

ForAug: Recombining Foregrounds and Backgrounds to Improve Vision Transformer Training with Bias Mitigation

Anonymous ICCV submission

Paper ID 6426

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 that limit their robustness
 005 and generalizability. This paper introduces ForAug, a novel
 006 data augmentation scheme that addresses these challenges
 007 and explicitly includes inductive biases, which commonly
 008 are part of the neural network architecture, into the training
 009 data. ForAug is constructed by using pretrained founda-
 010 tion models to separate and recombine foreground objects
 011 with different backgrounds, enabling fine-grained control
 012 over image composition during training. It thus increases
 013 the data diversity and effective number of training samples.
 014 We demonstrate that training on ForNet, the application of
 015 ForAug to ImageNet, significantly improves the accuracy of
 016 ViTs and other architectures by up to 4.5 percentage points
 017 (p.p.) on ImageNet and 7.3 p.p. on downstream tasks. Im-
 018 portantly, ForAug enables novel ways of analyzing model
 019 behavior and quantifying biases. Namely, we introduce met-
 020 rics for background robustness, foreground focus, center
 021 bias, and size bias and show that training on ForNet substan-
 022 tially reduces these biases compared to training on ImageNet.
 023 In summary, ForAug provides a valuable tool for analyzing
 024 and mitigating biases, enabling the development of more
 025 robust and reliable computer vision models. Our code and
 026 dataset are publicly available at [<url>](#).

027 1. Introduction

028 Image classification, a fundamental task in computer vi-
 029 sion (CV), involves assigning a label to an image from a
 030 predefined set of categories. This seemingly simple task
 031 underpins a wide range of applications, including medical di-
 032 agnosis [39, 50], autonomous driving [52], and object recog-
 033 nition [3, 15, 17]. Furthermore, image classification is used
 034 for large-scale pretraining of vision models [10, 31, 47] and
 035 to judge the progress of the field of CV [22, 36]. The advent

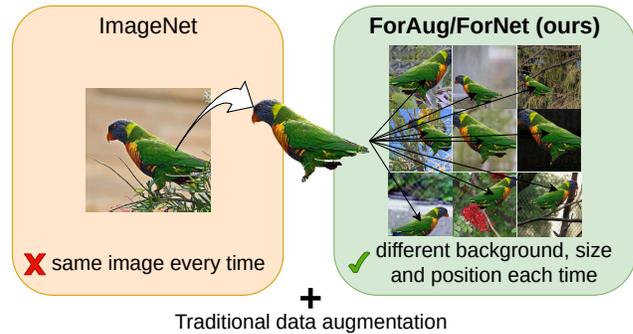


Figure 1. Comparison of ForNet and ImageNet. ForNet recombines foreground objects with different backgrounds each epoch, thus creating a more diverse training set. We still apply traditional data augmentation afterwards.

of large-scale datasets, particularly ImageNet [8], containing millions of labeled images across thousands of categories, has been instrumental in driving significant progress in this field. ImageNet served as a catalyst for the rise of large-scale CV models [16, 25] and remains the most important CV benchmark for more than a decade [16, 25, 48, 54].

While traditionally, convolutional neural networks (CNNs) have been the go-to architecture for image classification, Transformers [49], particularly the Vision Transformer (ViT) [10], have emerged as a powerful alternative. These attention-based models have demonstrated superior performance in various vision tasks, including image classification [3, 51, 54, 57, 62].

Data augmentation is a key technique for training image classification models. Traditional data augmentation methods, such as random cropping, flipping, and color jittering, are commonly employed to increase the diversity of the training data and improve the model’s performance [42, 56]. These basic transformations, originally designed for CNNs, change the input images in a way that preserves their semantic meaning [1]. However, the architectural differences of CNNs and Transformers suggest that the latter might benefit from different data augmentation strategies. In particular,

059 the Transformers self-attention mechanism is not translation
060 equivariant [9, 38], meaning that the model does not inher-
061 ently understand the spatial relationships between pixels.

062 Inspired by this inductive bias of CNNs, that is not inher-
063 ent to ViTs, we propose *ForAug*, a novel data augmentation
064 scheme for image classification which makes the transla-
065 tion equivariance of CNNs explicit in the training data by
066 recombining foreground objects at varying positions with
067 different backgrounds. Applying *ForAug* to ImageNet gives
068 rise to *ForNet*, a novel dataset that enables this data aug-
069 mentation with with fine-grained control over the image
070 composition. Recognizing that Transformers need to learn
071 the spatial relationships from data, since they are not inher-
072 ently translation invariant, and in general are usually trained
073 on larger datasets [24], we separate the foreground objects
074 in ImageNet from their backgrounds, using an open-world
075 object detector [37], and fill in the background in a plausible
076 way using an object removal model [43, 45]. This allows us
077 to recombine any foreground object with any background
078 on the fly, creating a highly diverse training set. During re-
079 combination, we can control important parameters, like the
080 size and position of the foreground object, to help the model
081 learn the spatial invariances necessary for image classifica-
082 tion. We show that training on *ForNet* instead of ImageNet
083 increases the model accuracy of Transformers by up to 4.5
084 p.p. on ImageNet and an up to 39.3% reduction in error rate
085 on downstream tasks.

086 Additionally, *ForAug* is a useful tool for analyzing model
087 behavior and biases, when used during the evaluation phase.
088 We utilize our control over the image distribution to quantify
089 a model’s background robustness (by varying the choice of
090 background), foreground focus (by leveraging our knowl-
091 edge about the placement of the foreground object), center
092 bias (by controlling the object’s position), and size bias (by
093 controlling object size). These analyses provide insights
094 into model behavior and biases, which is crucial for model
095 deployment and future robustness optimizations. We show
096 that training on *ForNet*, instead of ImageNet, significantly
097 reduces all of these biases, completely removing the models’
098 dependence on the background distribution. We make our
099 code for *ForAug* and the *ForNet*-dataset publicly available¹
100 to facilitate further research.

101 **Contributions**

- 102 • We propose *ForAug*, a novel data augmentation scheme,
103 that recombines objects and backgrounds to train Trans-
104 formers for image classification.
- 105 • We show that training on *ForNet*, the ImageNet instanti-
106 ation of *ForAug*, leads to 4.5 p.p. improved accuracy on
107 ImageNet and 7.3 p.p. on downstream tasks.
- 108 • We propose novel *ForAug*-based metrics to analyze and
109 quantify fine-grained biases trained models: Background

¹Link will go here.

Robustness, Foreground Focus, Center Bias, and Size Bias. 110
Training on *ForNet*, instead of ImageNet, significantly 111
reduces these biases. 112

2. Related Work 113

Data Augmentation for Image Classification Data aug- 114
mentation is a crucial technique for improving the perfor- 115
mance and generalization of image classification models. 116
Traditional augmentation strategies rely on simple geomet- 117
ric or color-space transformations like cropping, flipping, 118
rotation, blurring, color jittering, or random erasing [61] to 119
increase the diversity of the training data without changing 120
their semantic meaning. With the advent of Transformers, 121
new data augmentation operations like PatchDropout [30] 122
have been proposed. Other transformations like Mixup [60], 123
CutMix [58], or random cropping and patching [46] com- 124
bine multiple input images. These simple transformations 125
are usually bundled to form more complex augmentation 126
policies like AutoAugment [5] and RandAugment [6], which 127
automatically search for optimal augmentation policies or 128
3-augment [48] which is optimized to train a ViT. For a gen- 129
eral overview of data augmentation techniques for image 130
classification, we refer to [42, 56]. 131

We build upon these general augmentation techniques 132
by introducing a novel approach to explicitly separate and 133
recombine foregrounds and backgrounds for image classifi- 134
cation. Our approach is used in tandem with traditional data 135
augmentation techniques to improve model performance and 136
reduce biases. 137

Copy-Paste Augmentation The copy-paste augmentation 138
[14], which is used for object detection [14, 41] and instance 139
segmentation [28, 53], involves copying segmented objects 140
from one image and pasting them onto another. While typi- 141
cally human-annotated segmentation masks are used to ex- 142
tract the foreground objects, other foreground sources have 143
been explored, like 3D models [19] and pretrained object- 144
detection models for use on objects on white background 145
[11] or synthetic images [12]. DeePaste [53] focuses on us- 146
ing inpainting for a more seamless integration of the pasted 147
object. 148

Unlike these methods, *ForNet* focuses on image classifi- 149
cation. While for detection and segmentation, objects are 150
pasted onto another image (with a different foreground) or 151
on available or rendered background images of the target 152
scene, we extract foreground objects and fill in the resulting 153
holes in the background in a semantically neutral way. This 154
way, we can recombine any foreground object with a large 155
variety of neutral backgrounds from natural images, enabling 156
a controlled and diverse manipulation of image composition. 157

158 **Model robustness evaluation** Evaluating model robust- 209
 159 ness to various image variations is critical for understand- 210
 160 ing and improving model generalization. Datasets like 211
 161 ImageNet-C [18] and ImageNet-P [18] introduce common 212
 162 corruptions and perturbations. ImageNet-E [27] evaluates 213
 163 model robustness against a collection of distribution shifts. 214
 164 Other datasets, such as ImageNet-D [59], focus on varying 215
 165 background, texture, and material, but rely on synthetic data. 216
 166 Stylized ImageNet [13] investigates the impact of texture 217
 167 changes. ImageNet-9 [55] explores background variations 218
 168 using segmented images, but the backgrounds are often arti- 219
 169 ficial. 220

170 In contrast to these existing datasets, which are used only 221
 171 for evaluation, *ForNet* provides fine-grained control over 222
 172 foreground object placement, size, and background selection, 223
 173 enabling a precise and comprehensive analysis of specific 224
 174 model biases within the context of a large-scale, real-world 225
 175 image distribution. As *ForNet* also provides controllable 226
 176 training set generation, it goes beyond simply measuring 227
 177 robustness to actively improving it through training.

178 3. *ForAug* (Method)

179 We introduce *ForAug*, a data augmentation scheme designed 228
 180 to enhance Transformer training by explicitly separating and 229
 181 recombining foreground objects and backgrounds. *ForAug* 230
 182 involves two stages: Segmentation and Recombination, both 231
 183 visualized in Figure 2. 232

184 Segmentation

185 The segmentation stage isolates the foreground objects and 233
 186 their corresponding backgrounds. We then fill in the back- 234
 187 ground in a visually plausible way [43] using a pretrained 235
 188 object-removal model. This stage is computed once offline 236
 189 and the results are stored for the recombination stage. 237

190 First, foreground objects are detected and segmented 238
 191 from their backgrounds using a prompt-based segmen- 239
 192 tation model to exploit the classification datasets labels. 240
 193 We use the state-of-the-art Grounded SAM [37], which 241
 194 is based on Grounding DINO [29] and SAM [23]. The 242
 195 prompt we use is “a <class name>, a type of 243
 196 <object category>”, where <class name> is the 244
 197 specific name of the objects class as defined by the dataset 245
 198 and <object category> is a the broader category of 246
 199 the object. The <object category> guides the segmen- 247
 200 tation model towards the correct object in case the <class 248
 201 name> alone is too specific. This can be the case with 249
 202 prompts like “sorrel” or “guenon”, where the more general 250
 203 name “horse” or “monkey” is more helpful. We derive the 251
 204 <object category> from the WordNet hierarchy, us- 252
 205 ing the immediate hypernym. 253

206 We iteratively extract up to n foreground masks for 254
 207 each dataset-image, using different more and more general 255
 208 prompts based on the more general synsets of WordNet (e.g.

209 “a sorrel, a type of horse”, “a horse, a type of equine”, ...). 210
 211 Masks that are very similar, with a pairwise IoU of at least 212
 213 0.9, are merged. The output is a set of masks delineating 214
 215 the foreground objects and the backgrounds. We select the 216
 217 best mask per image (according to Equation (1)) in a later 218
 219 filtering step, described below. 220

221 An inpainting model that is specifically optimized to re- 222
 223 move objects from images, such as LaMa [45] or Attentive 224
 225 Eraser [43], is used to inpaint the foreground regions in the 226
 227 backgrounds. To ensure the quality of the foreground and 228
 229 background images (for each dataset-image), we select a 230
 231 foreground/background pair from the $\leq n$ variants we have 232
 233 extracted and infilled in the previous steps. Using an ensem- 234
 235 ble of six ViT, ResNet, and Swin Transformer models pre- 236
 237 trained on the original dataset, we select the foreground/back- 238
 239 ground pair that maximizes foreground performance while 240
 241 minimizing the performance on the background and size of 242
 243 the foreground according to: 244

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

245 Here, E is the ensemble of models and m is a pretrained 246
 247 model, c is the correct foreground class, fg, and bg are the 248
 249 foreground and background and $\text{size}(\cdot)$ is the size in number 249
 250 of pixels. We ran a hyperparameter search using a manually 251
 252 annotated subset of foreground/background variants to find 253
 254 the factors in Equation (1): $\lambda = 2$ and $\varepsilon = 0.1$. The *opti-* 255
 256 *mal foreground size* of 10% of the full image balances the 256
 257 smallest possible foreground size that encompasses all the 257
 258 respective class information in the image with still convey- 258
 259 ing the foreground information after pasting it onto another 259
 260 background. This filtering step ensures we segment all the 260
 261 relevant foreground objects. 261

262 Finally, we filter out backgrounds that are more than 80% 262
 263 infilled, as these tend to be overly synthetic, plain and don’t 263
 264 carry much information (see the supplementary material). 264
 265 We ablate this choice in Section 4.1. In summary, we factor- 265
 266 ize the dataset into a set of foreground objects with a 266
 267 transparent background and a set of diverse backgrounds per 267
 268 class. The next step is to recombine them as data augmenta- 268
 269 tion before applying common data augmentation operations 269
 270 during training. 270

249 Recombination

250 The recombination stage, which is performed online, com- 251
 252 bines the foreground objects with different backgrounds to 252
 253 create new training samples. For each object, we follow the 253
 254 pipeline of: Pick an appropriate background, resize it to a 254

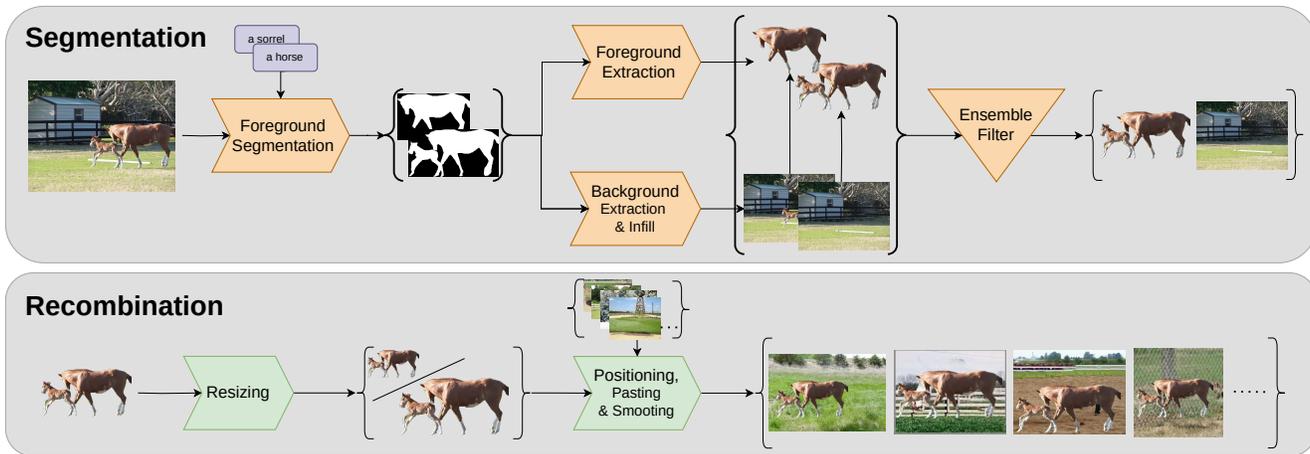


Figure 2. Overview of *ForNet*. The data creation consists of two stages: (1, offline) Segmentation, where we segment the foreground objects from the background and fill in the background. (2, online) Recombination, where we combine the foreground objects with different backgrounds to create new samples.

254 fitting size, place it in the background image, smooth the
255 transition edge, and apply other data augmentations.

256 For each foreground object, we sample a background using
257 one of the following strategies: (1) the original image
258 background, (2) the set of backgrounds from the same class,
259 or (3) the set of all possible backgrounds. These sets are trad-
260 ing off the amount of information the model can learn from
261 the background against the diversity of new images created.
262 In each epoch, each foreground object is seen exactly once,
263 but a background may appear multiple times.

264 The selected foreground is resized based on its relative
265 size within its original image and the relative size of the
266 original foreground in the selected background image. The
267 final size is randomly selected from a 30% range around
268 upper and lower limits (s_u and s_l), based on the original
269 sizes:

270
$$s \sim \mathcal{U}[(1 - 0.3)s_l, (1 + 0.3)s_u]. \quad (2)$$

271 To balance the size of the foreground and that of the back-
272 grounds original foreground, the upper and lower limit s_u
273 and s_l are set to the mean or range of both sizes, depending
274 on the foreground size strategy: *mean* or *range*.

275 The resized foreground is then placed at a random posi-
276 tion within the background image. This position is sampled
277 from a generalization of the Bates distribution [2] with pa-
278 rameter $\eta \in \mathbb{N}$, visualized in Figure 3. We choose the bates
279 distribution, as it presents an easy way to sample from a
280 bounded domain with just one hyperparameter that controls
281 the concentration of the distribution. $\eta = 1$ corresponds to
282 the uniform distribution; $\eta > 1$ concentrates the distribution
283 around the center; and for $\eta < -1$, the distribution is con-
284 centrated at the borders. To more seamlessly integrate the
285 foreground, we apply a Gaussian blur with $\sigma \in [\frac{\sigma_{\max}}{10}, \sigma_{\max}]$,

inspired by the standard range for the Gaussian blur opera-
286 tion in [48], to the foreground’s alpha-mask. 287

288 We can apply standard data augmentation techniques in
289 two modes: Either we apply all augmentations to the recom-
290 bined image, or we apply the cropping and resizing to the
291 background only and then apply the other augmentations af-
292 ter recombination. The second mode ensures the foreground
293 object remains fully visible, while the first mode mirrors
294 standard data augmentation practices. 295

296 We experiment with a constant mixing ratio, or a linear or
297 cosine anealing schedule that increases the amount of images
298 from the original dataset over time. The mixing ratio acts as
299 a probability of selecting an image from the original dataset;
300 otherwise, an image with the same foreground is recombined
301 using *ForAug*. Thus, we still ensure each foreground is seen
once per epoch.

4. Experiments 302

303 We conduct a comprehensive suit of experiments to validate
304 the effectiveness of our approach. We compare training on
305 *ForNet*, the ImageNet instantiation of *ForAug*, to training on
306 ImageNet for 7 different models. Furthermore, we assess
307 the impact of using *ForNet* for pretraining on multiple fine-
308 grained downstream datasets. Additionally, we use *ForAug*’s
309 control over the image distribution to quantify some model
310 behaviors and biases.

4.1. Design Choices of *ForAug* 311

312 We start by ablating the design choices of *ForAug*. For
313 this, we revert to TinyImageNet [26], a subset of Image-
314 Net containing 200 categories with 500 images each, and
315 *TinyForNet*, a version of *ForAug* derived from TinyImageNet.
316 Table 1 presents the results of these ablations.

Dataset	Detect. prompt	Infill Model	FG. size	Augmentation Order	BG. strategy	BG. pruning	edge smoothing	original image mixing	TinyImageNet Accuracy	
									ViT-Ti [%]	ViT-S [%]
TinyImageNet									66.1 ± 0.5	68.3 ± 0.7
<i>TinyForNet</i>	specific	LaMa [45]	mean	crop→paste→color	same	-	-	-	64.6 ± 0.5	70.0 ± 0.6
<i>TinyForNet</i>	specific	LaMa [45]	range	crop→paste→color	same	-	-	-	65.5 ± 0.4	71.2 ± 0.5
<i>TinyForNet</i>	general	LaMa [45]	range	crop→paste→color	same	-	-	-	66.4 ± 0.6	72.9 ± 0.6
<i>TinyForNet</i>	general	Att. Eraser [43]	range	crop→paste→color	same	-	-	-	67.5 ± 1.2	72.4 ± 0.5
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	-	-	-	67.1 ± 1.2	72.9 ± 0.5
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	1.0	-	-	67.0 ± 1.2	73.0 ± 0.3
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	-	-	67.2 ± 1.2	72.9 ± 0.8
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.6	-	-	67.5 ± 1.0	72.8 ± 0.7
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	$\sigma_{\max} = 2.0$	-	67.2 ± 0.4	72.9 ± 0.5
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	$\sigma_{\max} = 4.0$	-	65.9 ± 0.5	72.4 ± 0.6
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	-	$p = 0.2$	69.8 ± 0.5	75.0 ± 0.3
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	-	$p = 0.33$	69.5 ± 0.4	75.2 ± 1.0
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	-	$p = 0.5$	70.3 ± 1.0	74.2 ± 0.2
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	-	linear	70.1 ± 0.7	74.9 ± 0.8
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	-	reverse lin.	67.6 ± 0.2	73.2 ± 0.3
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	-	cos	71.3 ± 1.0	75.7 ± 0.8
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	$\sigma_{\max} = 4.0$	cos	70.0 ± 0.8	75.5 ± 0.7
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	orig.	0.8	$\sigma_{\max} = 4.0$	cos	67.2 ± 0.9	69.9 ± 1.0
<i>TinyForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	all	0.8	$\sigma_{\max} = 4.0$	cos	70.1 ± 0.7	77.5 ± 0.6
<i>ForNet</i>									-	80.5 ± 0.1
<i>ForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	same	0.8	-	cos	-	80.7 ± 0.1
<i>ForNet</i>	general	Att. Eraser [43]	range	paste→crop→color	all	0.8	$\sigma_{\max} = 4.0$	cos	-	81.3 ± 0.1

Table 1. Ablation of design decisions of *TinyForNet* on TinyImageNet and *ForNet* on ImageNet.

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

321 **Inpainting.** Attentive Eraser [43] produces superior re-
 322 sults compared to LaMa [45] (see the supplementary for
 323 examples).

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

328 **Order of data augmentation.** Applying all aug-
 329 mentations after foreground-background recombination
 330 (*paste→crop→color*) slightly improves ViT-S’s perfor-
 331 mance compared to applying crop-related augmentations
 332 before pasting (*crop→paste→color*). For ViT-Ti, the results
 333 are ambiguous.

334 **Background pruning.** When it comes to the choice of
 335 backgrounds to use, we test two pruning thresholds (t_{prune})
 336 to exclude backgrounds with excessive inpainting. A thresh-
 337 old of $t_{\text{prune}} = 1.0$ means that we use all backgrounds that
 338 are not fully infilled. Varying t_{prune} has minimal impact.
 339 Therefore, we choose $t_{\text{prune}} = 0.8$ to exclude predomi-
 340 nantly artificial backgrounds. Similarly, applying edge smooth-
 341 ing to foreground masks with Gaussian blurring actually hurts
 342 performance on *TinyForNet*, but slightly improves it on *For-*
 343 *Net*.

344 **Mixing *ForNet*** with the original ImageNet data proves
 345 crucial. While constant and linear mixing schedules improve

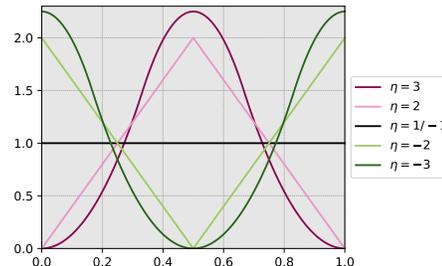


Figure 3. Plot of the probability distribution function (PDF) of the extended Bates distribution for different parameters η . Higher values of η concentrate the distribution around the center.

performance over no mixing by 2 – 3 p.p. compared to only
 using *TinyForNet*, the cosine annealing schedule yields the
 best results, boosting accuracy by another 0.5 – 1 p.p.

Background strategy. Another point is the allowed
 choice of background image for each foreground object.
 We compare using the original background, a background
 from the same class, and any background. These strategies
 go from low diversity and high shared information content
 between the foreground and background to high diversity
 and low shared information content. For *ViT-Ti*, the latter two
 strategies perform comparably, while *ViT-S* benefits from the
 added diversity of using any background. The same is true
 when training on the full (ImageNet) version of *ForNet*.

Foreground position. Finally, we analyze the foreground
 object’s positioning in the image. We utilize an extended
 Bates distribution to sample the position of the foreground

Training Set/ Bates Parameter	TIN	TinyForNet				
		$\eta = -3$	-2	$1/-1$	2	3
TinyImageNet	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 2. Accuracy of ViT-S trained on TinyImageNet (TIN) and TinyForNet with different foreground position distributions by varying the parameter of a Bates distribution η . The best performance is achieved using the uniform distribution ($\eta = 1$).

Dataset	Classes	Training Images	Validation Images
TinyImageNet	200	100,000	10,000
TinyForNet	200	99,404	9,915
ImageNet	1,000	1,281,167	50,000
ForNet	1,000	1,274,557	49,751

Table 3. Dataset statistics for TinyImageNet, TinyForNet, ImageNet, and ForNet. For ForNet and TinyForNet we report the number of foreground/background pairs.

362 object. The Bates distribution [2] with parameter $\eta \geq 1$ is the
 363 mean of η independent uniformly distributed random vari-
 364 ables [20]. Therefore, the larger η , the more concentrated the
 365 distribution is around the center. We extend this concept to
 366 $\eta \leq -1$ by defining $X \sim \text{Bates}(\eta) : \Leftrightarrow s(X) \sim \text{Bates}(-\eta)$
 367 for $\eta \leq 1$ with s being the sawtooth function on $[0, 1]$:

$$368 \quad s(x) = \begin{cases} x + 0.5 & \text{if } 0 < x < 0.5 \\ x - 0.5 & \text{if } 0.5 \leq x \leq 1 \end{cases} \quad (3)$$

369 Note that $s \circ s = \text{id}$ on $[0, 1]$. This way, distributions with
 370 $\eta \leq -1$ are more concentrated around the borders. $\eta = 1$
 371 and $\eta = -1$ both correspond to the uniform distribution.
 372 The PDF of this extended Bates distribution is visualized in
 373 Figure 3.

374 When sampling more towards the center of the image,
 375 the difficulty of the task is reduced, which then reduces
 376 the performance on TinyImageNet. This is reflected in the
 377 performance when evaluating on TinyForNet with $\eta = 2$
 378 and $\eta = 3$ compared to $\eta = -1/1$. We observe a similar
 379 reduction for $\eta < -1$. This experiment is conducted using
 380 the LaMa infill model.

381 After fixing the optimal design parameters in Table 1 (last
 382 row), we construct the full ForNet dataset using the entire
 383 ImageNet dataset. Table 3 compares the dataset statistics of
 384 ImageNet and ForNet. The slightly reduced image count in
 385 ForNet is due to instances where Grounded SAM failed to
 386 produce valid object detections.

Model	ImageNet Accuracy when trained on		Delta
	ImageNet	ForNet	
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
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 4. ImageNet results of models trained on ForNet and on ImageNet directly. ForNet improves the performance of all models in our test.

4.2. Image Classification Results

387 Table 4 compares the ImageNet performance of models
 388 trained on ForNet and ones trained directly on ImageNet.
 389 We adopt the training setup of [33] and [48] (details in the
 390 supplementary material) for training ViT [10], Swin [31] and
 391 ResNet [16] models. Notably, ForNet improves performance
 392 across all tested architectures, including the ResNet models
 393 (up to 1 p.p.), demonstrating benefits beyond Transformers.
 394 For Transformer models, we observe improvements from
 395 1.2 p.p. to 4.5 p.p. This improvement is more substantial for
 396 the larger models, with ViT-L gaining 4.5 p.p. in accuracy.
 397 ForNet’s improvements mostly counteract the drop in perfor-
 398 mance due to overfitting for large models. When training on
 399 ImageNet, this drop is 3.8 p.p. from ViT-S to ViT-L, while
 400 for ForNet it is reduced to 1.6 p.p.

401 To assess the transferability of ForNet-trained models,
 402 we finetune models pretrained on ImageNet and ForNet on
 403 five fine-grained datasets: FGVC-Aircraft [32], Stanford
 404 Cars [7], Oxford Flowers [34], Food-101 [21], and Oxford-
 405 IIIT Pets [35]. While for ResNets, the performance of both
 406 training datasets is about the same, for every Transformer,
 407 we see the accuracy improve on all downstream dataset by
 408 up to 7.3 p.p. and a reduction of error rate of up to 39.3%.
 409 In summary, these results demonstrate that the improved repre-
 410 sentation learning achieved by training on ForNet translates
 411 to superior performance not only on ImageNet, but also on a
 412 variety of fine-grained image classification tasks.

4.3. Further Model Evaluation

413 Beyond its use for training, ForNet’s unique properties and
 414 controlled data generation capabilities make it a powerful
 415 tool for analyzing model behavior and biases.

416 **Background Robustness** We assess the robustness of
 417 models to shifts in the background distribution from a class-
 418
 419

Model	Aircraft	Cars	Flowers	Food	Pets
ViT-S @ ImageNet	72.4 ± 1.0	89.8 ± 0.3	94.5 ± 0.2	89.1 ± 0.1	93.8 ± 0.2
ViT-S @ <i>ForNet</i>	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 @ ImageNet	71.7 ± 0.5	90.0 ± 0.2	94.8 ± 0.4	89.8 ± 0.2	94.1 ± 0.4
ViT-B @ <i>ForNet</i>	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
ViTL @ ImageNet	72.1 ± 1.0	88.8 ± 0.3	94.4 ± 0.3	90.1 ± 0.2	94.2 ± 0.4
ViTL @ <i>ForNet</i>	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
Swin-Ti @ ImageNet	77.0 ± 0.1	91.3 ± 0.6	95.9 ± 0.1	90.0 ± 0.2	94.2 ± 0.1
Swin-Ti @ <i>ForNet</i>	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 @ ImageNet	75.7 ± 1.4	91.0 ± 0.3	95.9 ± 0.5	91.1 ± 0.2	94.4 ± 0.1
Swin-S @ <i>ForNet</i>	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 @ ImageNet	78.2 ± 0.5	89.8 ± 0.2	91.7 ± 0.4	84.4 ± 0.2	93.7 ± 0.3
ResNet-50 @ <i>ForNet</i>	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	+0.1	±0
ResNet-101 @ ImageNet	78.4 ± 0.6	90.3 ± 0.1	91.2 ± 0.5	86.0 ± 0.2	94.3 ± 0.2
ResNet-101 @ <i>ForNet</i>	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

Table 5. Downstream accuracy in percent when finetuning on other datasets. Models were pretrained on *ForNet* and ImageNet. Pretraining on *ForNet* increases Transformer downstream accuracy on all datasets.

Model	Background Robustness when trained on		Delta
	ImageNet	<i>ForNet</i>	
ViT-S	0.73 ± 0.01	0.99 ± 0.01	+0.26
ViT-B	0.72 ± 0.01	1.00 ± 0.01	+0.28
ViT-L	0.70 ± 0.01	1.00 ± 0.01	+0.30
Swin-Ti	0.72 ± 0.01	1.00 ± 0.01	+0.28
Swin-S	0.72 ± 0.01	1.00 ± 0.01	+0.28
ResNet-50	0.79 ± 0.01	0.99 ± 0.01	+0.20
ResNet-101	0.79 ± 0.01	1.00 ± 0.01	+0.21

Table 6. Evaluation of the background robustness of models trained on *ForNet* and on ImageNet directly. Training on *ForNet* improves the background robustness of all model to ≈ 1.00 , meaning the model is indifferent to the choice of background.

related background to any background. Background robustness is defined to be the ratio of accuracy on *ForNet* with same-class backgrounds to accuracy with any background:

$$\text{Background Robustness} = \frac{\text{Acc}(\text{ForNet}_{\text{all}})}{\text{Acc}(\text{ForNet}_{\text{same}})} \quad (4)$$

It represents the relative drop in performance under a background distribution shift. Table 6 presents the background robustness of various models. When trained on ImageNet, smaller models generally exhibit greater robustness to changes in the background distribution than larger models and ResNet is more robust than the tested Transformer models. Crucially, training on *ForNet* instead of ImageNet improves the background robustness of all models to ≈ 1.00 ,

Model	Foreground Focus when trained on					
	IN		FN		IN	
	GradCam	GradCam++	IG	IN	FN	IN
ViT-S	1.2 ± 0.1	2.3 ± 0.3	1.2 ± 0.1	2.1 ± 0.4	1.9 ± 0.1	2.7 ± 0.1
ViT-B	1.2 ± 0.1	2.4 ± 0.7	1.1 ± 0.1	2.1 ± 0.1	1.7 ± 0.1	2.7 ± 0.1
ViT-L	1.3 ± 0.1	1.6 ± 0.1	1.1 ± 0.1	1.3 ± 0.1	1.3 ± 0.1	2.6 ± 0.1
Swin-Ti	0.9 ± 0.1	0.7 ± 0.1	1.0 ± 0.3	0.7 ± 0.3	2.5 ± 0.1	4.8 ± 0.3
Swin-S	0.8 ± 0.1	0.7 ± 0.1	0.7 ± 0.1	0.7 ± 0.4	2.4 ± 0.1	4.6 ± 0.3
ResNet-50	2.2 ± 0.1	2.7 ± 0.1	2.0 ± 0.1	2.9 ± 0.1	3.2 ± 0.1	4.9 ± 0.2
ResNet-101	2.3 ± 0.1	2.8 ± 0.1	2.2 ± 0.1	3.0 ± 0.1	3.2 ± 0.1	4.8 ± 0.1

Table 7. Evaluation of the foreground focus using GradCam, GradCam++ and IntegratedGradients of models trained on *ForNet* (FN) and on ImageNet (IN) directly. Training on *ForNet* improves the foreground focus of almost all models.

meaning that these models are agnostic to the choice of background and only classify based on the foreground. These findings highlight the generalization benefits of *ForNet*.

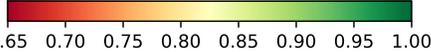
Foreground Focus Leveraging our inherent knowledge of the foreground masks when using *ForNet*, as well as common XAI techniques [4, 40, 44], we can evaluate a model’s focus on the foreground object. We can directly evaluate ImageNet trained models, but this technique can also be extended to other datasets without relying on manually annotated foreground-masks. To evaluate the foreground focus, we employ Grad-CAM [40], Grad-CAM++ [4] or IntegratedGradients (IG) [44] to compute the per-pixel importance of an image for the model’s prediction. The foreground focus is defined to be the ratio of the foreground’s relative importance to its relative size in the image:

$$\text{FG Focus}(\text{img}) = \frac{\text{Area}(\text{img}) \text{ Importance}(\text{fg})}{\text{Area}(\text{fg}) \text{ Importance}(\text{img})} \quad (5)$$

The foreground focus of a model is its average foreground focus over all test images. Table 7 presents our findings. Training on *ForNet* significantly increases the foreground focus of ViT and ResNet across all metrics used. For Swin, the foreground focus stagnates when measured using GradCam and GradCam++, but almost doubles when using IG.

Center Bias With *ForNet* we have unique control over the position of the foreground object in the image. This lets us quantify the center bias of ImageNet- and *ForNet*-trained models. We divide the image into a 3×3 grid and evaluate model accuracy when the foreground object is in each of the 9 grid cells. Each cell’s accuracy is divided by the accuracy in the center cell for normalization, which gives us the relative performance drop when the foreground is in each part of the image. The center bias is calculated as one minus the average of the minimum performance of a corner

Model	Center Bias when trained on		Delta
	ImageNet	ForNet	
ViT-S	 0.255 ± 0.008	 0.220 ± 0.003	-0.035
ViT-B	 0.254 ± 0.004	 0.190 ± 0.002	-0.064
ViT-L	 0.243 ± 0.011	 0.117 ± 0.007	-0.126
Swin-Ti	 0.250 ± 0.007	 0.165 ± 0.002	-0.085
Swin-S	 0.232 ± 0.001	 0.156 ± 0.002	-0.076
ResNet50	 0.263 ± 0.003	 0.197 ± 0.003	-0.066
ResNet101	 0.230 ± 0.003	 0.199 ± 0.002	-0.031



0.65 0.70 0.75 0.80 0.85 0.90 0.95 1.00

Table 8. Evaluation of the position bias. We plot the accuracy relative to the center accuracy of multiple instantiations of the models when the foreground objects is in different cells a 3 × 3 grid. Training on *ForNet* significantly reduces a models center bias.

464 cell and the minimum performance of a side cell:

Center Bias =

$$1 - \frac{\min_{a,b \in \{0,2\}} \text{Acc}(\text{cell}_{(a,b)}) + \min_{\substack{a=1 \text{ or } b=1 \\ a \neq b}} \text{Acc}(\text{cell}_{(a,b)})}{2\text{Acc}(\text{cell}_{(1,1)})} \quad (6)$$

466 Table 8 visualizes the center bias of three instantiations of each model. Performance is generally highest in the center and the center top and bottom and center left and right cells, and lowest in the four corners. Interestingly, ImageNet-trained models perform slightly better when the foreground object is on the right side of the image, compared to the left side, despite our use of random flipping with a probability of 0.5 during training. Training on *ForNet* significantly reduces center bias across all models. This demonstrates that *ForNet* promotes a more uniform spatial attention distribution. Their accuracy is higher in the center left and right cells than in the center top and bottom ones, which is not the case for ImageNet-trained models.

479 **Size Bias** Finally, we evaluate the impact of different-sized foreground objects on the accuracy. For this evaluation, we

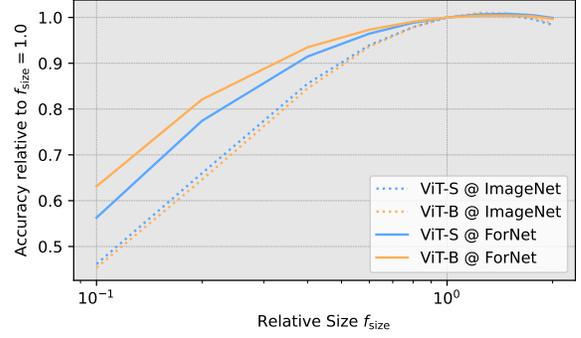


Figure 4. Evaluation of the size bias of models trained on *ForNet*. We plot the accuracy relative to the accuracy when using the mean foreground size.

use the *mean* foreground size strategy. We introduce a size factor f_{size} by which we additionally scale the foreground object before pasting it onto the background. Results are again normalized by the accuracy when using the mean foreground size ($f_{\text{size}} = 1.0$). Figure 4 shows the size bias curves of ViT-S and ViT-B when trained on ImageNet and *ForNet*. Models trained on *ForNet* maintain better performance even with smaller foreground objects, when ImageNet-trained models exhibit a more rapid performance decline. Therefore, *ForNet*-training improves robustness to variations in object scale.

5. Discussion & Conclusion

We introduce *ForAug*, a novel data augmentation scheme that facilitates improved Transformer training for image classification. By explicitly separating and recombining foreground objects and backgrounds, *ForAug* enables controlled data augmentation, leading to significant performance gains on ImageNet and downstream fine-grained classification tasks. Furthermore, *ForAug* provides a powerful framework for analyzing model behavior and quantifying biases, including background robustness, foreground focus, center bias, and size bias. Our experiments demonstrate that training on *ForNet*, the instantiation of *ForAug* on ImageNet, not only boosts accuracy but also significantly reduces these biases, resulting in more robust and generalizable models. In the future, we see *ForAug* be also applied to other datasets and tasks, like video recognition or segmentation. *ForAug*'s ability to both improve performance and provide insights into model behavior makes it a valuable tool for advancing CV research and developing more reliable AI systems.

Acknowledgements

Will be in the final paper.

513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568

References

[1] Khaled Alomar, Halil Ibrahim Aysel, and Xiaohao Cai. Data augmentation in classification and segmentation: A survey and new strategies. 9(2):46. 1

[2] G.E. Bates. Joint distributions of time intervals for the occurrence of successive accidents in a generalized polya urn scheme. *Annals of Mathematical Statistics*, 26:705–720, 1955. 4, 6

[3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. 2020. 1

[4] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 839–847, 2018. 7

[5] Ekin D. Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation policies from data. 2018. 2

[6] Ekin D. Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V. Le. Randaugment: Practical automated data augmentation with a reduced search space. 2019. 2

[7] Afshin Dehghan, Syed Zain Masood, Guang Shu, and Enrique G. Ortiz. View independent vehicle make, model and color recognition using convolutional neural network. 2017. 6

[8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2009. 1

[9] Peijian Ding, Davit Soselia, Thomas Armstrong, Jiahao Su, and Furong Huang. Reviving shift equivariance in vision transformers. 2023. 2

[10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. 1, 6

[11] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, paste and learn: Surprisingly easy synthesis for instance detection. 2017. 2

[12] Yunhao Ge, Jiashu Xu, Brian Nlong Zhao, Neel Joshi, Laurent Itti, and Vibhav Vineet. Beyond generation: Harnessing text to image models for object detection and segmentation. *ArXiv*, abs/2309.05956, 2023. 2

[13] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. 2018. 3

[14] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D. Cubuk, Quoc V. Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. 2020. 2

[15] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. 2013. 1

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 6

[17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. 2017. 1

[18] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. 2019. 3

[19] Stefan Hinterstoisser, Olivier Pauly, Hauke Heibel, Marek Martina, and Martin Bokeloh. An annotation saved is an annotation earned: Using fully synthetic training for object detection. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 2787–2796, 2019. 2

[20] Norman L. Johnson, Samuel Kotz, and N. Balakrishnan. *Continuous Univariate Distributions*. Wiley, 2 edition, 1995. Wiley series in probability and mathematical statistics. 6

[21] Parneet Kaur, Karan Sikka, and Ajay Divakaran. Combining weakly and weakly supervised learning for classifying food images. 2017. 6

[22] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM Computing Surveys*, 54(10s):1–41, 2022. 1

[23] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. 2023. 3

[24] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *Computer Vision – ECCV 2020*, pages 491–507, Cham, 2020. Springer International Publishing. 2

[25] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2012. 1

[26] Yann Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015. 4

[27] Xiaodan Li, Yuefeng Chen, Yao Zhu, Shuhui Wang, Rong Zhang, and Hui Xue. Imagenet-e: Benchmarking neural network robustness via attribute editing. 2023. 3

[28] Evan Ling, Dezhao Huang, and Minhoe Hur. Humans need not label more humans: Occlusion copy & paste for occluded human instance segmentation. 2022. 2

[29] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. 2023. 3

[30] Yue Liu, Christos Matsoukas, Fredrik Strand, Hossein Azizpour, and Kevin Smith. Patchdropout: Economizing vision transformers using patch dropout. 2022. 2

625 [31] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng
626 Zhang, Stephen Lin, and Baining Guo. Swin transformer:
627 Hierarchical vision transformer using shifted windows. In
628 *2021 IEEE/CVF International Conference on Computer Vi-*
629 *sion (ICCV)*, pages 9992–10002. Los Alamitos, CA, USA,
630 2021. IEEE Computer Society. 1, 6

631 [32] S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi.
632 Fine-grained visual classification of aircraft. Technical report,
633 2013. 6

634 [33] Tobias Christian Nauen, Sebastian Palacio, and Andreas Dengel.
635 Which transformer to favor: A comparative analysis of
636 efficiency in vision transformers, 2023. 6

637 [34] Maria-Elena Nilsback and Andrew Zisserman. Automated
638 flower classification over a large number of classes. In *Indian*
639 *Conference on Computer Vision, Graphics and Image*
640 *Processing*, 2008. 6

641 [35] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and
642 C. V. Jawahar. Cats and dogs. In *IEEE Conference on Com-*
643 *puter Vision and Pattern Recognition*, 2012. 6

644 [36] Gabriela Rangel, Juan C. Cuevas-Tello, Jose Nunez-Varela,
645 Cesar Puente, and Alejandra G. Silva-Trujillo. A survey on
646 convolutional neural networks and their performance limita-
647 tions in image recognition tasks. 2024(1). 1

648 [37] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li,
649 He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan,
650 Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang
651 Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling
652 open-world models for diverse visual tasks. 2024. 2, 3

653 [38] Renan A. Rojas-Gomez, Teck-Yian Lim, Minh N. Do, and
654 Raymond A. Yeh. Making vision transformers truly shift-
655 equivariant. 2023. 2

656 [39] Edward Sanderson and Bogdan J. Matuszewski. *FCN-*
657 *Transformer Feature Fusion for Polyp Segmentation*, pages
658 892–907. Springer International Publishing. 1

659 [40] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das,
660 Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-
661 cam: Visual explanations from deep networks via gradient-
662 based localization. 128(2):336–359, 2016. 7

663 [41] Ang Jia Ning Shermaine, Michalis Lazarou, and Tania
664 Stathaki. Image compositing is all you need for data aug-
665 mentation. 2025. 2

666 [42] Connor Shorten and Taghi M. Khoshgoftaar. A survey on
667 image data augmentation for deep learning. 6(1). 1, 2

668 [43] Wenhao Sun, Benlei Cui, Xue-Mei Dong, and Jingqun Tang.
669 Attentive eraser: Unleashing diffusion model’s object removal
670 potential via self-attention redirection guidance. 2024. 2, 3, 5

671 [44] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic
672 attribution for deep networks. 2017. 7

673 [45] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin,
674 Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov,
675 Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lem-
676 pitsky. Resolution-robust large mask inpainting with fourier
677 convolutions. 2021. 2, 3, 5

678 [46] Ryo Takahashi, Takashi Matsubara, and Kuniaki Uehara. Data
679 augmentation using random image cropping and patching for
680 deep cnns. 30(9):2917–2931, 2018. 2

[47] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco
681 Massa, Alexandre Sablayrolles, and Herve Jegou. Training
682 data-efficient image transformers & distillation through atten-
683 tion. In *Proceedings of the 38th International Conference on*
684 *Machine Learning*, pages 10347–10357. PMLR, 2021. 1
685

[48] Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii:
686 Revenge of the vit. In *Computer Vision – ECCV 2022*, pages
687 516–533, Cham, 2022. Springer Nature Switzerland. 1, 2, 4,
688 6
689

[49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszko-
690 reit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia
691 Polosukhin. Attention is all you need. In *Advances in Neu-*
692 *ral Information Processing Systems*. Curran Associates, Inc.,
693 2017. 1
694

[50] Ioannis A. Vezakis, Konstantinos Georgas, Dimitrios Fotiadis,
695 and George K. Matsopoulos. Effisegnet: Gastrointestinal
696 polyp segmentation through a pre-trained efficientnet-based
697 network with a simplified decoder. 2024. 1
698

[51] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang
699 Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed,
700 Saksham Singhal, Subhojit Som, and Furu Wei. Image as a
701 foreign language: Beit pretraining for all vision and vision-
702 language tasks. 2022. 1
703

[52] Wenhui Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi
704 Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng
705 Li, Xiaogang Wang, and Yu Qiao. Internimage: Exploring
706 large-scale vision foundation models with deformable convo-
707 lutions. 2022. 1
708

[53] Levi Kassel Michael Werman. Deepaste – inpainting for
709 pasting. 2021. 2
710

[54] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Re-
711becca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos,
712 Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon
713 Kornblith, and Ludwig Schmidt. Model soups: averaging
714 weights of multiple fine-tuned models improves accuracy
715 without increasing inference time. In *Proceedings of the*
716 *39th International Conference on Machine Learning*, pages
717 23965–23998. PMLR, 2022. 1
718

[55] Kai Xiao, Logan Engstrom, Andrew Ilyas, and Aleksander
719 Madry. Noise or signal: The role of image backgrounds in
720 object recognition. 2020. 3
721

[56] Mingle Xu, Sook Yoon, Alvaro Fuentes, and Dong Sun Park.
722 A comprehensive survey of image augmentation techniques
723 for deep learning. 137:109347. 1, 2
724

[57] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mo-
725 jtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive
726 captioners are image-text foundation models. *Transactions*
727 *on Machine Learning Research*, 2022. 1
728

[58] Sangdoon Yun, Dongyoon Han, Sanghyuk Chun, Seong Joon
729 Oh, Youngjoon Yoo, and Junsuk Choe. CutMix: Regular-
730 ization strategy to train strong classifiers with localizable
731 features. In *2019 IEEE/CVF International Conference on*
732 *Computer Vision (ICCV)*. IEEE. 2
733

[59] Chenshuang Zhang, Fei Pan, Junmo Kim, In So Kweon, and
734 Chengzhi Mao. Imagenet-d: Benchmarking neural network
735 robustness on diffusion synthetic object. 2024. 3
736

- 737 [60] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and
738 David Lopez-Paz. mixup: Beyond empirical risk minimiza-
739 tion. In *International Conference on Learning Representa-*
740 *tions*, 2018. 2
- 741 [61] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and
742 Yi Yang. Random erasing data augmentation. 2017. 2
- 743 [62] Zhuofan Zong, Guanglu Song, and Yu Liu. Detsr with collab-
744 orative hybrid assignments training. 2022. 1