
Generalization to translation shifts: a study in architectures and augmentations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We provide a detailed evaluation of various image classification architectures
2 (convolutional, vision transformer, and fully connected MLP networks) and data
3 augmentation techniques towards generalization to large translation shifts. We
4 make the following observations: (a) In the absence of data augmentation, all archi-
5 tectures, including convolutional networks suffer degradation in performance when
6 evaluated on translated test distributions. Understandably, both the in-distribution
7 accuracy and degradation to shifts is significantly worse for non-convolutional
8 architectures. (b) Across all architectures, even a minimal augmentation of 4 pixel
9 random crop improves the robustness of performance to much larger magnitude
10 shifts of up to $1/4$ of image size (8-16 pixels) in the test data – suggesting a form
11 of meta generalization from augmentation. For non-convolutional architectures,
12 while the absolute accuracy is still low, we see dramatic improvements in robust-
13 ness to large translation shifts. (c) With sufficiently advanced augmentation (4
14 pixel crop+RandAugmentation+Erasing+MixUp) pipeline all architectures can be
15 trained to have competitive performance, in terms of in-distribution accuracy as
16 well as generalization to large translation shifts.

17 1 Introduction

18 Convolutional neural networks (ConvNets) are a natural architectural choice for a variety of computer
19 vision tasks. The built-in structure from localization and translation equivariance of the convolutional
20 layers is intrinsically useful in many image processing scenarios [Krizhevsky et al., 2012, LeCun
21 et al., 1989, Fukushima and Miyake, 1982]. For over a decade ConvNets were the backbone of
22 computer vision and continue to be one of the most important class of models. At the same time,
23 advances in large datasets and data augmentation techniques have made it possible to train general
24 purpose architectures to be competitive with ConvNets despite lacking any image specific priors.
25 Most popular among these are the Vision Transformers (ViTs) and their variants [Vaswani et al.,
26 2017, Dosovitskiy et al., 2020, Touvron et al., 2021c,b]. When pretrained on ultra-large datasets
27 like ImageNet-21k (14 million images) or JFT-300/3B (300 million/3 billion weakly labeled images,
28 respectively), ViTs can outperform similarly pretrained ConvNets on diverse vision tasks. While the
29 scale of the data was original thought to be crucial, follow up work show that competitive accuracies
30 can also be achieved in small-to-medium data regimes using advanced data augmentation [Touvron
31 et al., 2021b] or optimization techniques [Chen et al., 2021]. Detailed experimentation on competitive
32 benchmarks by Steiner et al. [2021] showed that ViTs trained with extensive data augmentation can
33 recover the performance gains from $\sim 10x$ larger independently annotated dataset. In such large
34 data or extensive augmentation regimes, even simpler fully connected multi-layer perceptron (MLP)
35 architectures [Tolstikhin et al., 2021, Touvron et al., 2021a] can achieve competitive performance. In
36 this work we focus on the role of data augmentation in learning from general purpose architectures.

37 *Can data augmentation capture the inductive biases of carefully designed architectures?* Beyond
38 accuracy, the design of architectures is often motivated by domain knowledge of desired invariances.
39 One of the fundamental image priors is the invariance of object labels to spatial shift or translation of
40 its position. Indeed, the success and motivation of ConvNets is often attributed to their component
41 convolutional operators being definitionally translation equivariant. In this work, we study to what
42 extent the shift-equivariant convolutional layers make ConvNets robust to translation shifts compared
43 to ViTs and MLPs? How effective are data augmentation techniques in encouraging similar behaviors?

44 There is a rich literature on understanding the translation invariance properties of ConvNets (*cf.*
45 Section 1.1). It has been shown that despite the equivariance of the convolutional operator, other
46 architectural components like non-linearities and strides can cause the networks lose their invariance
47 to translations [Zhang, 2019, Azulay and Weiss, 2019, Chaman and Dokmanic, 2021, Engstrom et al.,
48 2019, Xiao et al., 2018, Alsallakh et al., 2020]. Despite these findings, it is reasonable to believe
49 that ConvNets will still have substantial edge in robustness to translation when compared to ViT or
50 MLP models. At the same time, the competitive performance of latter architectures suggests that
51 good image priors can also be learned from rich training data. The goal of this work is to quantify the
52 relative effectiveness of architectures design and data augmentation towards robustness to translation
53 shifts. Algorithmically, these are complementary tools for incorporating domain knowledge. On one
54 hand, it is conceptually simpler to generate data with desired invariant transformation rather than hard
55 code it in the architecture. On the other hand, data augmentation only provides a weak supervision
56 about the invariant properties and is limited by the biases in the training samples.

57 **Generalization vs invariance.** In our experiments, rather than chase the gold standard of *strict* in-
58 variance on all inputs, we work with a data-centric measure of *generalization to structured translation*
59 *shifts* in the test distribution. We evaluate our models for accuracy on specific out-of-distribution test
60 datasets where the object locations are systematically shifted without creating additional distortions
61 or domain gaps from training distribution. See Section 2 for full setup.

62 In the strict sense, a classifier $f : X \rightarrow Y$ on input space X is invariant to set of spatial translations
63 T if for all $t \in T$ and all inputs $x \in X$, $f(t(x)) \approx f(x)$. Often we are not concerned with
64 performance on all inputs but rather on typical samples from a task distribution, say $(x, y) \sim$
65 D . We thus evaluate our models a structured out-of-distribution accuracy, wherein we train on
66 samples from D , but test on translation shifted inputs $t(x)$, *i.e.*, generalized accuracy to shift t is
67 $\mathbb{E}_{(x,y) \sim D} \mathbf{1}[f(t(x)) = y]$. This evaluation also differs from adversarial robustness to translation shifts
68 as studied in Engstrom et al. [2019], Xiao et al. [2018]. In our notation, the adversarial accuracy
69 metric would be $\mathbb{E}_{(x,y) \sim D} \mathbf{1}[f(t_x(x)) \neq y]$, where $t_x \in T$ is adversarially chosen per-sample (x, y) .
70 Finally, we clarify that data augmentation does change the training distribution D . If inputs x were to
71 be augmented with their respective transformations $\{t(x) : t \in T\}$ (and no other augmentation), then
72 there is no distribution shift. However, our experiments are never in the no-distribution shift regime.
73 Importantly, all our augmentation pipelines uses random crop of at most 4 pixels (some are further
74 restricted to 1 or 2 pixels), but we evaluate our models on much larger translation shifts of up to 1/4th
75 of image dimension (8 pixels on CIFAR and 16 on TinyImageNet). Thus, any generalization to larger
76 translation shifts from limited augmentation can be thought of as a form of meta-generalization.

77 1.1 Related work

78 **Translation invariance in ConvNets.** ConvNets have been extensively studied on various ap-
79 proximate measures of invariance to translation shifts including, but not limited to, Goodfellow et al.
80 [2009], Zeiler and Fergus [2014], Fawzi and Frossard [2015], Kanbak et al. [2018], Azulay and Weiss
81 [2019], Zhang [2019], Engstrom et al. [2019], Kayhan and Gemert [2020], Chaman and Dokmanic
82 [2021]. For example, Zeiler and Fergus [2014] provide visualizations of the hidden layer filters in
83 early ConvNets that show their sensitivity to small changes in translation, scale, and rotation. Similar
84 visual evaluation was also more recently used in Alsallakh et al. [2020]. Other work Fawzi and
85 Frossard [2015], Kanbak et al. [2018], Azulay and Weiss [2019], Zhang [2019], Engstrom et al.
86 [2019] focus on more quantitative measures of invariance such as mean change in top-1 prediction or
87 class probabilities between images within a distortion range. These works collectively establish that
88 despite the built-in inductive bias, ConvNets are not “truly” invariant even to small translation. Some
89 recent work, importantly [Zhang, 2019, Chaman and Dokmanic, 2021], address this shortcoming
90 of ConvNets by designing new architectural modifications to ConvNets. These prior works show
91 that strict spatial invariance is a strong measure that even ConvNets with their built-in priors do

92 not satisfy. We choose a data-centric measure in our evaluation that is more practical for general
93 architectures. Our experimental and evaluation protocol has some key differences from these prior
94 works studying invariance: (a) Our evaluation is designed to introduce only translations shifts in
95 the test data without confounding with any other domain gap between train and test distributions.
96 For example, in the experiments of Azulay and Weiss [2019], Zhang [2019] the test samples are
97 padded and resized during evaluation, while the models were trained on regular dataset without such
98 padding or resizing. This leads the test dataset to differ from the training data in ways that are not
99 just translation of the objects. (b) Our evaluation is non-adversarial in that it measures degradation
100 in average test loss and not in the worst-case drop in performance on any single image. Further, by
101 comparing the drop in performance from translation shift to the unperturbed in-distribution accuracy,
102 we inherently down-weight the non-robustness on hard-to-learn inputs on which the classifier had
103 inaccurate prediction to begin with. (c) Finally, our evaluation does not penalize models from learning
104 position dependent features as observed by Kayhan and Gemert [2020]. For example, a network has
105 the flexibility to use its large representation power to create a separate model for an object (say a
106 cat) at each pixel location — while such a model will be inefficient, it would still do a good job at
107 detecting cats in translated test distribution. These differences are nuanced but significant, which
108 makes our results complementary to prior work in this space. Such differences may or may not be
109 important depending on the application.

110 **Robustness under adversarial perturbations and other OOD benchmarks.** In a work closest to
111 ours, Engstrom et al. [2019] study test accuracy degradation from adversarially chosen translations
112 and rotations on test images. While much of the work in adversarial robustness focus on ℓ_∞ or
113 ℓ_2 norm bounded perturbations to the inputs, Engstrom et al. [2019] show that ConvNets can be
114 effectively “attacked” even when the perturbed inputs are merely small rotations and/or translations
115 of the input. Similar adversarial attack on ConvNets based on more complex spatial transformations
116 was previously studied in [Xiao et al., 2018]. Our evaluation is closer to “random perturbation”
117 evaluation, which is only briefly explored in Engstrom et al. [2019]. In comparison to these studies,
118 our distribution shift is picked non-adversarially and independent of input samples.

119 Complementing the work on adversarial robustness, there has also been lot of interest in evaluating
120 models on other out-of-distribution robustness. Generalization to out-of-distribution test datasets is a
121 broad umbrella topic, and aside from adversarial robustness, many work also focus on robustness
122 performance on benchmarks for benign “natural” perturbations Recht et al. [2019], Hendrycks and
123 Dietterich [2018], Koh et al. [2021], Djolonga et al. [2021]. For ConvNets, inspired by the neural
124 networks *scaling laws* line of work, Djolonga et al. [2021] probe for relation between OOD robustness
125 and learning choices like data size, model size, optimization algorithm, and model choices like model
126 sizes, and normalization (they do not consider data augmentation in detail though). They also propose
127 a synthetic benchmark SI-SCORE for controlled image invariance evaluation. Expanding on this line
128 of work, Yung et al. [2021], Bhojanapalli et al. [2021], Bai et al. [2021], Mahmood et al. [2021], Shao
129 et al. [2021], Paul and Chen [2021], Pinto et al. [2022], compare ConvNets and ViTs on adversarial
130 robustness and/or out-of-distribution robustness benchmarks.

131 Aside from the methodological differences, much of the work in this space has focused on the relative
132 merits, demerits, and robustness of different model choices and the role of training data size. The
133 effects of data augmentation is only minimally considered, if at all. In contrast, our goal in this work
134 is to specifically quantify how much different data augmentation pipelines can capture the inductive
135 biases in a carefully designed architecture.

136 2 Experimental setup

137 All our experiments are conducted on the CIFAR-10, CIFAR-100, and TinyImageNet datasets. Since
138 our study involves training from scratch and testing on large models in numerous configurations, it is
139 beyond the scope of the paper to extend such a detailed study to larger benchmarks like full ImageNet.
140 Moreover, we emphasize that our goal here is not to get the state-of-the-art accuracy/robustness
141 on benchmarks, but rather to understand how much data augmentation captures the benefits of the
142 convolutional architecture. Arguably, it is also the small data regime where the inductive biases
143 from architecture and/or augmentations play more important roles. In ultra-large scale datasets,
144 accuracy/robustness might naturally come from dataset size itself rather than model priors.

145 In the main paper we focus on results from CIFAR-10 and CIFAR-100 datasets which consists of
 146 32×32 pixel RGB images balanced across 10 and 100 classes, respectively. We defer the discussion
 147 on TinyImageNet, which a subset of the more diverse ImageNet benchmark, to Appendix C. To study
 148 large translation shifts without introducing domain gaps, we modify the dataset by symmetrically
 149 padding all the CIFAR images with 8 pixels ($1/4$ of image size) on each side. The padded pixels
 150 contain the mean channel values of the entire training dataset, which ensures that the channel-wise
 151 means and standard deviations across training dataset remains the same as the original un-padded
 152 dataset (see illustration in Figure 1). This padded dataset allows us to evaluate large translation shifts
 153 of up to 16 pixels (Hamming distance) in the test dataset without creating additional confounding
 154 factors. *Importantly, in all the shifted test sets, there is no cropping or loss of the image content and*
 155 *the entire image is available to the network at the same scale as seen during training.* After padding
 156 with a mean-valued canvas, we resize the resulting $48 \times 48 \times 3$ images to $224 \times 224 \times 3$ (the standard
 157 input size for ImageNet) using bilinear interpolation. This up-sampling step helps avoid extensive
 158 hyperparameter tuning of the models, especially, the ViT and MLP models.

159 We briefly discuss the alternative evaluation methodologies. The more natural random cropping of
 160 images to evaluate robustness to translation is inherently limited by the number of pixels we can
 161 faithfully forgo without losing information and hence cannot capture large translation shifts. In prior
 162 work, Azulay and Weiss [2019] also used similar padded images to investigate translation shifts. A
 163 key difference in our methodology is that we have our entire training and testing pipeline on the
 164 preprocessed images (with padding), while the latter paper evaluated models pretrained on standard
 165 ImageNet without any padding – this creates an uncontrolled distribution shifts. Another technical
 166 difference is that Azulay and Weiss [2019] downsampled the images, which leads to loss in resolution,
 167 while our preprocessing is non-lossy. Finally, synthetic benchmarks such as SI-Score proposed
 168 in Djolonga et al. [2021] are a good alternative to our setup. However, for our simple controlled
 169 setting, the conceptual advantage of padding is that it does not change the natural distribution of
 170 foreground and background which would be lost in the cut-and-paste protocol of Djolonga et al.
 171 [2021]. Furthermore, the segmentation process in Djolonga et al. [2021] appears to be not perfect
 172 which might create additional confounders. To further ensure that our “synthetic” padding is benign,
 173 we verified that the test accuracy of our models trained on the padded dataset is comparable to the
 174 standard train-test pipeline on $32 \times 32 \times 3$ CIFAR inputs without any padding (see Table 1).



Figure 1: Sample images from the preprocessing steps

175 **Training** All the models are trained on the mean-padded datasets for 1600 epochs on $8 \times V100$
 176 GPUs. We use implementations (with suitable modifications) of the models from various open source
 177 repositories, most notably Wightman [2019] and Liu [2017]. We performed basic hyperparameter
 178 tuning in a small grid around the parameters reported in the respective papers. The exact value of
 179 hyperparameters used in experiments along with code are provided in the supplementary material.
 180 All the evaluation metrics reported in this paper are median performance over 3 runs.

181 **2.1 Architectures**

182 Our goal is to compare fundamentally different architectures on generalization to translation shifts.
 183 After initial experiments with different variants, we choose the following models in our evaluation.

- 184 • **resnet18_bn (11M parameters)**: We use ResNets He et al. [2016] as our representative ConvNet.
 185 In our initial experiments, larger ResNets and the other ConvNets like RegNet Radosavovic et al.
 186 [2020] did not yield qualitative difference in performance on our datasets.
- 187 • **resnet18_gn (11M parameters)**: It has been observed that batch normalization often leads to
 188 poor performance in transfer learning as the batch statistics from source task could be widely off
 189 for target task (see *e.g.*, Kolesnikov et al. [2020], Wu and Johnson [2021]). The same reasoning
 190 also applies when dealing with distribution shifts, where batch statistics could become irrelevant

191 even when the test distribution shifts systematically Djolonga et al. [2021]. To overcome this, we
192 consider a variant of *resnet18* with group normalization and weight standardization Wu and He
193 [2018], Qiao et al. [2019] in place of batch normalization. This modification indeed leads to more
194 stable performance in our experiments, specially in the absence of data augmentation.

195 • **antialiased_resnet18 (11M parameters)**: Azulay and Weiss [2019], Zhang [2019] show that
196 the downsampling layers (pooling/strides) make ConvNets non-shift invariant. To remedy this,
197 Zhang [2019] proposed a specialized modification to ConvNets to improve their invariance to
198 spacial shifts by introducing a *BlurPool* layer as an antialiasing-filter before downsampling. This
199 constitute a model with more specialized priors about translation invariance built into its design.

200 • **cait_xxs36 (17M parameters)**: CaiT architecture Touvron et al. [2021c] is a variant of the basic
201 vision transformer (ViT) Dosovitskiy et al. [2020] that leads to more efficient training of deeper
202 models. We use CaiT as our representative transformer model as it had the best performance in
203 initial experiments. Other ViT variants, including larger models and the distilled variant DeiT
204 Touvron et al. [2021b] did not provide significant performance boost on our small scale datasets.

205 • **resmlp_12 (18M parameter)**: Among the MLP models for image classification, we tried MLP-
206 mixer Tolstikhin et al. [2021] and ResMLP Touvron et al. [2021a] in our initial experiments. We
207 stick with ResMLP for detailed experimentation as it had slightly better performance.

208 Although the above model configurations are not the state-of-the-art on larger benchmarks, on smaller
209 scale CIFAR and TinyImageNet datasets, they have competitive performance as their larger or more
210 complex counterparts. Since our goal is to evaluate the relative degradation in performance with
211 translation shifts, we do not overly optimize for top-accuracy.

212 2.2 Augmentations

213 We first consider four data augmentation pipelines while training the models described above. In the
214 appendix, we look at more minimal augmentations to elaborate on our findings.

215 • **No Augmentation (*NoAug*)**: We use this setting as a baseline for purely evaluating the merits of
216 an architecture in generalization to translation shifts.

217 • **Basic augmentation (*BA*)**: The basic augmentation consists of a random flip and a random crop
218 with up to 4×4 pixel padding. This minimal augmentation has been a de-facto standard in many
219 vision tasks, and it already gives over 5% boost in accuracy even without considering its effects
220 on distribution shifts. Note that unlike in standard training pipeline, with our padding of training
221 images, the random crop does not lose any original image pixels.

222 • **Advanced augmentation (*AA*)**: The current slate of image data augmentation techniques are
223 more varied and less intuitive compared to the simple transformations described above. In our
224 version of advanced augmentation (*AA*), we use the following pipeline: first we apply (a) the
225 basic augmentation (*BA*) described above, then (b) RandAugment Cubuk et al. [2019], then (c)
226 random erasing Zhong et al. [2020], and finally use (d) MixUp Inoue [2018]. RandAugment uses a
227 randomly chosen composition of transformations from a predefined list. We use the standard Rand
228 Augment list from Wightman [2019] but without the `TranslationX` and `TranslationY`
229 as these are covered with more control within basic augmentation.

230 • **AA without translation (*AA(no tr)*)**: We also consider a variant of advanced augmentation (*AA*)
231 where we remove any augmentations that are explicitly related translation shifts. Specifically, we
232 remove random crop from basic augmentation (*BA*) and in the RandAugment transformations
233 list, along with previously removed `TranslationX` and `TranslationY`, we also remove
234 `ShearX` and `ShearY`. In this case, any improvements over *NoAug* arise only indirectly.

235 2.3 In-distribution test accuracy

236 Table 1 gives the test accuracies on the in-distribution test dataset, which is preprocessed with the
237 same padding configuration as the train dataset (*i.e.*, with 8 pixels padded symmetrically on all sides).
238 We use this as a reference performance without any distribution shifts (a.k.a. the *in-distribution*
239 accuracy). Ideally, we would expect a classifier that that learns good image priors to maintain their
240 reference performance even when the test dataset is shifted by object invariant properties.

		resent18_bn	resnet18_gn	antialiased_resnet18	cait_xxs36	resmlp_12
CIFAR-10	NoAug	90.85±0.19	91.48±0.09	92.55±0.21	77.58±0.11	79.99±0.15
	BA	96.10±0.05	95.96±0.06	95.63±0.15	87.69±0.50	87.73±0.07
	AA(no-tr)	96.35±0.06	96.06±0.07	96.98±0.07	95.09±0.19	91.90±0.10
	AA	97.74±0.03	98.03±0.06	97.77±0.10	97.25±0.01	96.03±0.09
CIFAR-100	NoAug	67.62±0.65	64.62±0.29	70.37±0.09	43.43±0.12	52.79±0.21
	BA	78.68±0.18	78.42±0.08	78.20±0.16	57.62±0.78	60.52±0.60
	AA(no-tr)	74.56±0.46	74.09±0.22	77.62±0.25	77.74±1.40	65.43±0.30
	AA	82.98±0.14	82.09±0.27	81.54±0.17	82.46±0.27	78.63±0.28

Table 1: Accuracy on in-distribution test dataset (*i.e.*, test and train images have same padding).

241 3 Generalization to translation shifts

242 With the flexibility of padded pixels in our preprocessed training data, we can now create test datasets
 243 with up to 16 pixel translations (in Hamming distance) from the training image distribution by moving
 244 the test images anywhere within the 48×48 frame. This allows us to evaluate large translation while
 245 not distorting the contents of the image itself. In the extreme locations (see, *e.g.*, corners of grid in
 246 Figure 2) there only 25% overlap with the training distribution.

247 For each trained model, we can look at a grid of 17×17 test evaluations on modified test datasets.
 248 Each cell in the evaluation grid corresponds to the position of the 32×32 test images within the
 249 48×48 frame (see illustration in Figure 2 for a *resnet18_bn* network). The center cell of the
 250 grid acts as the reference performance and corresponds to the no distribution shift, *i.e.*, the test
 251 images are centered on the frame, same as the train images. As we move away from the center, we
 252 analogously translate the position of the object image in the test dataset. The model is then evaluated
 253 for classification accuracy on the shifted test dataset. Thus, the generalization or robustness of trained
 models to translation shifts can be comprehensively summarized by such a grid.

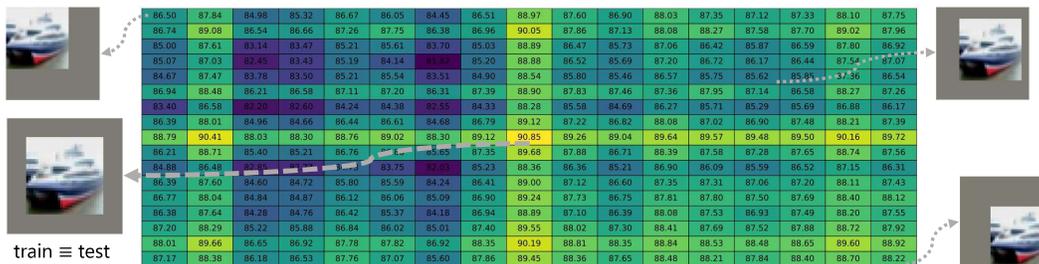


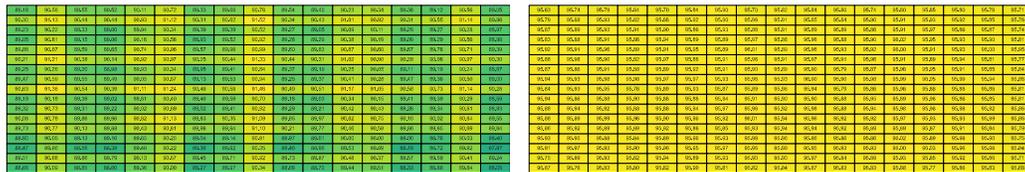
Figure 2: Generalization to translation shifts of a *resnet18_bn* trained without data augmentation (*NoAug*) for 1600 epochs on *CIFAR-10* dataset: Each cell in the grid corresponds to the model performance on a test dataset with specific positioning of the object image. The center cell corresponds to no translation shift from training; and the distance of the cells from the center corresponds to respective position shift between train and test dataset. The values in each cell is the accuracy on the shifted test dataset. The color corresponds to “relative” drop in performance from the in-distribution performance: yellow the maximum accuracy in the grid (typically the center cell), while the dark blue is saturated at 90% of the max-value in the grid, *i.e.*, 10% drop in accuracy. For *resnet18_bn* +*NoAug* on *CIFAR-10*, max accuracy on the grid is 90.85% while min accuracy is 81.82%.

254

255 3.1 Case study: Convolutional networks

256 We first describe some observations from our initial experimentation exclusively on ConvNets trained
 257 on *CIFAR-10*. Our evaluations on other architectures are provided in Section 4. In this subsection,
 258 we mainly compare variations of *resnet18* trained in different configurations. First, in Figure 2 along
 259 with the illustration of translation shift grid, we show the evaluation of the basic *resnet18_bn* (with
 260 batch normalization) on generalization to translation shifts in *CIFAR-10*. This model was trained
 261 without any augmentation to demonstrate the baseline performance of convolutional architecture.
 262 We observe that, *despite the built-in image priors in ConvNets and despite using a weaker notion*

263 *than translation invariance, there are significant drops in performance on even small translations*
 264 *shifts. For example, in the non-yellow cells close to center in Figure 2, we notice that a two pixel*
 265 *hamming distance translation can lead to > 5% drop in performance.* It is worth forward referencing
 266 at this point that in spite of these drops in performance, when compared to other architectures (see
 267 Section 4), we find that ConvNets are relatively more resilient to translation shifts in the absence of
 268 data augmentation – the worst case drop in performance is $\sim 10\%$ for ConvNets, while for other
 269 architectures, the drop in performance could be as large as 30 – 50% (see Figure 4).



(a) *resnet18_gh*+NoAug: max = 91.48, min = 87.97. (b) *resnet18_gh*+BA: max = 95.96, min = 95.63

Figure 3: Generalization to translation shifts of *resnet* variants using the same evaluation as in Figure 2. (a) *resnet18_gh* network group norm and weight standardization again trained without any augmentation. (b) *resnet18_gh* trained with basic augmentation (BA) consisting of random horizontal flip and random crop at most 4 pixels. For quick reference, the sub-captions mention the max and min accuracy of the models over the grid.

270 **Batchnorm vs groupnorm+weight standardization.** Prior work most notably Azulay and Weiss
 271 [2019], Zhang [2019], Chaman and Dokmanic [2021] attribute the lack of invariance to strides and
 272 ReLU non-linearities in the standard networks. We believe these factors also affect the weaker
 273 notion of generalization to translation shifts that we study. Indeed, the “near” periodic locations of
 274 high-performing (yellow) cells is in line with prior observations that strides in ResNets would tend
 275 to have periodic translation invariance (Azulay and Weiss [2019], Fukushima and Miyake [1982]).
 276 Additionally, our experiments suggest that batch normalization is yet another factor that contributes
 277 to lack of translation invariance. In Figure 3(a), we show the performance of *resnet18_gh* model
 278 trained with the same configuration as *resnet18_bn* in Figure 2, but with a modification that all the
 279 batch normalization layers replaced by group normalization and weight standardization Wu and He
 280 [2018], Qiao et al. [2019], as was also done in Kolesnikov et al. [2020]. We see that this simple
 281 modification already improves the generalization of the ResNet to translation shifts. In hindsight, it is
 282 understandable that batch normalization would have detrimental effects when the test distribution
 283 shifts from the training distribution as the batch statistics obtained from exponential moving average
 284 of training statistics no longer remains accurate Wu and Johnson [2021].

285 **Training with basic augmentation (BA).** Experiments in Azulay and Weiss [2019] show that even
 286 after using data augmentation, ConvNets are not translation invariant when tested on inputs that are
 287 not from the same manifold as training images. On the other hand, in our evaluation of a weaker
 288 notion of generalization to translation shifts, we observe a different conclusion. In Figure 3(b), we
 289 show that even a small nudge using basic augmentation (BA) can make the models remarkably robust.
 290 A simple augmentation using random crops of up to 4 pixels and horizontal flips, not only improve
 291 in-distribution accuracy by over 5%, but also make the networks near-perfectly robust to up to 8 pixel
 292 translations in test distribution – indicating a form of meta-generalization from augmentations. In the
 293 appendix we provide further evidence of such meta-generalization: (a) In (Appendix B) ConvNets
 294 trained with a more minimal *BA-lite* and *BA-liter* augmentation with smaller range of random crops
 295 of at most 2 and 1 pixels, respectively already help robustness to translation shifts of up to 8 pixels on
 296 CIFAR-10, and (b) In Appendix C our results on TinyImageNet dataset show robustness to a larger
 297 range of translation shifts of up to 16 pixel shifts, even though BA still uses only 4 pixel shifts.

298 4 Architectures and augmentations for generalization to translation shifts

299 The summary grid view of evaluations on translation shifts (as in Figure 2-3) is more comprehensive
 300 and we will revisit them in the appendix. However, it is not ideal for comparing different configura-
 301 tions of architectures and augmentations. In this section, we use an alternative visualization and plot
 302 the test accuracies as a function of Hamming distance between the position of images in the test and
 303 training datasets. The performance of all our models and augmentations are summarized in Figure 4.

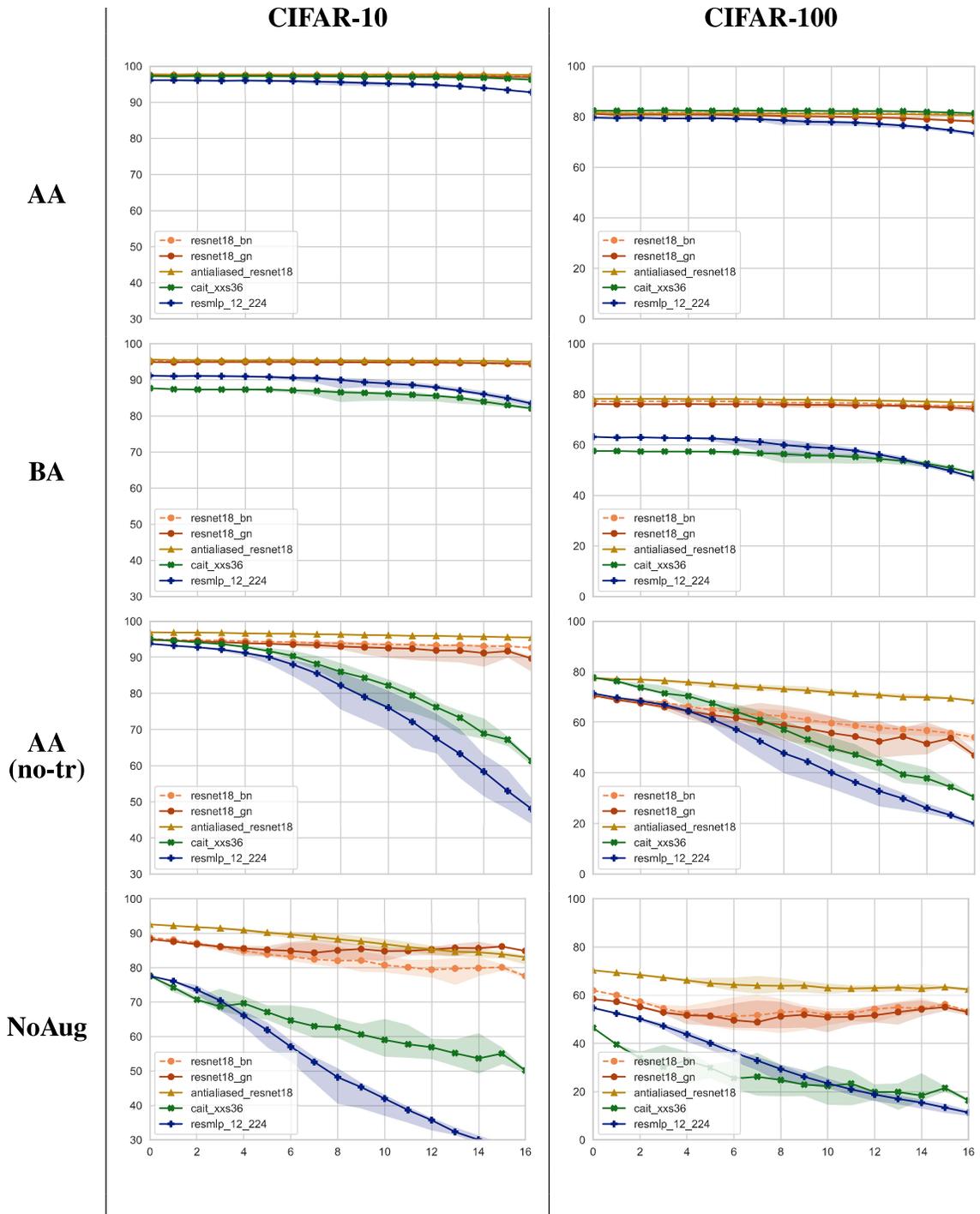


Figure 4: Generalization to translation shifts. The x -axis in each plot is the Hamming distance of the location of test images (within the 48×48 frame) relative to the location in the training images (center of the frame). The larger the x -axis value, the larger is the shift from train distribution. For a given value of x , there might be multiple test configurations that are x -Hamming distance away (like shift of 1 pixel to top, bottom, right, or left when $x = 1$). The median of these values is plotted as line plot, while the shaded region covers the min and max values of the list. The left and right plots correspond to accuracies on the CIFAR-10 and the more challenging CIFAR-100 dataset, respectively. The y -axis for each dataset (column) is normalized to be on the same scale.

304 4.1 Key takeaways

305 We make the following observations that are supported by our experiments. We provide additional
306 experiments and discussions from our study including evaluation on TinyImageNet in the appendix.

- 307 1. Without data augmentation, even ConvNets with designed architecture has a noticeable drop in
308 performance on spatially shifted test dataset. From Figures 2-3, one can see that merely switching
309 batchnorm with groupnorm appears to mitigate some performance gap. Using a specialized
310 antialiased modification by Zhang [2019] further improves the performance, albeit not to perfect
311 invariance as the non-linearities could be a source of breaking spacial invariance Chaman and
312 Dokmanic [2021]. Nonetheless, the architectural inductive biases in these ConvNets are still
313 useful, if not perfect, as we can predictably see that the drop in performance is much more dramatic
314 for the non-convolutional architectures without any image priors. In CaiT and ResMLP, without
315 augmentation, even 1 – 2 pixel translation shift can lead to dramatic drop in performance.
- 316 2. In the other extreme, with an Advanced Augmentation (AA) pipeline, all the architectures are
317 remarkably robust even to large translation shifts in test distribution. Note that even with AA, the
318 maximum translation augmentation we provide (in the form of random crop) is at most 4 pixels (8
319 pixels in Hamming distance), but we see robustness to up to 8 pixel shifts (16 pixels in Hamming
320 distance). This supports a notion of meta generalization in robustness performance. In Appendix C,
321 we see that ResNets continue to be robust to even larger translation shifts of up to 16 pixels on
322 each direction on TinyImageNet dataset. To further support the idea of meta-generalization, we
323 also show in Appendix B, even more minimal 1 or 2 pixel random crop augmentation already
324 boosts robustness to translation shifts.
- 325 3. Furthermore, with AA, the performance on in-distribution test error becomes significantly closer
326 for all architectures. Specially, the performance of resnet18 and cait_xxs36 are statistically
327 identical in this setting even though we trained on the small-medium scale CIFAR datasets. The
328 performance of resmlp_12 is however relatively suboptimal even with AA pipeline. First, there
329 is non-trivial gap in the in-distribution accuracies. Secondly, the robustness to translation shifts
330 is not nearly as good as with ResNet or CaiT. Despite these differences, even for ResMLP, the
331 augmentations dramatically boost the generalization to translation shifts and the differences in
332 relative drop start to appear only after 10-12 pixel hamming distance shifts in test distribution.
333 These experiments suggest that with sufficient augmentation, the relative benefits or shortcomings
334 of the architectures are effectively diminished.
- 335 4. Even with a minimal Basic Augmentation (BA), we see significant improvement in robustness to
336 translations. In fact, for ConvNets, BA is sufficient to achieve the near perfect generalization on
337 our canvas. This further highlights the benefit of the built-in inductive biases in ConvNet. For
338 non-convolutional architectures, this simple augmentation is not sufficient to achieve optimal
339 absolute test performance, but the relative robustness is still uniformly improved.
- 340 5. Finally, an intriguing phenomenon is observed when training with advance augmentation but
341 without translation related augmentations (AA(no-tr)). Here the absolute test accuracy of all the
342 models improves (presumably from learning some useful priors). For ConvNets on CIFAR-10,
343 even such indirect augmentation is effective in making the models robust to translation shifts – but
344 this does not appear to uniformly hold across datasets, so the conclusion might be spurious. For
345 transformer and MLP architectures, the robustness does not improve significantly, even though the
346 in-distribution accuracies are significantly higher.
- 347 6. Somewhat tangentially, our experiment slightly supports the position that data augmentation can
348 recover some of the benefits of large datasets even when learning with general architectures like
349 ViT and MLP are beneficial. In ImageNet scale datasets this was previously observed in Touvron
350 et al. [2021c], Steiner et al. [2021] for ViTs and Touvron et al. [2021a] for MLPs. Our experiments
351 show similar validation even on small CIFAR-10 datasets.

352 In summary, even though convolutional networks are not invariant or robust to translations in absolute
353 sense, they clearly fare much better compared to other general-purpose architectures. Specially,
354 ConvNets learn robust models with minimal augmentation, while it appears the transformer and MLP
355 architectures require more sophisticated augmentations. At the same time, our experiments suggest
356 that data augmentation can enforce learning of the right inductive bias with comparable or more
357 effectiveness than the network architecture.

358 **References**

- 359 Bilal Alsallakh, Narine Kokhlikyan, Vivek Miglani, Jun Yuan, and Orion Reblitz-Richardson. Mind
360 the pad-cnns can develop blind spots. *arXiv preprint arXiv:2010.02178*, 2020.
- 361 Aharon Azulay and Yair Weiss. Why do deep convolutional networks generalize so poorly to small
362 image transformations? *Journal of Machine Learning Research*, 20:1–25, 2019.
- 363 Yutong Bai, Jieru Mei, Alan L Yuille, and Cihang Xie. Are transformers more robust than cnns?
364 *Advances in Neural Information Processing Systems*, 34, 2021.
- 365 Srinadh Bhojanapalli, Ayan Chakrabarti, Daniel Glasner, Daliang Li, Thomas Unterthiner, and
366 Andreas Veit. Understanding robustness of transformers for image classification. In *Proceedings*
367 *of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10231–10241,
368 October 2021.
- 369 Anadi Chaman and Ivan Dokmanic. Truly shift-invariant convolutional neural networks. In *Proceed-*
370 *ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3773–3783,
371 2021.
- 372 Xiangning Chen, Cho-Jui Hsieh, and Boqing Gong. When vision transformers outperform resnets
373 without pretraining or strong data augmentations. *arXiv preprint arXiv:2106.01548*, 2021.
- 374 Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment:
375 Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF Conference on*
376 *Computer Vision and Pattern Recognition*, pages 113–123, 2019.
- 377 Josip Djolonga, Jessica Yung, Michael Tschannen, Rob Romijnders, Lucas Beyer, Alexander
378 Kolesnikov, Joan Puigcerver, Matthias Minderer, Alexander D’Amour, Dan Moldovan, et al. On
379 robustness and transferability of convolutional neural networks. In *Proceedings of the IEEE/CVF*
380 *Conference on Computer Vision and Pattern Recognition*, pages 16458–16468, 2021.
- 381 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
382 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
383 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
384 *arXiv:2010.11929*, 2020.
- 385 Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. Exploring
386 the landscape of spatial robustness. In *International Conference on Machine Learning*, pages
387 1802–1811. PMLR, 2019.
- 388 Alhussein Fawzi and Pascal Frossard. Manitest: Are classifiers really invariant? In *British Machine*
389 *Vision Conference (BMVC)*, 2015.
- 390 Kunihiko Fukushima and Sei Miyake. Neocognitron: A self-organizing neural network model for a
391 mechanism of visual pattern recognition. In *Competition and cooperation in neural nets*, pages
392 267–285. Springer, 1982.
- 393 Ian Goodfellow, Honglak Lee, Quoc Le, Andrew Saxe, and Andrew Ng. Measuring invariances in
394 deep networks. *Advances in neural information processing systems*, 22, 2009.
- 395 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
396 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
397 pages 770–778, 2016.
- 398 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
399 corruptions and perturbations. In *International Conference on Learning Representations*, 2018.
- 400 Elad Hoffer, Tal Ben-Nun, Itay Hubara, Niv Giladi, Torsten Hoefer, and Daniel Soudry. Augment
401 your batch: Improving generalization through instance repetition. In *Proceedings of the IEEE/CVF*
402 *Conference on Computer Vision and Pattern Recognition*, pages 8129–8138, 2020.
- 403 Hiroshi Inoue. Data augmentation by pairing samples for images classification. *arXiv preprint*
404 *arXiv:1801.02929*, 2018.

- 405 Can Kanbak, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. Geometric robustness of deep
406 networks: analysis and improvement. In *Proceedings of the IEEE Conference on Computer Vision
407 and Pattern Recognition*, pages 4441–4449, 2018.
- 408 Osman Semih Kayhan and Jan C van Gemert. On translation invariance in cnns: Convolutional layers
409 can exploit absolute spatial location. In *Proceedings of the IEEE/CVF Conference on Computer
410 Vision and Pattern Recognition*, pages 14274–14285, 2020.
- 411 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-
412 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A
413 benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*,
414 pages 5637–5664. PMLR, 2021.
- 415 Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and
416 Neil Houlsby. Big transfer (bit): General visual representation learning. In *Computer Vision–ECCV
417 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*,
418 pages 491–507. Springer, 2020.
- 419 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolu-
420 tional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 421 Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- 422 Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne
423 Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition.
424 *Neural computation*, 1(4):541–551, 1989.
- 425 Kuang Liu. Train cifar10 with pytorch. [https://github.com/kuangliu/
426 pytorch-cifar](https://github.com/kuangliu/pytorch-cifar), 2017.
- 427 Kaleel Mahmood, Rigel Mahmood, and Marten Van Dijk. On the robustness of vision transformers
428 to adversarial examples. In *Proceedings of the IEEE/CVF International Conference on Computer
429 Vision*, pages 7838–7847, 2021.
- 430 Sayak Paul and Pin-Yu Chen. Vision transformers are robust learners. *arXiv preprint
431 arXiv:2105.07581*, 2(3), 2021.
- 432 Francesco Pinto, Philip HS Torr, and Puneet K Dokania. An impartial take to the cnn vs transformer
433 robustness contest. *arXiv preprint arXiv:2207.11347*, 2022.
- 434 Siyuan Qiao, Huiyu Wang, Chenxi Liu, Wei Shen, and Alan Yuille. Weight standardization. 2019.
- 435 Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing
436 network design spaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
437 Pattern Recognition*, pages 10428–10436, 2020.
- 438 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers
439 generalize to imagenet? In *International Conference on Machine Learning*, pages 5389–5400.
440 PMLR, 2019.
- 441 Rulin Shao, Zhouxing Shi, Jinfeng Yi, Pin-Yu Chen, and Cho-Jui Hsieh. On the adversarial robustness
442 of visual transformers. *arXiv e-prints*, pages arXiv–2103, 2021.
- 443 Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas
444 Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *arXiv
445 preprint arXiv:2106.10270*, 2021.
- 446 Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Un-
447 terthiner, Jessica Yung, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, et al. Mlp-mixer: An
448 all-mlp architecture for vision. *arXiv preprint arXiv:2105.01601*, 2021.
- 449 Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaaeldin El-Nouby, Edouard
450 Grave, Gautier Izacard, Armand Joulin, Gabriel Synnaeve, Jakob Verbeek, et al. Resmlp:
451 Feedforward networks for image classification with data-efficient training. *arXiv preprint
452 arXiv:2105.03404*, 2021a.

- 453 Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve
454 Jegou. Training data-efficient image transformers & distillation through attention. In *International*
455 *Conference on Machine Learning*, volume 139, pages 10347–10357, July 2021b.
- 456 Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going
457 deeper with image transformers. In *Proceedings of the IEEE/CVF International Conference on*
458 *Computer Vision (ICCV)*, pages 32–42, October 2021c.
- 459 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
460 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
461 *systems*, 30, 2017.
- 462 Ross Wightman. Pytorch image models. [https://github.com/rwightman/
463 pytorch-image-models](https://github.com/rwightman/pytorch-image-models), 2019.
- 464 Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on*
465 *computer vision (ECCV)*, pages 3–19, 2018.
- 466 Yuxin Wu and Justin Johnson. Rethinking" batch" in batchnorm. *arXiv preprint arXiv:2105.07576*,
467 2021.
- 468 Chaowei Xiao, Jun-Yan Zhu, Bo Li, Warren He, Mingyan Liu, and Dawn Song. Spatially transformed
469 adversarial examples. In *International Conference on Learning Representations*, 2018.
- 470 Jessica Yung, Rob Romijnders, Alexander Kolesnikov, Lucas Beyer, Josip Djolonga, Neil Houlsby,
471 Sylvain Gelly, Mario Lucic, and Xiaohua Zhai. Si-score: An image dataset for fine-grained analysis
472 of robustness to object location, rotation and size. *arXiv preprint arXiv:2104.04191*, 2021.
- 473 Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In
474 *European conference on computer vision*, pages 818–833. Springer, 2014.
- 475 Richard Zhang. Making convolutional networks shift-invariant again. In *International conference on*
476 *machine learning*, pages 7324–7334. PMLR, 2019.
- 477 Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data aug-
478 mentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages
479 13001–13008, 2020.