Corner Cases: How Size and Position of Objects Challenge ImageNet-Trained Models

Anonymous CVPR submission

Paper ID *****

Abstract

Backgrounds in images play a major role in contributing 001 002 to spurious correlations among different data points. Owing 003 to aesthetic preferences of humans capturing the images, datasets can exhibit positional (location of the object within 004 005 a given frame) and size (region-of-interest to image ratio) 006 biases for different classes. In this paper, we show that these 007 biases can impact how much a model relies on spurious features in the background to make its predictions. To better 008 009 illustrate our findings, we propose a synthetic dataset derived from ImageNet1k, Hard-Spurious-ImageNet, which contains 010 011 images with various backgrounds, object positions, and object sizes. By evaluating the dataset on different pretrained 012 models, we find that most models rely heavily on spurious 013 features in the background when the region-of-interest (ROI) 014 to image ratio is small and the object is far from the center 015 of the image. Moreover, we also show that current methods 016 017 that aim to mitigate harmful spurious features, do not take into account these factors, hence fail to achieve considerable 018 019 performance gains for worst-group accuracies when the size 020 and location of core features in an image change.

021 1. Introduction

Spurious features are defined as features that are predictive 022 of the class label without being directly related to it. Such 023 features are usually helpful for object recognition when the 024 025 object is placed in a *perfect* environment or context. An example of that would be a sea lion near a body of water. 026 This is because most models learn to associate water with sea 027 lions and vice versa. On the contrary, spurious features can 028 029 be extremely harmful when the object or the "core" features 030 are observed in an unusual environment or against a spurious background. This scenario can happen when the model is 031 deployed in the wild. Deep neural networks can be fooled 032 easily to predict the label from the spurious cues in the back-033 ground without relying on "object" or "core" features in the 034 035 image itself. Recently, a plethora of techniques have been



Figure 1. Gradcam visualizations for Pre-trained ConvNext-Base. a) Model predicts core class "Tench" when the object is located in the center of the image, b) Spurious class "Zucchini" is predicted when the "core" class moves away from the center, c) Class "GoldFish" is predicted when the size of the core object is large (112×112) , d) Spurious class "Sea Lion" is predicted when size of core object reduces to 84×84 .

proposed to mitigate the reliance on unnecessary cues for 036 image classification. Sagawa et al. [20] introduced a distribu-037 tionally robust optimization technique which, coupled with 038 strong regularization, helped in achieving high accuracies for 039 data groups that have strong spurious feature reliance. Simi-040 larly, Kirichenko et al. [6] address this problem by retraining 041 the last layer of a DNN using equal data points from different 042 groups with core and spurious backgrounds. These methods 043 are helpful when the test set exhibits similar biases as the 044 training data, yet they fail to achieve similar performance 045 gains when these biases are explicitly removed. 046

Biases in datasets can hugely impact a deep neural net-047work's performance. Earlier works have proven that convo-048lutional neural networks are not entirely translation invariant049



Figure 2. ImageNet classes and their center and size scores. *Toyshop* has largest center and size scores, whereas *Volleyball* has smallest center score and *Balance Beam* has smallest size score. Other classes are sampled randomly for visualization.

050 and have the capacity to learn location information about objects [2]. Some studies have found that models perform 051 poorly on untrained locations [1]. Similarly, object size 052 within an input frame can lead to models performing badly 053 054 when the sizes differ at inference time. The deep learning community has tried to mitigate the effect of these biases by 055 proposing different data augmentation techniques that ensure 056 that models are robust to changes in size and locations of the 057 objects. However, the impact of the aforementioned factors 058 in the presence of spurious features remains less explored. 059

In this work, we try to answer the questions: In the absence of biases mentioned above, namely position and size
of objects, how much do pre-trained models rely on spurious
backgrounds to make their predictions, and are the current
techniques that mitigate harmful spurious features, enough
to tackle this problem? Specifically, the contributions of our
work are as follows:

- We calculate centeredness and size scores of different classes in ImageNet [4], and analyze their relation with the level of spuriousity present in that class.
- We derive a dataset from ImageNet1k, called Hard Spurious-ImageNet, containing objects against spurious
 backgrounds with varying sizes and positions. The code
 to generate the dataset will be provided.
- With the help of experimentation and ablation, we conclude that the size and location of the object should be taken into account when trying to mitigate harmful spurious correlations in the dataset.

078 2. Related Work

079 2.1. Spurious Features

Moayeri et al. [13] show that adversarial training increases
model reliance on spurious features. They also show that
increased spurious feature reliance occurs when the perturbations added to core features are too small to break spurious
correlations. Murali et al. [17] show that spurious features
are related with a model's learning dynamics. Specifically,
"easier" features learnt in the start of model training can hurt



Figure 3. Counts in log scale of relative centers of ground truth bounding boxes containing the object corresponding to the image class (ImageNet1k validation set). Most object centers are concentrated around the image center, while some are present along the main axes. Objects of interest are rarely present in image corners.

generalization. Neuhaus et al. [18] proposed a method to
identify spurious features in the ImageNet dataset and intro-
duced a fix to mitigate a model's dependence on these fea-
tures without requiring additional labels. While the proposed
methods to mitigate spurious feature reliance are helpful in
many cases, their efficacy is less known when factors such
as size and location of core features in an image change.087
088
089

2.2. Existing Datasets

Xiao et al. [25] present an analysis of model's performance 095 as a function of varying backgrounds and foregrounds for 096 ImageNet. They conclude that more accurate models have 097 less reliance on backgrounds. They also propose a dataset 098 called ImageNet-9 with mixed foregrounds and backgrounds. 099 Moayeri et al. [14] propose a dataset derived from ImageNet 100 with segmentation masks for a subset of images. These 101 masks label entire objects and various visual attributes. They 102 name this dataset RIVAL10 and also test different models' 103 sensitivity to noise in backgrounds and foregrounds. Moay-104 eri et al. [15] propose a dataset with segmentation masks 105 for images in 15 classes of ImageNet1k. These images have 106 high spurious features. They attribute this to objects being 107 small and less centered in these images. Singla and Feizi 108 [21] label spurious and core features for ImageNet samples. 109 They achieve this by making use of activation maps as soft 110 masks. Moayeri et al. [16] rank images in ImageNet dataset 111 based on spurious cues present. They show that spurious 112 feature reliance is influenced more by the data a model is 113 trained on rather than how a model is trained. Lynch et al. 114 [12] propose a photo-realistic dataset with many-to-many 115 spurious correlations between different groups of spurious 116 attributes and classes. One work closely related to ours is 117 [27]. They do a fine-grained analysis of the robustness of 118 different models by varying factors such as object size, lo-119 cation, and rotation. Our technical contributions differ from 120 theirs because we take into account the spuriosity level of 121 backgrounds and correlate it with the above factors as well. 122

147

161



Figure 4. Correlation between the validation accuracy on inpainted ImageNet and, from left to right, center scores, size scores, and their product, respectively. Jointly considering center and size score shows strongest negative correlation with the accuracy.



Figure 5. Histograms showing distribution of scores in different classes of ImageNet1k dataset.

123 2.3. Biases in Datasets

124 While capturing images through a camera, humans often 125 tend to place the region of interest in the center. Due to 126 this, there often exists a bias in classification datasets where objects are mostly located in the center of images and away 127 from the boundary of the image. Exploiting the center bias 128 in ImageNet, resizing and center cropping has been usually 129 130 used for testing image classification models. Taesiri et al. [23] show that there exists a strong center bias in out-of-131 132 distribution benchmarks such as ImageNet-A and ObjectNet 133 by using resize and center crop operations only. They resize the image to multiple scales and patchify it, followed by a 134 center crop operation at every patch. Doing this, they end up 135 136 with different zoomed-in versions of the input images. The computed accuracy of the center crop is maximum showing 137 the presence of a strong center bias in the dataset. In this pa-138 per, we do an in-depth analysis of the presence of center and 139 size bias in every class of ImageNet by computing distinct 140 141 scores. The detailed explanation of these scores are given in 142 following sections.

3. Biases in ImageNet

In this section, we quantitatively analyze positional and size biases present in ImageNet1k. To get a better sense of these biases, we propose *centeredness* and *size* scores.

3.1. Centeredness Score

In the majority of images in ImageNet1k, the objects of
interest are located in the image's center (see Figure 3).148Hence, in this paper, we use "positional" and "center" as
synonyms. To understand the extent of center bias prevalent
in ImageNet1k, we propose a *Center Score* defined as150

$$C_c = \frac{1}{M} \frac{1}{N} \sum_{i=1}^{M} \sum_{j=1}^{N} 1 - (\|I_{i,c} - O_{i,j,c}\|_{\infty}), \quad (1) \quad 153$$

where C_c is the centeredness score for class c, M is total number of images in the class, N is total number of objects within a frame, I is image center, and O is object center. The distance between image center and object center is calculated by the ℓ_{∞} norm. It is subtracted from 1 to establish a direct relationship between the score and center bias prevalent in the class c. 154 155 156 157 158 158 159 159 159 159

3.2. Size Score

To measure the average sizes of objects within images, we define a size score as 163

$$S_c = \frac{1}{M} \frac{1}{N} \sum_{i=1}^{M} \sum_{j=1}^{N} \frac{h_j w_j}{H_i W_i},$$
(2) 164

where S_c is the size score for class c, h and w refer to the 165 height and width of object j in image i. H and W are the 166 height and width of the image itself. Figure 2 shows the 167 center and size scores of different classes, with Toyshop 168 having the maximum center and size scores. The histograms 169 in Figure 5 show the distribution of center and size scores 170 of all the classes in the ImageNet1k validation data. It can 171 be seen that the majority of the classes in ImageNet1k are 172 highly centered with objects of interest occupying half of 173

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

the image pixels on average. These scores are calculated byusing Ground Truth bounding boxes of ImageNet.

176 3.3. Relationship with the Level of Spuriosity

To establish a correlation between centeredness and size 177 scores of every class to spurious feature reliance in Ima-178 179 geNet, we first calculate the validation accuracies of different 180 classes in ImageNet with object information removed. We achieve this by using Inpaint-Anything [26] with the goal of 181 182 creating a more realistic effect when the region of interest is removed from the image. The input to Inpaint Anything are 183 the object bounding boxes and it makes use of Segment Any-184 thing [7] to predict masks for objects within these bounding 185 186 boxes. These predicted masks are then input to the inpainting model LaMa [22] which fills the masked region predicted by 187 SAM. Finally, we resize the inpainted images to 224×224 . 188 We use ConvNext-Base [11] pre-trained on ImageNet22k 189 and fine-tuned on ImageNet1k, to compute the validation 190 191 accuracies for the inpainted dataset. Classes with higher val-192 idation accuracies indicate higher spurious feature reliance, since the model has learnt to associate the class label not 193 just with the core object, but also with the background infor-194 mation. In order to assess the correlation present between 195 center and size scores and the level of spuriousity present in 196 197 different classes of ImageNet, we use Kendall's τ coefficient 198 and Spearman's correlation coefficient. The negative corre-199 lation values (see Figure 4) depict that there is an inverse relationship between both inpainted data's accuracy and the 200 different considered scores, which validates the hypothesis 201 that a higher spurious feature reliance is observed in case of 202 203 non-centered small object sizes. The correlation is overall rather weak, which is to be expected since different classes 204 are differently hard to classify, even from their core features. 205

4. Dataset

Similar to the waterbirds dataset [20], we assume that ev-207 208 ery datapoint (x, y) has an attribute $a(x) \in A$ which is 209 spuriously correlated with label y. We conjecture that the strength of the correlation between attribute a(x) and label 210 y is controlled by two factors: size s and position p of the 211 core features in the input image. To investigate this corre-212 213 lation, we propose Hard-Spurious-ImageNet, a synthetic 214 dataset to illustrate the problem of spurious feature reliance in the presence of varying object bounding box sizes, loca-215 216 tions, and backgrounds. The prime motivation of creating the dataset is to have precise control over these factors and 217 218 help the community build robust models against stronger 219 spurious cues.

We consider the image content within the provided ground truth object bounding boxes for ImageNet as core features and the features outside the bounding box as the background. In ImageNet, bounding boxes are available for all images in the validation data, yet only a subset of images

in training data are annotated. The images are annotated 225 and verified through Amazon Mechanical Turk. We employ 226 these annotations to provide us an estimate of the location of 227 core features in any image. A brief analysis of these anno-228 tations using Grounding DINO is provided in the appendix 229 in section 12. As a first step, we want to disentangle core 230 features from the rest of the image. We achieve this by crop-231 ping out the core objects from the images and inpainting the 232 resulting image, as explained in the previous section. Next, 233 we resize core object bounding boxes to different sizes, and 234 place them in two different locations against inpainted back-235 grounds. The size and location of core objects and the kind 236 of background chosen, gives rise to different groups in the 237 data. To efficiently gauge the performance of these different 238 groups, we categorize them as follows: 239

- **Group CeO**: Core object in the *Ce*nter of image against its *O*riginal inpainted background.
- **Group CoO**: Core object in the top right *Co*rner of image against its *O*riginal inpainted background.
- **Group CeR**: Core object in the *Ce*nter of image against *R* andom inpainted background.
- **Group CoR**: Core object in the top right *Co*rner of image against *R* andom inpainted background.

We consider three core object sizes: 56×56 , 84×84 , and 112×112 . It is important to note that all the inpainted backgrounds have already been resized to 224×224 , so the core object sizes mentioned above represent $\frac{4}{64}$ th, $\frac{9}{64}$ th, and $\frac{16}{64}$ th of the whole image.

We also experimented with object masks obtained from the Segment Anything [7] model rather than the provided bounding boxes as foreground objects (see Table 7 in supplementary). We observed that the mask quality for some objects was not good enough, hence, we used provided bounding boxes for this work.

4.1. Hard-Spurious-ImageNet-v2

Randomly chosen backgrounds have varying levels of spu-260 riosity based on the classes they are taken from. We derive a 261 variant of the proposed dataset where, instead of choosing 262 backgrounds in a random fashion, they are chosen based on 263 the level of spurious features present in them. To achieve 264 this, we first analyze the level of spuriosity present in every 265 class. We give inpainted images without the core objects, as 266 input to the pretrained ConvNext-Base model, and record 267 the accuracies of every class. The classes where accuracies 268 are high indicate that the model has learnt to predict the class 269 label without the presence of core objects. On the contrary, 270 classes for which the accuracy is low are highly reliant on 271 core features to make predictions. We choose 10 classes 272 that are highly spurious, namely: snorkel, bobsled, maypole, 273 potter's wheel, gondola, bearskin, volleyball, basketball, ca-274 noe, geyser, and yellow lady's slipper as backgrounds. For 275 foreground objects, we choose 10 classes with high core 276

294



Figure 6. Different samples from Hard-Spurious-ImageNet. Image size remains same in all images i.e. 224×224 , whereas object size changes. The label of every image is same as the label of the foreground object.

Model	Clean Accuracy
ConvNext-Base	85.86
ResNet-50	80.20
CoAtNet	83.59

Table 1. Clean accuracies of standard ImageNet validation data with different pre-trained models.

277 features such as: *bluetick*, *box turtle*, *Chihuahua*, *Japanese* 278 spaniel, Maltese dog, Shih-Tzu, Blenheim spaniel, papillon, Rhodesian ridgeback, and basset. We combine the above-279 mentioned foregrounds and backgrounds to create a dataset 280 281 with 10 classes of foreground objects and highly spurious 282 backgrounds. Similar to before, for every class, the chosen 283 background class remains same for all images belonging to 284 that class, but the backgrounds can differ from one image to another. Finally, we create four groups for the dataset as 285 before and test on pre-trained models. 286

287 5. Experimental Results

288 We test the robustness of different models with the two pro-289 posed two variants of Hard-Spurious-ImageNet. The images 290 are already resized to 224×224 , so no additional resizing is 291 applied to the images when giving as input to the pre-trained 292 models. Images are normalized with mean and standard deviation of the ImageNet dataset. We use HuggingFace PyTorch models to test the dataset.

Figure 7 shows test accuracies of the proposed data 295 and its variant on three pretrained models. We consider 296 ConvNext-Base trained on ImageNet22k and fine-tuned with 297 ImageNet1k, ResNet-50 [5] and CoAtNet [3] pretrained on 298 ImageNet1k to test the performance of proposed dataset. De-299 tailed results are given in Tables 5 and 6 in supplementary 300 section. ConvNext Base performs best across all groups and 301 datasets. This can be attributed to the fact that the data aug-302 mentation pipeline of ConvNext-Base consists of rigorous 303 steps, which ensures it stays robust to varying object sizes 304 and locations. The difference in accuracy between groups 305 CeR and CoR, when the core object size is 112×112 is less 306 across all the models. This indicates that the core feature 307 size is big enough for the model to ignore changes in loca-308 tion. Moreover, $\frac{1}{4}$ th of the number of pixels in the image 309 are occupied by core features in this case, so backgrounds 310 are less exposed as compared to when the core object size is 311 even less. Another interesting observation is that the impact 312 of size change is far stronger on model performance than the 313 location of core features. We also see that Hard-Spurious-314 ImageNet-v2 has far worse performance on groups CeR and 315 CoR across all architectures and sizes. This indicates that the 316 strength of spurious backgrounds is far greater than that of 317 core features when the size of core features starts to decrease. 318





Figure 7. Benchmarking results of different models. Performance for our Hard-Spurious-ImageNet-v2 is the worst across all groups.

We also observe that in almost all the groups, there is significant drop in performance compared with clean accuracies
on standard validation dataset (see Table 1).

322 Based on the above observations, we divide all the 12 groups consisting of different core feature sizes and loca-323 324 tions into three distinct categories: Easy: This set con-325 sists of Groups CeO and CoO for larger core feature sizes, i.e. 84×84 and 112×112 , as these groups seem to be doing 326 considerably better than the rest. Hard: Groups CeR and 327 CoR are the worst performing across all architecture for core 328 329 feature sizes 56×56 and 84×84 . We categorize them as 330 Hard group. The remaining groups, i.e. groups CeO and CoO for size 56×56 , and groups CeR and CoR for size 331 112×112 seem to be performing moderately, we put them 332 in Medium category. 333

334 Following the analysis done earlier (see Figure 5), we find that most of the images in ImageNet are centered with an 335 336 estimated size score of ≈ 0.5 , indicating that on average, the core features in an image occupy half the number of pixels of 337 the entire image. Keeping this in mind, we create the training 338 339 data of Hard-Spurious-ImageNet consisting of majority and minority groups, where the number of images belonging to 340 341 majority groups are far more than in minority groups. This is done to replicate the long-tailed distribution nature of the 342 ImageNet dataset in terms of hardness. For the training data, 343 we consider 80 images per group in the Easy category and 344 10 images from groups in Medium and Hard categories. This 345 346 brings the total to 400 images per class in the training data. 347 Out of the 400 images, 320 images belong to the Easy group and 80 to the Medium and Hard groups. For the validation 348 set, we use a balanced dataset having equal data points from 349 350 every group. We use 20 images per group, resulting in 240 images per class. Both training and validation set of 351 Hard-Spurious-ImageNet are derived from training data of 352 353 ImageNet, whereas the test set is derived from the validation data. The test set is also balanced, comprising 50 images per 354 355 group, totaling 600 images in every class.

356 5.1. Effects of Data Augmentation and Self-357 Supervised Models

To measure the effect of data augmentations, we compared vanilla ResNet-50 trained without any augmentations on ImageNet1K with an advanced training recipe involving auto-360 augment, random erase, mixup, and cutmix. The results 361 (shown in Table 2) indicate that while data augmentation 362 increases accuracy across groups CeO, CoO, and CeR, the 363 performance decreases in case of group CoR for all sizes. 364 This indicates that standard data augmentation approaches 365 do not take into account the presence of spurious features in 366 the data while augmenting, hence, may end up highlighting 367 them instead. Moreover, the gap in performance still persists 368 across all four groups for a given core object size. This hints 369 that mere data augmentation strategies are insufficient to deal 370 with this problem. In the supplementary materials provided 371 (see Table 5 and Table 6), we test the model on Hiera-Base 372 with Masked Autoencoder which has been trained in a self-373 supervised manner. The results follow a similar trend across 374 groups as other methods shown in the paper, although the 375 Group CoR for size shows the worst performance when com-376 pared with all the other architectures. Moreover, we also 377 computed the performance of different groups in the pro-378 posed dataset on a ViT pretrained on WIT-400M image-text 379 pairs by OpenAI using CLIP and fine-tuned on ImageNet1k. 380 The results are given in Table 2 and show similar trends as 381 reported earlier. 382

5.2. Group Robustness Methods

We measure the performance of the proposed dataset using 384 simple fine-tuning and two state-of-the-art group robustness 385 methods. Empirical Risk Minimization or ERM [24] is 386 conventional training to optimize average training accuracy 387 without specialized methods for optimizing worst-group ac-388 curacy. Deep Feature Reweighting or **DFR** [6] tackles the 389 problem of spurious correlations by retraining the last layer 390 of a pre-trained model with equal data points from different 391 groups present in the training data. Just Train Twice or JTT 392 [9] upsamples the training images which were wrongly pre-393 dicted by the ERM trained model by a certain factor λ_{up} , 394 and trains the classifier again. We experiment with different 395 variations of the above methods. The implementation details 396 are given in the supplementary material. 397

Madal	Clean	Object	Group Accuracies			
wodel	Accuracy	Resolution	CeO	CoO	CeR	CoR
		56^{2}	38.62	32.53	8.71	7.03
ResNet-50 (Baseline)	76.13	84^{2}	56.46	52.47	28.44	27.37
		112^{2}	65.87	64.16	46.58	46.57
	80.33	56^{2}	49.14	38.47	13.19	4.49
ResNet-50 (Data Augmentations)		84^{2}	65.74	58.48	33.30	20.01
		112^{2}	72.93	68.19	45.12	36.40
		56^{2}	44.78	39.37	7.15	5.57
ViT Base	81.92	84^{2}	63.65	58.83	28.56	25.46
		112^{2}	71.31	70.46	46.43	47.93

Table 2. The first two rows show the impact of data augmentation on the proposed dataset. Performance across group CoR becomes worse, indicating that just augmenting the data might not be enough to deal with spurious correlations. The third row shows the performance on ViT-Base pre-trained using CLIP and fine-tuned on IN-1K, highlighting similar trends observed earlier.

Methods	Easy	Medium	Hard	Average
Pretrained	65.39	48.50	16.54	43.48
\mathbf{ERM}	74.84	66.67	57.56	65.94
$_{\rm JTT}$	60.90	53.09	46.49	53.50
DFR	72.47	65.65	59.79	65.97

Table 3. Test Performance of different methods on Easy, Medium, and Hard categories in Hard-Spurious-ImageNet. Average accuracy is the average test performance of all the groups combined.



Figure 8. (left) The effect of training epochs of ERM model on the performance of DFR. ERM model trained with 20 epochs gives the highest performance for DFR. (right) ERM_{all} narrows the gap between easy, medium, and hard groups.

5.3. Results 398

399 400 401

The results in Table 3 show that pretrained ImageNet models perform worst on the hard group. This could be attributed to the fact that the model has very little exposure to small core features against spurious backgrounds in the training 402 data. The ERM model does better across easy, medium and 403 404 hard groups, but there still exists a disparity in performance 405 among the three groups.

DFR is able to perform slightly better in the Hard group 406 by sacrificing some accuracy in Easy and Medium groups. 407 408 The average test accuracy is similar for ERM and DFR. The 409 performance with JTT also decreases, which hints that the

Size	CeO	CoO	CeR	CoR
56^{2}	62.25	60.8	54.56	54.45
84^{2}	73.35	72.60	69.76	69.96
112^{2}	77.19	77.13	75.34	75.48

Table 4. Breakdown of test accuracies with $\mathrm{ERM}^{\mathrm{all}}$ model. The network architecture is ResNet-50.

task of learning data has become difficult for the model in 410 the presence of upsampled images. Since the embeddings 411 in DFR are dependent on the ERM-trained model, we also 412 analyze how the number of training epochs the ERM model 413 is trained for, impacts the DFR performance. The epochs for 414 retraining the last layer remain fixed to 1000, all other hy-415 perparameters also remain the same for DFR models trained 416 with different ERM-trained embeddings. The left plot in Fig-417 ure 8 indicates that, when the base model is fine-tuned for 20 418 epochs, the performance of DFR on the test set increases. As 419 the training time increases for ERM, performance by DFR 420 decreases, whereas the ERM model continues to improve. 421

In case of ERM, we also analyze the effect of the percent-422 age of training data in minority groups i.e. easy and hard 423 groups on model's test performance. We refer to ERM_{easy} 424 as the model that has been fine-tuned with data from the 425 majority group only i.e. 0% of data from medium and hard 426 group. Conversely, we refer to ERM_{all} as the model that has 427 been fine-tuned with equal data points from all the groups, 428 and ERM as the standard training data consisting of 20% of 429 data from minority groups. The results are depicted in the 430 right plot in Figure 8. We see that training with the Easy 431 group has worst performance on the Hard group. ERM_{all} 432 seems to narrow the gap between all groups. The accuracy 433 of the Easy group remains similar across the three models. 434

Table 4 shows the breakdown of accuracies for all the 435



Figure 9. Gradcam visualizations showing regions of the image the model pays attention to in order to make the classification decision. Labels in red show false predictions and labels in green indicate correct prediction.

436 sub-groups for ERM_{all} model. As compared to accuracies 437 shown in Figure 7, there is a considerable improvement in 438 case of CoR and CeR for size 56×56 and 84×84 . The 439 closest to clean accuracy for ResNet-50 is observed in case 440 of size 112×112 and group CeO.

56x56 Mushroom 112x112 Goldfish

441 **5.4.** Analysing Classifications with Saliency Maps

We use Gradcam to visualize the predictions on the ResNet-442 50 model. Figure 9 shows the visualizations on the ImageNet 443 pretrained model and two variations of ERM: ERM^{all} which 444 is fine-tuned with equal data points from all the groups and 445 ERM^{easy} which is fine-tuned only with images from the 446 Easy category, consisting of subgroups CeO and CoO for 447 size 54×54 and 112×112 respectively. The images on 448 449 the left side of Figure 9 show a stick insect of size 56×56 450 placed in the center against an outdoor environment. The pre-trained and ERM^{easy} model make their predictions by 451 picking up cues from the backgrounds and predicting class 452 Whiptail and Partridge respectively. Upon inspection, we 453 454 find that most of the images in these classes are set in similar 455 environments, hence the model has learnt to associate the 456 given outdoor environment with these classes and are ignoring the core features. ERMall, however, is more robust to 457 changes in environment and makes the correct prediction of 458 459 class *Stick Insect*. The images on the right show that, while 460 the pre-trained model is confused by the spurious cues in the background, ERM^{easy} makes the wrong predictions based 461 on the cues in the core features and the background together. 462 However, ERM^{all} makes the correct prediction by mostly 463 relying on core features. Figure 10 highlights the effect of 464 the size of core features on the ability of the ERM^{all} model 465 to make correct predictions. Having a smaller core feature 466 size results in the model making incorrect prediction of class 467 Goldfish. 468

6. Challenges and Future Work

The dataset variants of Hard-Spurious-ImageNet are proposed to understand the extent of background reliance as a
function of size and location of core features. One of the
limitations of the datasets is that they rely on ground truth

Figure 10. Effect of core feature size on model performance. A larger core feature makes the model ignore spurious cues in the background and surroundings. Both the predictions are for the $\rm ERM^{all}$ model.

bounding boxes of objects. In case of images where core 474 features are not labeled by bounding boxes, no inpainting is 475 performed on them, subsequently leading to core features 476 in background and foreground occurring simultaneously. 477 Moreover, the presence of secondary objects and clutter 478 in the background makes it difficult for the models to learn 479 small core feature sizes. The lack of segmentation bounding 480 boxes for all images in ImageNet restricted us to using ob-481 ject bounding boxes instead of masks. Currently, we have 482 only experimented with one location per core object. For 483 future work, we plan to experiment with different locations 484 of core objects in the images and analyze the impact of using 485 different network architectures with the dataset. Moreover, it 486 would be interesting to extend this analysis to other datasets 487 and models trained in different ways such as with contrastive 488 learning, and various data augmentation techniques. 489

7. Conclusion

In this paper, we propose a variant of ImageNet, Hard-491 Spurious-ImageNet, to help the deep learning community to 492 better understand spurious feature reliance. We show that 493 ImageNet is center-biased and exhibits a bias towards large 494 object sizes. We also provide an analysis showing that there 495 exists a negative correlation between size and location of 496 core features in an image and the strength of spurious cues in 497 the background. We experiment with different group robust-498 ness methods and highlight the need for specialized methods 499 to solve this problem. 500

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

References 501

- 502 [1] Valerio Biscione and Jeffrey Bowers. Learning translation 503 invariance in cnns. arXiv preprint arXiv:2011.11757, 2020. 2
- [2] Valerio Biscione and Jeffrey S Bowers. Convolutional neural 504 505 networks are not invariant to translation, but they can learn 506 to be. Journal of Machine Learning Research, 22(229):1-28, 507 2021. 2
- 508 [3] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. 509 Coatnet: Marrying convolution and attention for all data sizes. 510 Advances in neural information processing systems, 34:3965-511 3977, 2021. 5, 1
- 512 [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. 513 514 In 2009 IEEE conference on computer vision and pattern 515 recognition, pages 248-255. Ieee, 2009. 2
- 516 [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 517 Deep residual learning for image recognition. In Proceed-518 ings of the IEEE conference on computer vision and pattern 519 recognition, pages 770-778, 2016. 5, 1
- 520 [6] Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient for robustness to spuri-522 ous correlations. arXiv preprint arXiv:2204.02937, 2022. 1, 523 6,3
- 524 [7] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, 525 Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer White-526 head, Alexander C Berg, Wan-Yen Lo, et al. Segment any-527 thing. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 4015-4026, 2023. 4, 1 528
- 529 [8] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeva Man-530 galam, Bo Xiong, Jitendra Malik, and Christoph Feichten-531 hofer. Mvitv2: Improved multiscale vision transformers for 532 classification and detection. In Proceedings of the IEEE/CVF 533 conference on computer vision and pattern recognition, pages 534 4804-4814, 2022. 1
- 535 [9] Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghu-536 nathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness 537 538 without training group information. In International Confer-539 ence on Machine Learning, pages 6781-6792. PMLR, 2021. 540
- 541 [10] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao 542 Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, 543 Hang Su, Jun Zhu, and Lei Zhang. Grounding dino: Marrying 544 dino with grounded pre-training for open-set object detection. 545 arXiv preprint arXiv:2303.05499, 2023. 1
- 546 [11] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feicht-547 enhofer, Trevor Darrell, and Saining Xie. A convnet for the 548 2020s. In Proceedings of the IEEE/CVF conference on com-549 puter vision and pattern recognition, pages 11976-11986, 2022. 4, 1 550
- 551 [12] Aengus Lynch, Gbètondji JS Dovonon, Jean Kaddour, and 552 Ricardo Silva. Spawrious: A benchmark for fine control of 553 spurious correlation biases. arXiv preprint arXiv:2303.05470, 554 2023. 2
- 555 [13] Mazda Moayeri, Kiarash Banihashem, and Soheil Feizi. Ex-556 plicit tradeoffs between adversarial and natural distributional

robustness. Advances in Neural Information Processing Sys-557 tems, 35:38761-38774, 2022, 2 558

- [14] Mazda Moayeri, Phillip Pope, Yogesh Balaji, and Soheil 559 Feizi. A comprehensive study of image classification model 560 sensitivity to foregrounds, backgrounds, and visual attributes. 561 In Proceedings of the IEEE/CVF Conference on Computer 562 Vision and Pattern Recognition, pages 19087–19097, 2022. 2 563
- [15] Mazda Moayeri, Sahil Singla, and Soheil Feizi. Hard ima-564 genet: Segmentations for objects with strong spurious cues. 565 Advances in Neural Information Processing Systems, 35: 566 10068-10077, 2022. 2 567
- [16] Mazda Moayeri, Wenxiao Wang, Sahil Singla, and Soheil Feizi. Spuriosity rankings: sorting data to measure and mitigate biases. Advances in Neural Information Processing Systems, 36:41572-41600, 2023. 2
- [17] Nihal Murali, Aahlad Puli, Ke Yu, Rajesh Ranganath, and Kayhan Batmanghelich. Beyond distribution shift: Spurious features through the lens of training dynamics. Transactions on machine learning research, 2023, 2023. 2
- [18] Yannic Neuhaus, Maximilian Augustin, Valentyn Boreiko, and Matthias Hein. Spurious features everywhere-large-scale detection of harmful spurious features in imagenet. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 20235-20246, 2023. 2
- [19] Chaitanya Ryali, Yuan-Ting Hu, Daniel Bolya, Chen Wei, Haoqi Fan, Po-Yao Huang, Vaibhav Aggarwal, Arkabandhu Chowdhury, Omid Poursaeed, Judy Hoffman, et al. Hiera: A hierarchical vision transformer without the bells-and-whistles. In International Conference on Machine Learning, pages 29441-29454. PMLR, 2023. 1
- [20] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worstcase generalization. arXiv preprint arXiv:1911.08731, 2019. 1.4
- [21] Sahil Singla and Soheil Feizi. Salient imagenet: How to discover spurious features in deep learning? arXiv preprint arXiv:2110.04301, 2021. 2
- [22] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. arXiv preprint arXiv:2109.07161, 2021. 4
- [23] Mohammad Reza Taesiri, Giang Nguyen, Sarra Habchi, Cor-Paul Bezemer, and Anh Nguyen. Imagenet-hard: The hardest images remaining from a study of the power of zoom and spatial biases in image classification. Advances in Neural Information Processing Systems, 36, 2024. 3
- [24] Vladimir Vapnik. Principles of risk minimization for learning theory. Advances in neural information processing systems, 4, 1991. 6
- [25] Kai Xiao, Logan Engstrom, Andrew Ilyas, and Aleksander 608 Madry. Noise or signal: The role of image backgrounds in 609 object recognition. arXiv preprint arXiv:2006.09994, 2020. 610 2.3 611
- [26] Tao Yu, Runseng Feng, Ruoyu Feng, Jinming Liu, Xin Jin, 612 Wenjun Zeng, and Zhibo Chen. Inpaint anything: Seg-613

- 614ment anything meets image inpainting.arXiv preprint615arXiv:2304.06790, 2023.4
- 616 [27] Jessica Yung, Rob Romijnders, Alexander Kolesnikov, Lucas
 617 Beyer, Josip Djolonga, Neil Houlsby, Sylvain Gelly, Mario
- 618Lucic, and Xiaohua Zhai. Si-score: An image dataset for619fine-grained analysis of robustness to object location, rotation
- 620 and size. *arXiv preprint arXiv:2104.04191*, 2021. 2

Corner Cases: How Size and Position of Objects Challenge ImageNet-Trained Models

Supplementary Material



Figure 11. Histograms showing distribution of scores in different classes of train data in ImageNet1k dataset.

621 8. Benchmark Results

The results for Hard-Spurios-ImageNet and its variant are 622 given in Tables 5 and 6 respectively. We test the perfor-623 mance of the datasets on 5 different pre-trained architectures: 624 ConvNext-Base [11] trained on ImageNet21k and fine-tuned 625 626 on ImageNet1k, ResNet-50 [5], CoATNet [3], Hiera-Base with MAE [19], and MVit2-small [8]. Except for ConvNext, 627 628 these models are pretrained on ImageNet1k only. Across all models, the performance on Group CoR for size 56×56 629 is the worst. Benchmark results for different groups along 630 631 with clean accuracies are given in Tables 5 and 6. Clean accuracies in Table 6 are for 10 Hard-Spurious-ImageNet-v2 632 633 classes only.

634 9. Biases in ImageNet

Figure 11 shows the distribution of center and size scores
for different classes in the training data of ImageNet. We
calculate these scores using the available bounding boxes for
ImageNet training data. Figure 5 refers to the distribution
for the validation data.

640 **10. Inpaint Anything**

641 The predicted masks from Segment Anything are dilated by
642 a kernel size of 15 to avoid edge effects when the "hole" is
643 filled by LaMa. Some examples of the inpainted data are
644 given in Figure 12.

645 11. Hard-Spurious-ImageNet-v2

646 Despite inpainting, the background (Bg) consists of cues that647 help the model predict the background label (see Figure 13).



Figure 12. Original images with their resized inpainted versions.



Core: Japanese Spaniel Bg: Snorkel

Bg: Busby Hat

Figure 13. Examples from Hard-Spurious-ImageNet-v2.

Bg: Potter's Wheel

648

649

664

12. True Objects in Background

Ensuring that the backgrounds do not contain true objects 650 depends on the fidelity of provided ImageNet annotations. 651 We perform an additional analysis with a foundation model, 652 Grounding DINO [10], to extract bounding boxes from the 653 images. We consider similarity scores between Grounding 654 DINO predictions and the ImageNet annotations to analyze 655 the correctness of ImageNet annotations. For ImageNet 656 validation data, we get an overall mIOU of 0.8675 across 657 all classes between both sets of bounding boxes with 139 658 classes having mIOU value less than 0.8 (see Figure 14 659 for a histogram by mIOU). This shows that the majority of 660 the classes in ImageNet data have correct bounding boxes 661 and the amount of objects from the foreground class in the 662 background is negligible. 663

13. Hard-Spurious-ImageNet with SAM

We also experiment with using the Segment Anything [7] 665 model to obtain masks for the objects inside a bounding box and resize it to 3 different sizes (56, 84, and 112). The resized masks are then placed in the center and corner of the inpainted image, similar to the setting described in the 669

CVPR 2025 Submission #*****. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

Madal	Clean Object		Group Accuracies			
wiodei	Accuracy	Resolution	CeO	CoO	CeR	CoR
		56^{2}	59.62	50.14	20.15	14.98
Convnext-Base	85.86	84^{2}	74.42	70.05	50.81	46.47
		112^2	79.82	76.54	67.43	60.92
		56^{2}	45.79	36.74	11.83	6.78
ResNet-50	80.20	84^{2}	62.19	56.80	23.37	24.19
		112^2	72.43	70.22	55.85	55.62
		56^{2}	40.79	35.78	5.27	3.28
CoATNet	83.59	84^{2}	66.58	50.14	25.92	17.02
		112^2	72.70	72.51	45.45	44.44
		56^{2}	49.45	34.34	4.61	1.31
Hiera	84.48	84^{2}	67.64	55.09	21.81	12.49
		112^2	74.07	69.32	47.26	37.82
		56^{2}	41.44	31.38	5.41	1.38
MVitv2	83.77	84^{2}	67.38	51.12	29.53	14.17
		112^2	70.51	64.86	48.00	37.81

Table 5	Test Ac	curacies o	n Hard-	Spurious	-ImageNet
rable 5.	1030 110	curacies o	ii mara-	Spurious	-mager vet.

Madal	Clean Object		Group Accuracies			
Widdei	Accuracy	Resolution	CeO	CoO	CeR	CoR
		56^{2}	57.4	38.6	8.2	1.0
Convnext-Base	85.8	84^{2}	79.0	71.0	27.0	11.4
		112^{2}	83.2	79.8	45.4	17.0
		56^{2}	45.00	29.6	6.4	0.0
ResNet-50	82.2	84^{2}	70.8	58.4	17.0	9.8
		112^{2}	79.4	78.6	35.6	28.6
		56^{2}	20.9	21.8	0.8	0.0
CoATNet	83.4	84^{2}	71.8	49.0	11.0	4.0
		112^{2}	75.40	80.8	18.2	23.0
		56^{2}	46.8	22.8	1.0	0.0
Hiera	85.8	84^{2}	75.6	55.6	4.0	1.8
		112^{2}	78.6	74.6	16.0	10.6
MVitv2		56^{2}	29.0	15.0	0.4	0.0
	86.6	84^{2}	72.6	47.8	11.2	1.2
		112^{2}	72.4	66.0	18.4	12.6

Table 6. Test Accuracies on Hard-Spurious-ImageNet-v2 with highly spurious backgrounds.

670 main paper. At the moment, we only consider one object 671 per image. Since we have access to ImageNet-annotated bounding boxes, we use them as prompts to be given to SAM. 672 The results are shown in Table 7. Compared to the results in 673 Table 5, the results with SAM are worse, mainly because the 674 675 resized SAM object masks are not entirely accurate in cases where objects are small and thin, such as insects, etc. Hence, 676 677 we preferred human-annotated ImageNet bounding boxes.

14. Group Robustness Methods

We use pretrained ResNet-50 trained on ImageNet1k 679 for our experiments. The Base model is fine-tuned with 680 batch size 256, constant learning rate of 0.001 for 20 681 epochs. The input images are randomly cropped with an aspect ratio in the bounds (0.75, 1.33) and finally resized 683 to 224×224 . Horizontal flipping is applied afterward. 684

CVPR 2025 Submission #*****. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

Madal	Clean Object		Group Accuracies			
Iviodei	Accuracy	Resolution	CeO	CoO	CeR	CoR
		56^{2}	46.07	36.07	13.86	6.21
Convnext-Base	85.8	84^{2}	61.18	53.92	31.04	22.30
		112^2	67.78	64.69	42.91	13.84
		56^{2}	29.33	24.34	6.68	4.36
ResNet-50	82.2	84^{2}	45.17	40.63	19.09	16.24
		112^2	55.24	52.56	31.34	29.87
		56^{2}	30.57	27.61	7.91	3.93
CoATNet	83.4	84^{2}	50.94	44.66	21.03	15.63
		112^2	60.60	56.73	33.00	29.30
		56^{2}	37.94	25.88	9.08	2.92
MVit2	85.8	84^{2}	54.89	44.73	24.74	15.15
		112^2	63.73	57.94	36.80	30.60
		56^{2}	39.88	27.06	10.34	3.198
Hiera	86.6	84^{2}	56.36	46.18	25.13	15.72
		112^2	66.14	60.38	39.15	31.63

Table 7. Test Accuracies on Hard-Spurious-ImageNet with SAM Masks.



Figure 14. Class-wise mIOU scores between Grounding DINO predictions and ImageNet annotations on the validation set. Averaged mIOU is 0.875.

A momentum of 0.9 and weight decay of 0.001 is used. 685 For DFR, we normalize the embeddings using mean and 686 687 standard deviation of validation data used to train the last layer, and use the same statistics to normalize embeddings 688 of test data. We re-train the last layer for 1000 epochs, 689 learning rate of 1, cosine learning rate scheduler and SGD 690 691 optimizer with full-batch. We use ℓ_2 regularization with λ set to 100. These hyperparamters are similar to the ones 692 set by Kirichenko et al. [6] for optimizing the last layer for 693 ImageNet-9 dataset [25]. Since, the data distribution in the 694 proposed dataset and ImageNet-9 is similar, we assumed 695 the same hyperparamteres. In case of JTT, models have the 696 697 same hyperparameters as the ERM trained model. λ_{up} is set

Methods	Easy	Medium	Hard	Average
Pretrained	71.14	54.93	29.21	51.75
ERM	76.91	70.63	63.48	70.34
$\mathrm{ERM}^{\mathrm{easy}}$	77.82	68.33	51.39	65.85
DFR	74.82	68.66	61.68	68.39

Table 8. Test Performance of different methods on Easy, Medium, and Hard categories in Hard-Spurious-ImageNet. Average accuracy is the average test performance of all the groups combined. The model is Convnext-tiny.

to 50.

After extracting the embeddings from the pre-trained 699 ERM model, the embeddings are normalized us-700 ing fit_transform() and transform() functions of 701 sklearn.preprocessing.StandardScaler for val and test 702 data, respectively. For the JTT model, the images are applied 703 with random resized cropping followed by horizontal 704 No additional data augmentation is applied flipping. 705 afterward. We also experimented with ConvNext-tiny 706 pre-trained on ImageNet-22k and fine-tuned on ImageNet1k. 707 We fine-tune the pre-trained model on the proposed data 708 under various settings. ERM is trained by replicating 709 the long-tailed distribution of the data, while ERMeasy 710 is trained only with the easy group. ERM^{all} is trained 711 with equal data points from all groups. DFR is trained by 712 extracting embeddings from ERM, and re-training the last 713 layer only. The number of train and test images is similar to 714 the data setting described in the main paper. 715