

# ForAug: Mitigating Biases and Improving Vision Transformer Training by Recombining Foregrounds and Backgrounds

Anonymous CVPR submission

Paper ID 4792

## Abstract

001 Transformers, particularly Vision Transformers (ViTs), have  
 002 achieved state-of-the-art performance in large-scale image  
 003 classification. However, they often require large amounts  
 004 of data and can exhibit biases, such as center or size bias,  
 005 that limit their robustness and generalizability. This paper  
 006 introduces ForAug, a novel data augmentation operation  
 007 that addresses these challenges by explicitly imposing invari-  
 008 ances into the training data, which are otherwise part of the  
 009 neural network architecture. ForAug is constructed by using  
 010 pretrained foundation models to separate and recombine  
 011 foreground objects with different backgrounds. This recom-  
 012 bination step enables us to take fine-grained control over  
 013 object position and size, as well as background selection.  
 014 We demonstrate that using ForAug significantly improves  
 015 the accuracy of ViTs and other architectures by up to 4.5  
 016 percentage points (p.p.) on ImageNet, which translates to  
 017 7.3 p.p. on downstream tasks. Importantly, ForAug not  
 018 only improves accuracy but also opens new ways to analyze  
 019 model behavior and quantify biases. Namely, we introduce  
 020 metrics for background robustness, foreground focus, center  
 021 bias, and size bias and show that using ForAug during train-  
 022 ing substantially reduces these biases. In summary, ForAug  
 023 provides a valuable tool for analyzing and mitigating bi-  
 024 ases, enabling the development of more robust and reliable  
 025 computer vision models. Our code and dataset are publicly  
 026 available at [<url>](#).

## 027 1. Introduction

028 Image classification, a fundamental task in computer vision  
 029 (CV), involves assigning labels to images from a set of cate-  
 030 gories. It underpins a wide range of applications, like medi-  
 031 cal diagnosis [37, 48], autonomous driving [51], and object  
 032 recognition [2, 14, 16] and facilitates large-scale pretrain-  
 033 ing [9, 29, 45], and progress evaluation in CV [21, 34]. The  
 034 advent of large-scale datasets, particularly ImageNet [7],  
 035 served as a catalyst for the rise of large-scale CV mod-

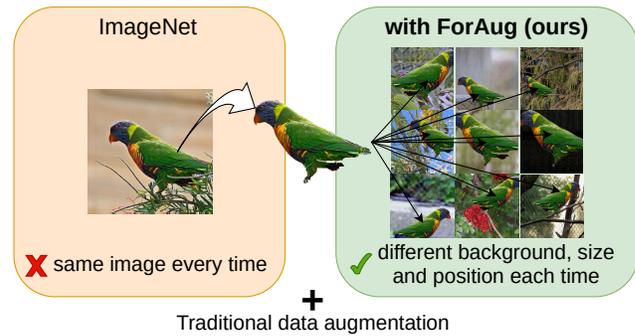


Figure 1. Comparison of traditional image classification training and training when using ForAug. ForAug recombines foreground objects with different backgrounds each epoch, thus creating a more diverse training set. We still apply strong traditional data augmentation afterwards.

els [15, 23] and remains the most important CV benchmark  
 for more than a decade [15, 23, 46, 53]. While tradition-  
 ally, convolutional neural networks (CNNs) have been the  
 go-to architecture in CV, Transformers [47], particularly the  
 Vision Transformer (ViT) [9], have emerged as a powerful  
 alternative and go-to architecture, demonstrating superior  
 performance in various vision tasks, including image classi-  
 fication [2, 50, 53, 56, 61].

Data augmentation is a key technique for training image  
 classification models. Traditional augmentation methods,  
 such as cropping, flipping, or color shifts, are commonly  
 employed to increase data diversity [40, 55], but remain  
 bound to existing image compositions. While these preserve  
 the images’ semantic meaning, their ability to teach spatial  
 invariances is limited. While combinations of these data aug-  
 mentations are still used today, they originally were proposed  
 to benefit CNNs. However, the architectural differences of  
 CNNs and Transformers suggest that the latter might benefit  
 from different data augmentation strategies. In particular, the  
 self-attention mechanism, unlike a CNN, is not translation  
 equivariant [8, 36], meaning that the model is not designed  
 to understand the spatial relationships between pixels.

058 Recognizing that Transformers need to learn spatial rela-  
059 tionships directly from data, we propose *ForAug*, a data aug-  
060 mentation method that makes these relationships explicit by  
061 recombining foreground objects with diverse backgrounds.  
062 Thus, *ForAug* goes beyond existing image compositions and  
063 encodes desired invariances directly into the training data  
064 (see Figure 1). Applying *ForAug* to a dataset like ImageNet  
065 is a two-step process: (1) We separate the foreground objects  
066 in ImageNet from their backgrounds, using an open-world  
067 object detector [35] and fill in the background in a neutral  
068 way using an object removal model [41, 43]. (2) This al-  
069 lows us to then recombine any foreground object with any  
070 background on the fly, creating a highly diverse training set.  
071 By exploiting the control over foreground size and position  
072 during recombination, *ForAug* explicitly teaches spatial in-  
073 variances that image classification models typically must  
074 learn implicitly. We show that using *ForAug* additionally  
075 to strong traditional data augmentation increases the model  
076 accuracy of Transformers by up to 4.5 p.p. on ImageNet and  
077 reduces the error rate by up to 7.3 p.p. in downstream tasks.

078 Beyond training, *ForAug* becomes a diagnostic tool for  
079 analyzing model behavior and biases, when used during eval-  
080 uation. We utilize our control over the image distribution  
081 to measure a model’s background robustness (by varying  
082 the choice of background), foreground focus (by leverag-  
083 ing our knowledge about the placement of the foreground  
084 object), center bias (by controlling position), and size bias  
085 (by controlling size). These analyses provide valuable in-  
086 sights into model behavior and biases, which is crucial for  
087 model deployment and future robustness optimizations. We  
088 show that training using *ForAug* significantly reduces all of  
089 these biases. We make our code for *ForAug* and the out-  
090 put of *ForAug*’s segmentation phase on ImageNet publicly  
091 available<sup>1</sup> to facilitate further research.

## 092 Contributions

- 093 • We propose *ForAug*, a novel data augmentation scheme,  
094 that recombines objects and backgrounds. *ForAug* allows  
095 us to move beyond the (possibly biased) image composi-  
096 tions in the dataset while preserving label integrity.
- 097 • We show that training a standard ViT using *ForAug* leads  
098 to up to 4.5 p.p. improved accuracy on ImageNet-1k and  
099 7.3 p.p. on downstream tasks.
- 100 • We propose novel *ForAug*-based metrics to analyze and  
101 quantify fine-grained biases of trained models: Back-  
102 ground Robustness, Foreground Focus, Center Bias, and  
103 Size Bias. We show that *ForAug* significantly reduces  
104 these biases by encoding invariance that benefits ViT into  
105 the training data.

<sup>1</sup>Link will go here.

## 2. Related Work 106

**Data Augmentation for Image Classification** Data aug- 107  
mentation is a crucial technique for improving the perfor- 108  
mance and generalization of image classification models. 109  
Traditional augmentation strategies rely on simple geomet- 110  
ric or color-space transformations like cropping, flipping, 111  
rotation, blurring, color jittering, or random erasing [60] to 112  
increase the diversity of the training data without chang- 113  
ing their semantic meaning. With the advent of Vision 114  
Transformers, new data augmentation operations like Patch- 115  
Dropout [28] have been proposed. Other transformations like 116  
Mixup [59], CutMix [57], or random cropping and patching 117  
[44] combine multiple input images. These simple transfor- 118  
mations are usually bundled to form more complex augmen- 119  
tation policies like AutoAugment [4] and RandAugment [5], 120  
or 3-augment [46] which is optimized to train a ViT. For a 121  
general overview of data augmentation techniques for image 122  
classification, we refer to Shorten and Khoshgoftaar [40], Xu 123  
et al. [55]. 124

We build upon these general augmentations by introduc- 125  
ing a novel approach to explicitly separate objects and back- 126  
grounds for image classification, allowing us to – unlike 127  
these basic transformations – move beyond dataset image 128  
compositions. Our approach is used additionally to strong 129  
traditional techniques to improve performance and reduce 130  
biases. 131

**Copy-Paste Augmentation** The copy-paste augmentation 132  
[13], which is used only for object detection [13, 39] and 133  
instance segmentation [26, 52], involves copying segmented 134  
objects from one image and pasting them onto another. 135  
While typically human annotated segmentation masks are 136  
used to extract the foreground objects, other foreground 137  
sources have been explored, like 3D models [18] and pre- 138  
trained object-detection models for use on objects on white 139  
background [10] or synthetic images [11]. [19] apply copy- 140  
paste as an alternative to CutMix in image classification, but 141  
they do not shift the size or position of the foregrounds and 142  
use normal dataset images as backgrounds. 143

Unlike prior copy-paste methods that overlay objects, 144  
*ForAug* extracts foregrounds and replaces their backgrounds 145  
with semantically neutral fills, thereby preserving label in- 146  
tegrity while enabling controlled and diverse recombination. 147

**Model robustness evaluation** Evaluating model robust- 148  
ness to various image variations is critical for understand- 149  
ing and improving model generalization. Datasets like 150  
ImageNet-C [17] and ImageNet-P [17] introduce common 151  
corruptions and perturbations. ImageNet-E [25] evaluates 152  
model robustness against a collection of distribution shifts. 153  
Other datasets, such as ImageNet-D [58], focus on varying 154  
background, texture, and material, but rely on synthetic data. 155

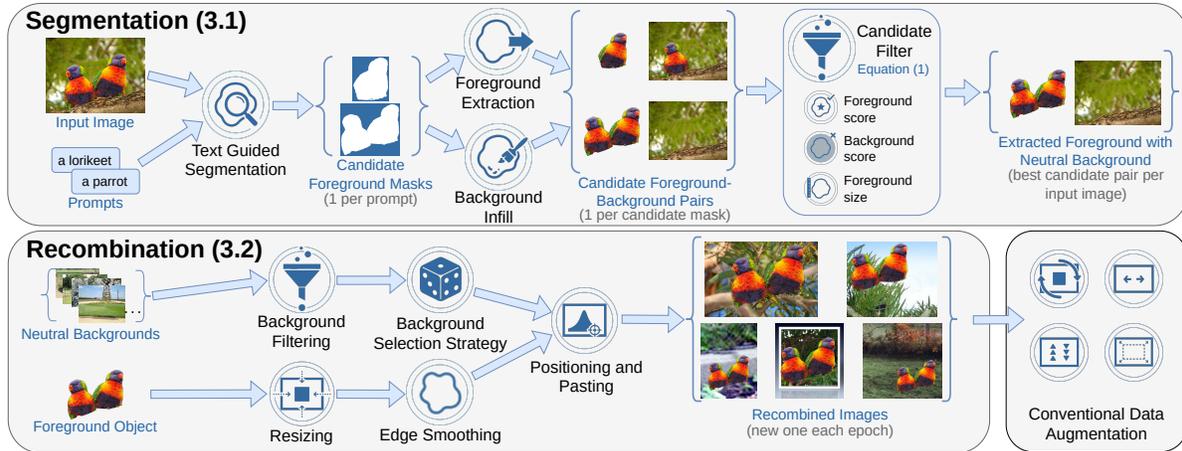


Figure 2. Overview of *ForAug*. The data creation consists of two stages: Segmentation (offline, Section 3.1), where we segment the foreground objects from the background and fill in the background. Recombination (online, Section 3.2), where we combine the foreground objects with different backgrounds to create new samples. After recombination, we apply strong, commonly used augmentation policies.

Stylized ImageNet [12] investigates the impact of texture changes. ImageNet-9 [54] explores background variations using segmented images, but backgrounds are often artificial.

In contrast to these existing datasets, which are used only for evaluation, *ForAug* provides fine-grained control over foreground object placement, size, and background selection, enabling a precise and comprehensive analysis of specific model biases within the context of a large-scale, real-world image distribution. As *ForAug* also provides controllable training set generation, it goes beyond simply measuring robustness to actively improving it through training.

### 3. *ForAug* (Method)

We introduce *ForAug*, a data augmentation designed to enhance Transformer training by embedding spatial invariances—which Transformers would otherwise need to learn implicitly—directly into the training data. *ForAug* comprises two distinct stages: Segmentation and Recombination. Both stages are illustrated in Figure 2.

#### 3.1. Segmentation

The segmentation stage isolates the foreground objects and their corresponding backgrounds. We then fill the background using a pretrained object-removal model, producing visually plausible [41], neutral scenes ready for recombination. This stage is computed once offline and the results are stored for the recombination stage.

First, foreground objects are detected and segmented from their backgrounds using a prompt-based segmentation model to exploit the classification datasets labels. We use the state-of-the-art Grounded SAM [35], which is based on Grounding DINO [27] and SAM [22]. The prompt we use is “a <class name>, a type of

<object category>”, where <class name> is the specific name of the objects class as defined by the dataset and <object category> is the broader category of the object. The <object category> guides the segmentation model towards the correct object in case the <class name> alone is too specific. This can be the case with prompts like “sorrel” or “guenon”, where the more general name “horse” or “monkey” is more helpful. We derive the <object category> from the WordNet hierarchy, using the immediate hypernym.

We iteratively extract  $n$  foreground masks for each dataset-image, creating prompts by going one hypernym up the WordNet-tree each step (e.g. “a sorrel, a type of horse”, “a horse, a type of equine”, ...). Masks that are very similar, with a pairwise IoU of at least 0.9, are merged. The output is a set of masks delineating the foreground objects and the backgrounds. We select the best mask per image (according to Equation (1)) in a later filtering step, described below.

First, an inpainting model that is specifically optimized to remove objects from images, such as LaMa [43] or Attentive Eraser [41], is used to inpaint the foreground regions in the backgrounds. Then, to ensure the quality of the foregrounds and the neutral background images, we select a foreground/background pair (for each dataset-image) from the  $\leq n$  variants we have extracted and infilled in the previous steps. Using an ensemble  $E$  of six ViT, ResNet, and Swin Transformer models pretrained on the original dataset, we select the foreground/background pair that maximizes foreground performance while minimizing the performance on the background and size of the foreground. For each model  $m \in E$ , we predict the score of the ground truth class  $c$  on the foreground  $fg$  and background  $bg$  and weigh these

220 with the size  $\text{size}(\cdot)$  in number of pixels according to:

$$\begin{aligned}
 \text{score}(\text{fg}, \text{bg}, c) &= \log \left( \frac{1}{|E|} \sum_{m \in E} \mathbb{P}[m(\text{fg}) = c] \right) \\
 &+ \log \left( 1 - \frac{1}{|E|} \sum_{m \in E} \mathbb{P}[m(\text{bg}) = c] \right) \quad (1) \\
 &+ \lambda \log \left( 1 - \left| \frac{\text{size}(\text{fg})}{\text{size}(\text{bg})} - \varepsilon \right| \right).
 \end{aligned}$$

222 We run a hyperparameter search using a manually annotated  
223 subset of foreground/background variants to find the factors  
224 in Equation (1):  $\lambda = 2$  and  $\varepsilon = 0.1$ .

225 Finally, we filter out backgrounds that are largely infilled,  
226 as these tend to be overly synthetic and do not carry much  
227 information (see the supplementary material). Although the  
228 segmentation stage is computational overhead, it is a one-  
229 time cost with results that can be reused across experiments  
230 (see the supplementary material for details). In summary, we  
231 factorize the dataset into a set of foreground objects with a  
232 transparent background and a set of diverse backgrounds per  
233 class. The next step is to recombine these, before applying  
234 other common data augmentation operations during training.

### 235 3.2. Recombination

236 The recombination stage, performed online during training,  
237 combines the foreground objects with different backgrounds  
238 to create new training samples. For each object, we follow  
239 the pipeline of: Pick an appropriate background, resize it to  
240 a fitting size, and place it in the background image. Through  
241 this step, we expose the model to variations beyond the  
242 image compositions of the dataset.

243 For each foreground object, we sample a background using  
244 one of the following strategies: (1) the original image  
245 background, (2) the set of backgrounds from the same class,  
246 or (3) the set of all possible backgrounds. These sets are trad-  
247 ing off the amount of information the model can learn from  
248 the background against the diversity of new images created.  
249 In each epoch, each foreground object is seen exactly once,  
250 but a background may appear multiple times.

251 The selected foreground is resized based on its relative  
252 size within its original image and the relative size of the  
253 original foreground in the selected background image. The  
254 final size is randomly selected from a 30% range around  
255 upper and lower limits ( $s_u$  and  $s_l$ ), based on the original  
256 sizes. To balance the size of the foreground and that of  
257 the backgrounds original foreground, the upper and lower  
258 limit  $s_u$  and  $s_l$  are set to the mean or range of both sizes,  
259 depending on the foreground size strategy: *mean* or *range*.

260 The resized foreground is then placed at a random po-  
261 sition within the background image. To more seamlessly  
262 integrate the foreground, we apply a Gaussian blur with  
263  $\sigma \in [\frac{\sigma_{\max}}{10}, \sigma_{\max}]$ , inspired by the standard range for the Gaus-  
264 sian blur operation in [46], to the foreground’s alpha-mask.

Table 1. Ablation of the design decisions in the segmentation phase of *ForAug* on TinyImageNet. The first line is our baseline, while the other lines are using *ForAug*. We use basic settings with the *same* background strategy during recombination for this experiment.

Detect. Prompt	Infill Model	TinyImageNet Accuracy [%]	
		ViT-Ti	ViT-S
<b>TinyImageNet</b>		66.1 ± 0.5	68.3 ± 0.7
specific	LaMa [43]	65.5 ± 0.4	71.2 ± 0.5
general	LaMa [43]	66.4 ± 0.6	72.9 ± 0.6
general	Att. Eraser [41]	67.5 ± 1.2	72.4 ± 0.5

265 We can apply standard data augmentation techniques in  
266 two modes: Either we apply all augmentations to the re-  
267 combined image, or we apply the cropping and resizing to  
268 the background only and then apply the other augmenta-  
269 tions after recombination. The first mode mirrors standard  
270 augmentation practice, whereas the second one ensures the  
271 foreground object remains fully visible.

272 We experiment with a constant mixing ratio, or a linear  
273 or cosine annealing schedule that increases the amount of  
274 images from the original dataset over time. The mixing ratio  
275 acts as a probability of selecting an image from the original  
276 dataset; otherwise, an image with the same foreground is  
277 recombined using *ForAug*, ensuring each object is seen once  
278 per epoch. The recombination stage is designed to be paral-  
279 lelized on the CPU during training and thus does not impact  
280 training time (see supplementary material for details).

## 281 4. Experiments

282 We conduct a comprehensive suit of experiments to vali-  
283 date the effectiveness of our approach, comparing ImageNet-  
284 training with and without *ForAug* for 10 different models.  
285 Furthermore, we assess the impact of using *ForAug* for pre-  
286 training on multiple fine-grained downstream datasets. Fi-  
287 nally, we exploit *ForAug*’s control over the image distribu-  
288 tion to quantify model behaviors and biases. We always  
289 report the mean and standard deviation of three independent  
290 training runs.

### 291 4.1. Design Choices of ForAug

292 We start by ablating the design choices of *ForAug* on TinyIm-  
293 ageNet [24], a subset of ImageNet containing 200 categories  
294 with 500 images each. Table 1 presents ablations for seg-  
295 mentation and Table 2 for recombination.

296 **Prompt.** First, we evaluate the type of prompt used to de-  
297 tect the foreground object. Here, the *general* prompt, which  
298 contains the class and the more general object category, out-  
299 performs only having the class name (*specific*).

300 **Inpainting.** Among inpainting models, Attentive  
301 Eraser [41] produces slightly better results compared to  
302 LaMa [43] (+0.5 p.p. on average). For inpainting examples,

Table 2. Ablation of the recombination phase of *ForAug* on TinyImageNet (top) and ImageNet (bottom). The first experiments use the initial segmentation settings with LaMa [43].

FG. size	Augment. Order	BG Strat.	BG. Prune	Original Mixing	Edge Smooth.	Accuracy [%]	
						ViT-Ti	ViT-S
<b>TinyImageNet</b>							
mean	crop→paste	same	-	-	-	66.1 ± 0.5	68.3 ± 0.7
range	crop→paste	same	-	-	-	64.6 ± 0.5	70.0 ± 0.6
range	crop→paste	same	-	-	-	65.5 ± 0.4	71.2 ± 0.5
range	crop→paste	same	-	-	-	67.5 ± 1.2	72.4 ± 0.5
range	paste→crop	same	-	-	-	67.1 ± 1.2	72.9 ± 0.5
range	paste→crop	same	1.0	-	-	67.0 ± 1.2	73.0 ± 0.3
range	paste→crop	same	0.8	-	-	67.2 ± 1.2	72.9 ± 0.8
range	paste→crop	same	0.6	-	-	67.5 ± 1.0	72.8 ± 0.7
range	paste→crop	same	0.8	$p = 0.2$	-	69.8 ± 0.5	75.0 ± 0.3
range	paste→crop	same	0.8	$p = 0.33$	-	69.5 ± 0.4	75.2 ± 1.0
range	paste→crop	same	0.8	$p = 0.5$	-	70.3 ± 1.0	74.2 ± 0.2
range	paste→crop	same	0.8	linear	-	70.1 ± 0.7	74.9 ± 0.8
range	paste→crop	same	0.8	reverse lin.	-	67.6 ± 0.2	73.2 ± 0.3
range	paste→crop	same	0.8	cos	-	71.3 ± 1.0	75.7 ± 0.8
range	paste→crop	same	0.8	cos	$\sigma_{\max} = 4.0$	70.0 ± 0.8	75.5 ± 0.7
range	paste→crop	orig.	0.8	cos	$\sigma_{\max} = 4.0$	67.2 ± 0.9	69.9 ± 1.0
range	paste→crop	all	0.8	cos	$\sigma_{\max} = 4.0$	70.1 ± 0.7	77.5 ± 0.6
<b>ImageNet</b>							
range	paste→crop	same	0.8	cos	-	-	79.1 ± 0.1
range	paste→crop	same	0.8	cos	$\sigma_{\max} = 4.0$	-	80.5 ± 0.1
range	paste→crop	same	0.8	cos	$\sigma_{\max} = 4.0$	-	80.7 ± 0.1
range	paste→crop	all	0.8	cos	$\sigma_{\max} = 4.0$	-	81.4 ± 0.1

303 see the supplementary material.

304 **Foreground size** significantly impacts performance. Em-  
 305 ploying a *range* of sizes during recombination, rather than  
 306 a fixed *mean* size, boosts accuracy by approximately 1 p.p.  
 307 This suggests that the added variability is beneficial.

308 **Order of data augmentation.** Applying all aug-  
 309 mentations after foreground-background recombination  
 310 (*paste→crop→color*) improves ViT-S’s performance com-  
 311 pared to applying crop-related augmentations before pasting  
 312 (*crop→paste→color*). ViT-Ti results are ambiguous.

313 **Background pruning.** When it comes to the back-  
 314 grounds to use, we test different pruning thresholds ( $t_{\text{prune}}$ )  
 315 to exclude backgrounds with large inpainting. A threshold of  
 316  $t_{\text{prune}} = 1.0$  means that we use all backgrounds that are not  
 317 fully infilled. Varying  $t_{\text{prune}}$  has minimal impact. We choose  
 318  $t_{\text{prune}} = 0.8$  to exclude predominantly artificial backgrounds.

319 **Mixing** *ForAug*-augmented samples with the original Im-  
 320 ageNet data proves crucial. While constant and linear mixing  
 321 schedules improve performance over no mixing by 2 – 3 p.p.  
 322 compared to only augmented samples, the cosine annealing  
 323 schedule proves optimal, boosting accuracy by 3 – 4 p.p.

324 **Edge smoothing.** We evaluate the impact of using Gaus-  
 325 sian blurring to smooth the edges of the foreground masks.  
 326 For larger models, this gives us a slight performance boost  
 327 on the full ImageNet (second to last line in Table 2).

328 **Background strategy.** Another point is the allowed  
 329 choice of background image for each foreground object.  
 330 We compare using the original background, a background  
 331 from the same class, and any background. These strategies  
 332 go from low diversity and high shared information content  
 333 between the foreground and background to high diversity  
 334 and low shared information content. For *ViT-Ti*, the latter two

Table 3. Accuracy of ViT-S on TinyImageNet (TIN) in percent using *ForAug* with different foreground position distributions by varying the Bates parameter  $\eta$ . The best performance is achieved when using the uniform distribution ( $\eta = 1$ ) for training.

Bates Parameter during training	TIN w/o <i>ForAug</i>	TIN w/ <i>ForAug</i>				
		$\eta = -3$	$-2$	$1/-1$	$2$	$3$
Baseline	68.9	60.5	60.2	60.8	62.6	63.1
$\eta = -3$	71.3	79.3	79.5	79.1	79.3	79.1
$\eta = -2$	71.5	80.0	78.7	79.3	79.1	78.8
$\eta = 1/-1$	72.3	79.5	78.9	80.2	79.7	80.4
$\eta = 2$	71.3	78.2	77.8	79.1	79.6	79.9
$\eta = 3$	71.4	77.2	76.9	78.6	79.6	79.7

Table 4. Dataset statistics for TinyImageNet and ImageNet with and without *ForAug*. For *ForAug* we report the number of foreground/background pairs.

Dataset	Classes	Training Images	Validation Images
TinyImageNet	200	100 000	10 000
TinyImageNet + <i>ForAug</i>	200	99 404	9915
ImageNet	1000	1 281 167	50 000
ImageNet + <i>ForAug</i>	1000	1 274 557	49 751

335 strategies perform comparably, while *ViT-S* benefits from the  
 336 added diversity of using any background. The same is true  
 337 when training on the full ImageNet.

338 **Foreground position.** Finally, we analyze the foreground  
 339 object’s positioning in the image, using a generalization of  
 340 the Bates distribution [1] with parameter  $\eta \in \mathbb{Z}$ . The Bates  
 341 distribution presents an easy way to sample from a bounded  
 342 domain with just one hyperparameter that controls its concen-  
 343 tration.  $\eta = 1/-1$  corresponds to the uniform distribution;  
 344  $\eta > 1$  concentrates the distribution around the center; and  
 345 for  $\eta < -1$ , the distribution is concentrated at the borders  
 346 (see supplementary material for details). When sampling  
 347 more towards the center of the image, the difficulty of the  
 348 task is reduced, which reduces performance on TinyIma-  
 349 geNet (Table 3). This is reflected in the performance when  
 350 evaluating using *ForAug* with  $\eta = 2$  and  $\eta = 3$  compared to  
 351  $\eta = -1/1$ . We observe a similar reduction for  $\eta < -1$ .

352 After fixing the optimal design parameters in Tables 1  
 353 and 2 (last rows), we run *ForAug*’s segmentation step on the  
 354 entire ImageNet dataset. Table 4 shows the resulting dataset  
 355 statistics. The slightly reduced image count for *ForAug* is  
 356 due to instances where Grounded SAM fails to produce valid  
 357 segmentation masks.

## 4.2. Image Classification Results

358 Table 5 compares the ImageNet performance of models  
 359 trained with and without *ForAug*. We adopt the training  
 360 setup of [31] and [46] for training ViT [9], Swin [29] and  
 361 ResNet [15] (representing CNNs) models as well as the setup  
 362 of DeiT [45] for that model. Both setups are using strong  
 363

Table 5. ImageNet results of models trained on ImageNet with and without *ForAug*. *ForAug* improves the performance of most models, with a larger gain for larger models.

Model	ImageNet Accuracy [%]		Delta
	w/o <i>ForAug</i>	w/ <i>ForAug</i>	
ViT-S	79.1 ± 0.1	81.4 ± 0.1	+2.3
ViT-B	77.6 ± 0.2	81.1 ± 0.4	+3.5
ViT-L	75.3 ± 0.4	79.8 ± 0.1	+4.5
DeiT-S	80.1 ± 0.1	80.0 ± 0.3	-0.1
DeiT-B	81.9 ± 0.3	81.9 ± 0.2	±0.0
DeiT-L	79.3 ± 2.3	82.4 ± 0.1	+3.1
Swin-Ti	77.9 ± 0.2	79.7 ± 0.1	+1.8
Swin-S	79.4 ± 0.1	80.6 ± 0.1	+1.2
ResNet-50	78.3 ± 0.1	78.8 ± 0.1	+0.5
ResNet-101	79.4 ± 0.1	80.4 ± 0.1	+1.0

Table 6. Comparison of *ForAug* and simple Copy-Paste methods. We train ViT-S on ImageNet using the same 3-augment data augmentation on top of the copy-paste augmentation.

Augmentation	labels	Accuracy [%]	Delta to Prev.
Baseline + <b>Simple Copy-Paste</b>	bg	31.3 ± 0.6	
+ mixed labels	fg + bg	32.0 ± 0.8	+0.7
+ fg labels	fg	31.6 ± 0.9	-0.4
+ <i>range</i> foreground size variation	fg	43.0 ± 1.2	+11.4
+ infilled backgrounds	fg	68.7 ± 0.2	+25.7
+ <i>cos</i> mixing strategy	fg	81.2 ± 0.1	+12.5
+ edge smoothing	fg	81.3 ± 0.1	+0.1
+ background pruning= <i>ForAug</i>	fg	81.4 ± 0.1	+0.1

364 data augmentations like RandAugment, CutMix, and Mixup  
 365 optimized for Transformers (details in supplementary  
 366 material). Notably, *ForAug* improves performance across all  
 367 tested architectures, including the ResNet models, demon-  
 368 strating benefits beyond Transformers. For DeiT we only  
 369 observe benefits on ImageNet for the larger models. For  
 370 other transformers, we observe improvements from 1.2 p.p.  
 371 to 4.5 p.p. with increasing gains for larger models. *ForAug*'s  
 372 improvements counteract the drop in performance for in-  
 373 creasing model sizes. Without *ForAug* this drop is 3.8 p.p.  
 374 (ViT-S to L), while with *ForAug* it is reduced to 1.6 p.p. For  
 375 DeiT there is a drop of 0.8 p.p. from small to large while  
 376 when using *ForAug* there is a *gain* of 2.4 p.p.

377 **Comparison to Simple Copy-Paste.** We compare  
 378 *ForAug* to a simple adaption of the Copy-Paste augmentation  
 379 inspired by [11, 13, 39] in Table 6. Contrary to semantic seg-  
 380 mentation we do not have foreground masks available. Thus,  
 381 we paste the extracted foreground objects from *ForAug*'s *seg-*  
 382 *mentation stage* onto normal ImageNet images. We observe  
 383 3 large jumps in accuracy: (1) From our *range* foreground  
 384 size variation (+11.4%), (2) from using our infilled back-

Table 7. Downstream accuracy in percent when finetuning on other datasets. Models are pretrained on ImageNet with and without *ForAug*. Pretraining using *ForAug* increases transformer downstream accuracy.

Model	<i>ForAug</i>	Aircraft	Cars	Flowers	Food	Pets
ViT-S	✗	72.4 ± 1.0	89.8 ± 0.3	94.5 ± 0.2	89.1 ± 0.1	93.8 ± 0.2
ViT-S	✓	78.6 ± 0.5	92.2 ± 0.2	95.5 ± 0.2	89.6 ± 0.1	94.5 ± 0.2
		+6.2	+2.4	+1.0	+0.5	+0.7
ViT-B	✗	71.7 ± 0.5	90.0 ± 0.2	94.8 ± 0.4	89.8 ± 0.2	94.1 ± 0.4
ViT-B	✓	79.0 ± 2.2	93.3 ± 0.1	96.5 ± 0.1	90.9 ± 0.1	95.1 ± 0.4
		+7.3	+3.3	+1.7	+1.1	+1.0
ViT-L	✗	72.1 ± 1.0	88.8 ± 0.3	94.4 ± 0.3	90.1 ± 0.2	94.2 ± 0.4
ViT-L	✓	77.6 ± 1.2	89.1 ± 0.2	96.6 ± 0.1	91.3 ± 0.1	95.1 ± 0.1
		+5.5	+0.3	+2.2	+1.2	+0.9
DeiT-S	✗	75.3 ± 0.4	91.1 ± 0.2	94.8 ± 0.4	89.2 ± 0.2	92.4 ± 0.2
DeiT-S	✓	76.8 ± 0.8	91.9 ± 0.2	95.2 ± 0.3	89.1 ± 0.2	92.3 ± 0.4
		+1.5	+0.8	+0.4	-0.1	-0.1
DeiT-B	✗	77.0 ± 1.2	92.9 ± 0.2	96.1 ± 0.2	91.2 ± 0.1	93.3 ± 0.4
DeiT-B	✓	79.3 ± 0.3	93.1 ± 0.1	96.4 ± 0.2	91.3 ± 0.1	93.3 ± 0.1
		+2.3	+0.2	+0.3	+0.1	±0.0
DeiT-L	✗	72.8 ± 5.5	92.8 ± 1.0	95.8 ± 1.5	90.5 ± 2.6	92.4 ± 2.0
DeiT-L	✓	78.8 ± 0.8	93.8 ± 0.2	97.0 ± 0.2	92.0 ± 0.2	93.5 ± 0.2
		+6.0	+1.0	+1.2	+1.5	+1.1
Swin-Ti	✗	77.0 ± 0.1	91.3 ± 0.6	95.9 ± 0.1	90.0 ± 0.2	94.2 ± 0.1
Swin-Ti	✓	81.1 ± 0.8	92.8 ± 0.4	96.2 ± 0.1	90.4 ± 0.3	94.8 ± 0.5
		+4.1	+2.5	+0.3	+0.4	+0.6
Swin-S	✗	75.7 ± 1.4	91.0 ± 0.3	95.9 ± 0.5	91.1 ± 0.2	94.4 ± 0.1
Swin-S	✓	81.4 ± 0.2	93.1 ± 0.2	96.3 ± 0.3	91.2 ± 0.2	94.9 ± 0.3
		+5.7	+2.1	+1.4	+0.1	+0.5
ResNet-50	✗	78.2 ± 0.5	89.8 ± 0.2	91.7 ± 0.4	84.4 ± 0.2	93.7 ± 0.3
ResNet-50	✓	80.3 ± 0.4	90.4 ± 0.2	91.7 ± 0.2	84.5 ± 0.2	93.7 ± 0.3
		+2.1	+0.6	±0.0	+0.1	±0.0
ResNet-101	✗	78.4 ± 0.6	90.3 ± 0.1	91.2 ± 0.5	86.0 ± 0.2	94.3 ± 0.2
ResNet-101	✓	81.4 ± 0.5	91.3 ± 0.1	92.9 ± 0.2	86.3 ± 0.1	94.0 ± 0.3
		+3.0	+1.3	+1.7	+0.3	-0.3

grounds instead of images from the dataset (+25.7%), and  
 (3) from our *cos* mixing strategy with non-augmented im-  
 ages (+12.5%). *ForAug*'s changes to the naive copy-paste  
 augmentation are thus imperative for good classification per-  
 formance.

**Downstream tasks.** To assess the transferability of  
*ForAug*-trained models, we finetune models pretrained on  
ImageNet with and without *ForAug* on five fine-grained  
datasets: FGVC-Aircraft [30], Stanford Cars [6], Oxford  
Flowers [32], Food-101 [20], and Oxford-IIIT Pets [33]. In  
Table 7 we see transformer accuracies improve on all these  
datasets by up to 7.3 p.p. Notably, training with *ForAug*  
boosts the downstream performance of DeiT-S and DeiT-B,  
despite similar ImageNet results. This shows the improved  
representations from training with *ForAug* translate to gains  
beyond better ImageNet scores.

### 4.3. Bias and Robustness Evaluation

Beyond its use for training, *ForAug*'s unique properties and  
controlled data generation capabilities make it a powerful  
tool for analyzing behavior and biases of black-box models.

**Background Robustness.** We assess the robustness of  
models to shifts in the background distribution from a class-  
related background to any background. Figure 3 presents  
the background robustness results for three datasets: Im-

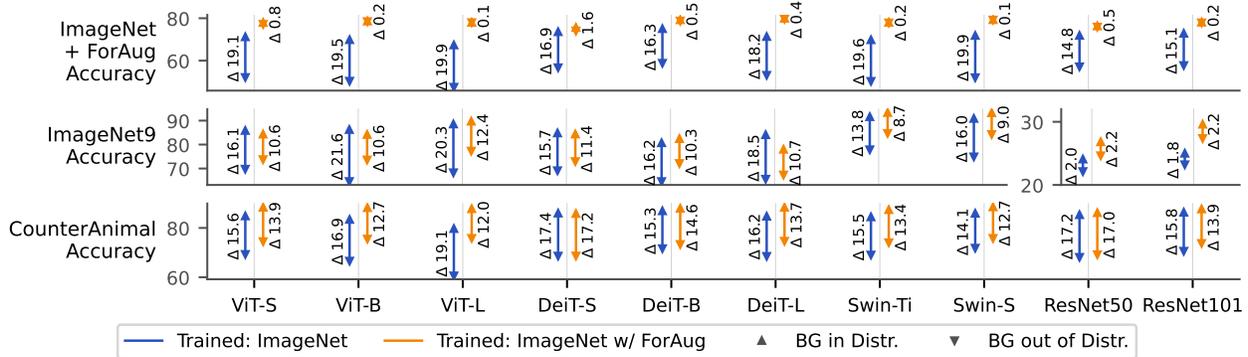


Figure 3. Evaluation of background robustness on ImageNet + *ForAug*, ImageNet9 and CounterAnimal. We plot the in-distribution (top of arrow) and the out-of-distribution (bottom of arrow) accuracy when training with and without *ForAug*. We annotate each arrow with its length  $\Delta$ . Training with *ForAug* improves the background robustness of all transformers by mostly boosting the out-of-distribution accuracy.

409 ageNet with *ForAug* (all backgrounds vs. backgrounds  
410 of same class), ImageNet9 [54] (random backgrounds vs.  
411 original backgrounds), and CounterAnimal [49] (counter  
412 vs. common background). The top triangle of each arrow  
413 represents the in-distribution backgrounds and the bot-  
414 tom triangle represents the out-of-distribution ones. We  
415 follow ImageNet9 and CounterAnimal and assess the back-  
416 ground robustness in terms of the accuracy gap when evalu-  
417 ating a model on images of normal background distribu-  
418 tion compared to out-of-distribution backgrounds (length  
419 of each arrow;  $\Delta$ ). Crucially, *ForAug* improves the back-  
420 ground robustness of all models and across datasets, reduc-  
421 ing the background-gap by boosting the performance on  
422 the out-of-background-distribution samples more than the in-  
423 distribution ones. These findings highlight the generalization  
424 benefits of *ForAug* to unusual image compositions.

425 **Foreground Focus.** Leveraging our inherent knowledge  
426 of the foreground masks when using *ForAug*, as well as com-  
427 mon XAI techniques [3, 38, 42], we can evaluate a model’s  
428 focus on the foreground object. We can directly evaluate  
429 ImageNet-trained models, but this technique can also be  
430 extended to other datasets without relying on manually an-  
431 notated foreground masks. To evaluate the foreground focus,  
432 we employ Grad-CAM [38], Grad-CAM++ [3] and Integrat-  
433 edGradients (IG) [42] to compute the per-pixel importance  
434 of an image for the model’s prediction. The foreground  
435 focus is defined to be the ratio of the foreground’s relative  
436 importance to its relative size in the image:

$$437 \quad \text{FG Focus}(\text{img}) = \frac{\text{Area}(\text{img}) \text{ Importance}(\text{fg})}{\text{Area}(\text{fg}) \text{ Importance}(\text{img})} \quad (2)$$

438 If all pixels uniformly receive the same importance value,  
439 the foreground focus is one. The foreground focus of a  
440 model is its average focus over all test images. Figure 4  
441 presents our findings. Using *ForAug* significantly increases  
442 the foreground focus of ViT, DeiT and ResNet across all XAI

Table 8. Accuracy relative to the center accuracy of multiple instan-  
tiations of the models when the foreground objects is in different  
cells of a  $3 \times 3$  grid. We calculate center bias according to Equa-  
tion (3). Using *ForAug* significantly reduces models’ center bias.

Model	Center Bias [%] when trained		Delta
	w/o <i>ForAug</i>	w/ <i>ForAug</i>	
ViT-S	25.5 ± 0.8	22.0 ± 0.3	-3.5
ViT-B	25.4 ± 0.4	19.0 ± 0.2	-6.4
ViT-L	24.3 ± 1.1	11.7 ± 0.7	-12.6
DeiT-S	20.4 ± 0.2	21.2 ± 0.1	+0.8
DeiT-B	19.0 ± 0.7	19.0 ± 0.2	±0.0
DeiT-L	21.2 ± 0.2	18.0 ± 0.2	-3.2
Swin-Ti	25.0 ± 0.7	16.5 ± 0.2	-8.5
Swin-S	23.2 ± 0.1	15.6 ± 0.2	-7.6
ResNet50	26.3 ± 0.3	19.7 ± 0.3	-6.6
ResNet101	23.0 ± 0.3	19.9 ± 0.2	-3.1

0.65 0.70 0.75 0.80 0.85 0.90 0.95 1.00

443 metrics. We hypothesize Swin’s below-uniform foreground  
444 focus with GradCam is due to its specific implementation.

445 **Center Bias.** With *ForAug* we have unique control over  
446 the position of the foreground object in the image. This  
447 lets us quantify the center bias of models trained with and

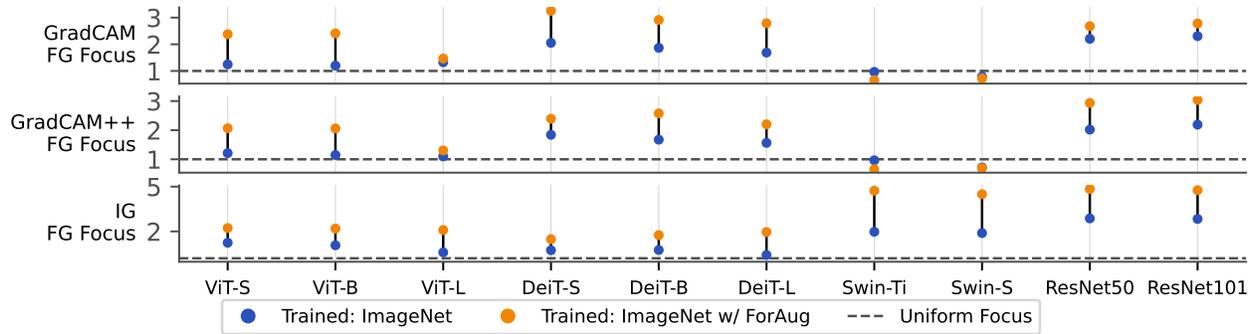


Figure 4. Evaluation of the foreground focus (Equation (2)) using GradCam, GradCam++ and IntegratedGradients (IG) of models trained on ImageNet. Training with *ForAug* improves the foreground focus of almost all models.

448 without *ForAug*. We divide the image into a  $3 \times 3$  grid and  
 449 evaluate model accuracy when the (scaled-down) foreground  
 450 object is in each of the 9 grid cells. Each cell’s accuracy  
 451 is divided by the accuracy in the center cell for normaliza-  
 452 tion, which gives us the relative performance drop when the  
 453 foreground is in each part of the image. The center bias is  
 454 calculated as one minus the average of the minimum perfor-  
 455 mance of a corner cell and the minimum performance of a  
 456 side cell:

$$457 \text{ Center Bias} = 1 - \frac{\min_{c \in \text{sides}} \text{Acc}(c) + \min_{c \in \text{corners}} \text{Acc}(c)}{2\text{Acc}(c_{\text{center}})} \quad (3)$$

458 Table 8 visualizes the center bias of three instantiations of  
 459 each model. Performance is generally highest in the center  
 460 and lowest in the four corners. Interestingly, ImageNet-  
 461 trained models perform slightly better when the foreground  
 462 object is on the right side of the image, compared to the left  
 463 side, despite our use of random flipping with a probability  
 464 of 0.5 during training. Using *ForAug* significantly reduces  
 465 center bias across models, with a more uniform performance  
 466 especially across the middle row. Thus, *ForAug* makes the  
 467 model recognize objects across a wider spatial distribution,  
 468 counteracting the center-bias of ImageNet.

469 **Size Bias.** Finally, we evaluate the impact of different  
 470 sized foreground objects on the accuracy. For this evaluation,  
 471 we use the *mean* foreground size strategy. We introduce  
 472 a size factor  $f_{\text{size}}$  by which we additionally scale the fore-  
 473 ground object before pasting it onto the background. Results  
 474 are normalized by the accuracy when using  $f_{\text{size}} = 1.0$ . Fig-  
 475 ure 5 shows the size bias curves of models trained with  
 476 and without *ForAug*. Models trained using *ForAug* main-  
 477 tain perform better, especially with smaller foreground objects.  
 478 Therefore, *ForAug*-training improves robustness to variations  
 479 in object scale, especially for larger models.

## 480 5. Discussion & Conclusion

481 We introduce *ForAug*, a novel data augmentation scheme  
 482 that facilitates improved Transformer training for image clas-

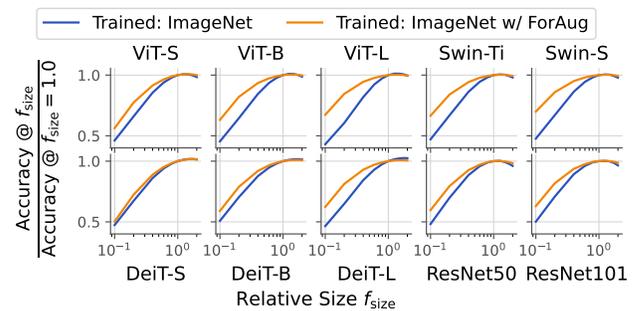


Figure 5. Evaluation of the size bias of models trained on ImageNet. We plot the accuracy relative to the accuracy when using the default size ( $f_{\text{size}} = 1.0$ ).

483 sification. By explicitly separating and recombining fore-  
 484 ground objects and backgrounds, *ForAug* enables controlled  
 485 data augmentation beyond existing image compositions,  
 486 leading to significant performance gains on ImageNet and  
 487 downstream fine-grained classification tasks. Furthermore,  
 488 *ForAug* provides a powerful framework for analyzing model  
 489 behavior and quantifying biases, including background ro-  
 490 bustness, foreground focus, center bias, and size bias. Our  
 491 experiments demonstrate that training using *ForAug* not only  
 492 boosts accuracy but also significantly reduces these biases,  
 493 resulting in more robust and generalizable models. In the  
 494 future, we see *ForAug* be also applied to other datasets and  
 495 tasks, like video recognition or segmentation. *ForAug*’s abil-  
 496 ity to both improve performance and provide insights into  
 497 model behavior makes it a valuable tool for advancing CV  
 498 research and developing more reliable AI systems.

## 499 Acknowledgements

500 Will be in the final paper.

## 501 References

- 502 [1] G.E. Bates. Joint distributions of time intervals for the oc-  
 503 currence of successive accidents in a generalized polya urn

- 504 scheme. *Annals of Mathematical Statistics*, 26:705–720, 1955. 5
- 505
- 506 [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas
- 507 Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-
- 508 end object detection with transformers. 2020. 1
- 509 [3] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and
- 510 Vineeth N Balasubramanian. Grad-cam++: Generalized
- 511 gradient-based visual explanations for deep convolutional
- 512 networks. In *2018 IEEE Winter Conference on Applications*
- 513 *of Computer Vision (WACV)*, pages 839–847, 2018. 7
- 514 [4] Ekin D. Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasude-
- 515 van, and Quoc V. Le. Autoaugment: Learning augmentation
- 516 policies from data. 2018. 2
- 517 [5] Ekin D. Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V.
- 518 Le. Randaugment: Practical automated data augmentation
- 519 with a reduced search space. 2019. 2
- 520 [6] Afshin Dehghan, Syed Zain Masood, Guang Shu, and En-
- 521 rique. G. Ortiz. View independent vehicle make, model and
- 522 color recognition using convolutional neural network. 2017.
- 523 6
- 524 [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li
- 525 Fei-Fei. ImageNet: A large-scale hierarchical image database.
- 526 In *2009 IEEE Conference on Computer Vision and Pattern*
- 527 *Recognition*. IEEE, 2009. 1
- 528 [8] Peijian Ding, Davit Soselia, Thomas Armstrong, Jiahao Su,
- 529 and Furong Huang. Reviving shift equivariance in vision
- 530 transformers. 2023. 1
- 531 [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov,
- 532 Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
- 533 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl-
- 534 vain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is
- 535 worth 16x16 words: Transformers for image recognition at
- 536 scale. In *9th International Conference on Learning Represen-*
- 537 *tations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*.
- 538 OpenReview.net, 2021. 1, 5
- 539 [10] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut,
- 540 paste and learn: Surprisingly easy synthesis for instance de-
- 541 tection. 2017. 2
- 542 [11] Yunhao Ge, Jiashu Xu, Brian Nlong Zhao, Neel Joshi, Laurent
- 543 Itti, and Vibhav Vineet. Beyond generation: Harnessing text
- 544 to image models for object detection and segmentation. *ArXiv*,
- 545 abs/2309.05956, 2023. 2, 6
- 546 [12] Robert Geirhos, Patricia Rubisch, Claudio Michaelis,
- 547 Matthias Bethge, Felix A. Wichmann, and Wieland Brendel.
- 548 Imagenet-trained cnns are biased towards texture; increasing
- 549 shape bias improves accuracy and robustness. 2018. 3
- 550 [13] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-
- 551 Yi Lin, Ekin D. Cubuk, Quoc V. Le, and Barret Zoph. Simple
- 552 copy-paste is a strong data augmentation method for instance
- 553 segmentation. 2020. 2, 6
- 554 [14] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra
- 555 Malik. Rich feature hierarchies for accurate object detection
- 556 and semantic segmentation. 2013. 1
- 557 [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.
- 558 Deep residual learning for image recognition. In *Proceed-*
- 559 *ings of the IEEE conference on computer vision and pattern*
- 560 *recognition*, pages 770–778, 2016. 1, 5
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Gir- 561  
shick. Mask r-cnn. 2017. 1 562
- [17] Dan Hendrycks and Thomas Dietterich. Benchmarking neural 563  
network robustness to common corruptions and perturbations. 564  
2019. 2 565
- [18] Stefan Hinterstoisser, Olivier Pauly, Hauke Heibel, Marek 566  
Martina, and Martin Bokeloh. An annotation saved is an an- 567  
notation earned: Using fully synthetic training for object de- 568  
tection. In *2019 IEEE/CVF International Conference on Com- 569  
puter Vision Workshop (ICCVW)*, pages 2787–2796, 2019. 2 570
- [19] Ji-Soo Kang and Kyungyong Chung. Staug: Copy-paste 571  
based image augmentation technique using salient target. 10: 572  
123605–123613. 2 573
- [20] Parneet Kaur, Karan Sikka, and Ajay Divakaran. Combining 574  
weakly and weakly supervised learning for classifying food 575  
images. 2017. 6 576
- [21] Salman Khan, Muzammal Naseer, Munawar Hayat, 577  
Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. 578  
Transformers in vision: A survey. *ACM Computing Surveys*, 579  
54(10s):1–41, 2022. 1 580
- [22] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, 581  
Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer White- 582  
head, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross 583  
Girshick. Segment anything. 2023. 3 584
- [23] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Im- 585  
agenet classification with deep convolutional neural networks. 586  
In *Advances in Neural Information Processing Systems*. Cur- 587  
ran Associates, Inc., 2012. 1 588
- [24] Yann Le and Xuan Yang. Tiny imagenet visual recognition 589  
challenge. *CS 231N*, 7(7):3, 2015. 4 590
- [25] Xiaodan Li, Yuefeng Chen, Yao Zhu, Shuhui Wang, Rong 591  
Zhang, and Hui Xue. Imagenet-e: Benchmarking neural 592  
network robustness via attribute editing. 2023. 2 593
- [26] Evan Ling, Dezhao Huang, and Minhoe Hur. Humans need 594  
not label more humans: Occlusion copy & paste for occluded 595  
human instance segmentation. 2022. 2 596
- [27] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao 597  
Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, 598  
Hang Su, Jun Zhu, and Lei Zhang. Grounding dino: Marrying 599  
dino with grounded pre-training for open-set object detection. 600  
2023. 3 601
- [28] Yue Liu, Christos Matsoukas, Fredrik Strand, Hossein Az- 602  
izpour, and Kevin Smith. Patchdropout: Economizing vision 603  
transformers using patch dropout. 2022. 2 604
- [29] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng 605  
Zhang, Stephen Lin, and Baining Guo. Swin transformer: 606  
Hierarchical vision transformer using shifted windows. In 607  
*2021 IEEE/CVF International Conference on Computer Vi-* 608  
*sion (ICCV)*, pages 9992–10002, Los Alamitos, CA, USA, 609  
2021. IEEE Computer Society. 1, 5 610
- [30] S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. 611  
Fine-grained visual classification of aircraft. Technical report, 612  
2013. 6 613
- [31] Tobias Christian Nauen, Sebastian Palacio, and Andreas Dengel. 614  
Which transformer to favor: A comparative analysis of 615  
efficiency in vision transformers. In *Proceedings of the Win-* 616  
*ter Conference on Applications of Computer Vision (WACV)*, 617  
pages 6955–6966, 2025. 5 618

- 619 [32] Maria-Elena Nilsback and Andrew Zisserman. Automated  
620 flower classification over a large number of classes. In *Indian Conference on Computer Vision, Graphics and Image*  
621 *Processing*, 2008. 6 677
- 623 [33] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and  
624 C. V. Jawahar. Cats and dogs. In *IEEE Conference on Com-*  
625 *puter Vision and Pattern Recognition*, 2012. 6 678
- 626 [34] Gabriela Rangel, Juan C. Cuevas-Tello, Jose Nunez-Varela,  
627 Cesar Puente, and Alejandra G. Silva-Trujillo. A survey on  
628 convolutional neural networks and their performance limita-  
629 tions in image recognition tasks. 2024(1), 2024. 1 679
- 630 [35] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li,  
631 He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan,  
632 Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang  
633 Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling  
634 open-world models for diverse visual tasks. 2024. 2, 3 680
- 635 [36] Renan A. Rojas-Gomez, Teck-Yian Lim, Minh N. Do, and  
636 Raymond A. Yeh. Making vision transformers truly shift-  
637 equivariant. 2023. 1 681
- 638 [37] Edward Sanderson and Bogdan J. Matuszewski. *FCN-*  
639 *Transformer Feature Fusion for Polyp Segmentation*, pages  
640 892–907. Springer International Publishing, 2022. 1 682
- 641 [38] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das,  
642 Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-  
643 cam: Visual explanations from deep networks via gradient-  
644 based localization. 128(2):336–359, 2016. 7 683
- 645 [39] Ang Jia Ning Shermaine, Michalis Lazarou, and Tania  
646 Stathaki. Image compositing is all you need for data aug-  
647 mentation. 2025. 2, 6 684
- 648 [40] Connor Shorten and Taghi M. Khoshgoftaar. A survey on  
649 image data augmentation for deep learning. 6(1), 2019. 1, 2 685
- 650 [41] Wenhao Sun, Benlei Cui, Xue-Mei Dong, and Jingqun Tang.  
651 Attentive eraser: Unleashing diffusion model’s object removal  
652 potential via self-attention redirection guidance. 2024. 2, 3, 4 686
- 653 [42] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic  
654 attribution for deep networks. 2017. 7 687
- 655 [43] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin,  
656 Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov,  
657 Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lem-  
658 pitsky. Resolution-robust large mask inpainting with fourier  
659 convolutions. 2021. 2, 3, 4, 5 688
- 660 [44] Ryo Takahashi, Takashi Matsubara, and Kuniaki Uehara. Data  
661 augmentation using random image cropping and patching for  
662 deep cnns. 30(9):2917–2931, 2018. 2 689
- 663 [45] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco  
664 Massa, Alexandre Sablayrolles, and Herve Jegou. Training  
665 data-efficient image transformers & distillation through atten-  
666 tion. In *Proceedings of the 38th International Conference on*  
667 *Machine Learning*, pages 10347–10357. PMLR, 2021. 1, 5 690
- 668 [46] Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii:  
669 Revenge of the vit. In *Computer Vision – ECCV 2022*, pages  
670 516–533, Cham, 2022. Springer Nature Switzerland. 1, 2, 4,  
671 5 691
- 672 [47] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszko-  
673 reit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia  
674 Polosukhin. Attention is all you need. In *Advances in Neu-*  
675 *ral Information Processing Systems*. Curran Associates, Inc.,  
676 2017. 1 692
- [48] Ioannis A. Vezakis, Konstantinos Georgas, Dimitrios Fotiadis,  
and George K. Matsopoulos. Effisegnet: Gastrointestinal  
polyp segmentation through a pre-trained efficientnet-based  
network with a simplified decoder. 2024. 1 693
- [49] Qizhou Wang, Yong Lin, Yongqiang Chen, Ludwig Schmidt,  
Bo Han, and Tong Zhang. A sober look at the robustness  
of CLIPs to spurious features. In *The Thirty-eighth Annual*  
*Conference on Neural Information Processing Systems*, 2024.  
7 694
- [50] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang  
Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed,  
Saksham Singhal, Subhojit Som, and Furu Wei. Image as a  
foreign language: Beit pretraining for all vision and vision-  
language tasks. 2022. 1 695
- [51] Wenhui Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi  
Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng  
Li, Xiaogang Wang, and Yu Qiao. Internimage: Exploring  
large-scale vision foundation models with deformable convo-  
lutions. 2022. 1 696
- [52] Levi Kassel Michael Werman. Deepaste – inpainting for  
pasting. 2021. 2 697
- [53] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Re-  
becca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos,  
Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon  
Kornblith, and Ludwig Schmidt. Model soups: averaging  
weights of multiple fine-tuned models improves accuracy  
without increasing inference time. In *Proceedings of the*  
*39th International Conference on Machine Learning*, pages  
23965–23998. PMLR, 2022. 1 698
- [54] Kai Xiao, Logan Engstrom, Andrew Ilyas, and Aleksander  
Madry. Noise or signal: The role of image backgrounds in  
object recognition. 2020. 3, 7 699
- [55] Mingle Xu, Sook Yoon, Alvaro Fuentes, and Dong Sun Park.  
A comprehensive survey of image augmentation techniques  
for deep learning. 137:109347, 2023. 1, 2 700
- [56] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mo-  
jtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive  
captioners are image-text foundation models. *Transactions*  
*on Machine Learning Research*, 2022. 1 701
- [57] Sangdoon Yun, Dongyoon Han, Sanghyuk Chun, Seong Joon  
Oh, Youngjoon Yoo, and Junsuk Choe. CutMix: Regular-  
ization strategy to train strong classifiers with localizable  
features. In *2019 IEEE/CVF International Conference on*  
*Computer Vision (ICCV)*. IEEE, 2019. 2 702
- [58] Chenshuang Zhang, Fei Pan, Junmo Kim, In So Kweon, and  
Chengzhi Mao. Imagenet-d: Benchmarking neural network  
robustness on diffusion synthetic object. 2024. 2 703
- [59] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and  
David Lopez-Paz. mixup: Beyond empirical risk minimiza-  
tion. In *International Conference on Learning Representa-*  
*tions*, 2018. 2 704
- [60] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and  
Yi Yang. Random erasing data augmentation. 2017. 2 705
- [61] Zhuofan Zong, Guanglu Song, and Yu Liu. Detsr with collab-  
orative hybrid assignments training. 2022. 1 706