]	IMAGENET	EFFICIENTNET	18	1
EFFICIENTNET-NOISY-STUDENT-B4	SUPERVISED & DISTILLATION [70]	IMAGENET	EfficientNet	37	1
EFFICIENTNET-NOISY-STUDENT-B7	SUPERVISED & DISTILLATION [70]	IMAGENET	EFFICIENTNET	64	2
VIVI-1X	Self-sup. & cotraining [66]	YT8M, IMAGENET	RESNET	50	1
VIVI-3X	Self-sup. & cotraining [66]	YT8M, IMAGENET	RESNET	50	3
BIGBIGAN-LINEAR	BIDIRECTIONAL ADVERSARIAL [14]	IMAGENET	RESNET	50	1
BIGBIGAN-FINETUNE	BIDIRECTIONAL ADVERSARIAL [14]	IMAGENET	RESNET	50	1

Table 2: Overview of models used in this study. SUP. abbreviates for supervised pre-training. CH. refers to the width multiplier for the number of channels.