
On the Effects of Data Distortion on Model Analysis and Training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Data modification can introduce artificial information. It is often assumed that
2 the resulting artefacts are detrimental to training, whilst being negligible when
3 analysing models. We investigate these assumptions and conclude that in some
4 cases they are unfounded and lead to incorrect results. Specifically, we show current
5 shape bias identification methods and occlusion robustness measures are biased
6 and propose a fairer alternative for the latter. Subsequently, through a series of
7 experiments we seek to correct and strengthen the community’s perception of how
8 distorting data affects learning. Based on our empirical results we argue that the
9 impact of the artefacts must be understood and exploited rather than eliminated.

10 1 Motivation

11 Modifying data has become commonplace both when training and analysing models, yet the wider
12 implications are often disregarded. We delve into some of the side-effects and point out that this
13 practice has resulted in the creation of biased model interpretation tools and poorly informed theories.
14 On the analysis side, we take as examples occlusion robustness and shape bias identification methods.
15 On the training side, we focus on some instances of Mixed Sample Data Augmentation (MSDA),
16 where two images are combined to obtain a new training sample. Visual examples of each can
17 be found in Figure 1. In this paper we study a number of assumptions that lie at the heart of the
18 aforementioned methods, which we briefly introduce below.

19 **Shape-texture bias:** Deep models are known to be sensitive to interventions that are imperceptible
20 to humans [35, 13], as well as to other forms of distribution shifts [1, 6, 8]. It has been argued that
21 this is intimately linked to networks tending to use texture rather than shape information [2, 11].
22 Recently, input distortions have become a popular way of assessing a model’s texture bias. To this end,
23 images are divided into a grid and the resulting patches are randomly shuffled such that information
24 is preserved locally, while the global shape is altered [32, 27, 25, 41]. It is implicitly assumed that
25 patch-shuffling does not introduce misleading shape or texture that could affect model evaluation.

26 **Occlusion robustness:** A widely adopted method for measuring occlusion robustness is through the
27 accuracy obtained after superimposing a rectangular patch on an image [5, 9, 39, 42, 20]. We refer
28 to this approach as CutOcclusion throughout the paper. Just as for shape bias, this method relies on
29 information introduced not to interfere with a model’s learnt representations such that a decrease in
30 performance can be directly attributed to lack of robustness.

31 **Data augmentation studies:** In statistical learning, training with augmented data is termed Vicinal
32 Risk Minimisation (VRM) [37, 4] and it is seen as injecting prior knowledge about the neighbourhood
33 of the data samples. The intuition behind augmentation caused researchers to interpret its effect
34 through the similarity between original and augmented data distributions. This perspective is often
35 challenged by methods which, despite generating samples that do not appear to fall under the

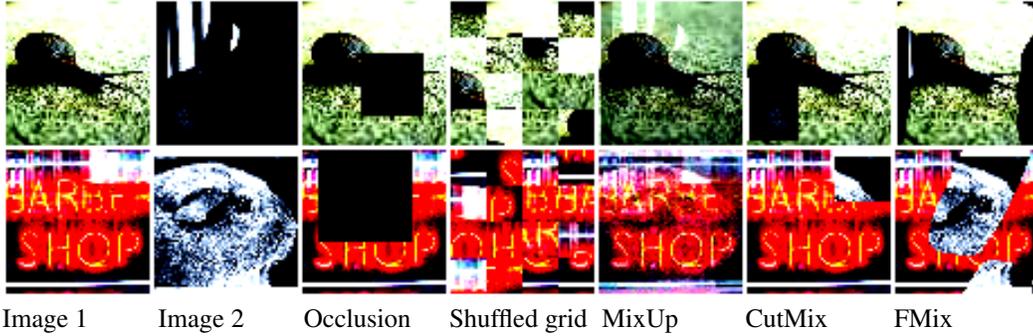


Figure 1: Examples of image distortions. For test-time distortions (Occlusion and Shuffled grid) only Image 1 was used. For mixing augmentations, the first row was generated with a mixing factor of 0.2, while the second one with 0.5.

36 distribution of natural images, lead to strong learners. Gontijo-Lopes et al. [12] argue it is the
 37 *perceived* distribution shift that needs to be minimised, while maximising the sample vicinity.
 38 Formalising these concepts, they introduce augmentation “diversity” and “affinity”. Diversity is
 39 defined as the training loss when learning with artificial samples, while affinity quantifies the
 40 difference between the accuracy on original test data and augmented test data for a reference model.
 41 The latter penalises augmentations that introduce artificial information to which the model is not
 42 invariant, implicitly assuming that training with that information is detrimental to generalisation.

43 In summary, it is currently assumed that the artefacts introduced by changes in the data are negligible
 44 when evaluating models, while those introduced when training are important and undesirable. Does
 45 the artificial information added by analysis methods not have major side-effects or does it lead to
 46 biased results? Conversely, are the artefacts important when training with modified data? Do they
 47 cause models to learn better or worse representations?

48 We set out to answer these questions and find that when the secondary effects of data manipulation
 49 are not accounted for, results can be misleading, especially in comparative studies. Subsequently, we
 50 construct empirical counter-examples which disprove common beliefs in the literature and highlight
 51 the importance of understanding the changes MSDAs introduce. Our contributions are:

- 52 • We show that increasingly popular model interpretation and analysis methods are biased,
 53 relying on unfounded assumptions (Section 2);
- 54 • For measuring occlusion robustness, we propose a fairer alternative (Section 3);
- 55 • We show that, in contrast to what is widely assumed, not preserving the data distribution
 56 can lead to learning better representations (Section 4).

57 2 Are artefacts negligible when analysing classifiers?

58 To verify whether distorting data at evaluation time could have side-effects previously not considered,
 59 we look at the increase in misclassifications per category. That is, from the number of incorrect
 60 predictions of a model evaluated on modified data, we subtract the incorrect predictions when testing
 61 on original data. If there is a significant increase for a specific class, it indicates that the distortion
 62 introduces features the model associates with that class. We refer to this phenomenon as “data
 63 interference”. Considering only positive differences, we denote the increase in the percentage of
 64 misclassifications for class c of a given model m by i_c^m . We define the Data Interference DI index as

$$DI = \frac{i_{c_{max}}}{\sum_c i_c} i_{c_{max}}, \quad (1)$$

65 where c_{max} is the class with highest mean increase across all runs. The DI index measures the
 66 percentage represented by the dominant class c_{max} weighted by its increase. A high index value
 67 indicates a sharp increase for a particular class which is consistent across runs. We associate this with
 68 an overlap between introduced artefacts and learnt representations. In Appendix B.1 we experiment

Table 1: DI index (%) for PreAct-ResNet18 on grid-shuffled images for four different types of models. Results with the highest average are given in italic and the lowest in bold. Information introduced when shuffling tends to interfere less with the representations of FMix and CutMix models.

	basic	MixUp	FMix	CutMix
CIFAR-10	<i>2.90\pm1.10</i>	2.54 \pm 1.29	0.60 \pm 0.26	0.33\pm0.17
CIFAR-100	<i>1.05\pm0.59</i>	0.93 \pm 0.44	0.23 \pm 0.29	0.11\pm0.10
FashionMNIST	1.12 \pm 0.63	<i>2.73\pm1.64</i>	1.12 \pm 0.61	0.70\pm0.22
Tiny ImageNet	<i>2.58\pm4.73</i>	0.54 \pm 0.27	0.38 \pm 0.12	0.14\pm0.12
ImageNet	0.82	<i>1.49</i>	0.58	–

69 with an alternative index, where we weight by the highest increase of a model across the 5 runs, so as
70 to obtain a worst-case analysis. As expected, we observe a more accentuated bias in this case.

71 To obtain models with different behaviours in a controlled manner, we make use of data augmentation.
72 Since it is sufficient to identify some common cases in which models are disfavoured, we choose to
73 reduce our environmental impact by restricting the analysis to simple MSDAs that combine images
74 without incurring additional computation time or external models. As will be argued in Section 3, we
75 expect the unfairness to be present in most settings, thus the exact choice of augmentation is irrelevant.
76 We focus on two popular MSDAs, MixUp [40] and CutMix [39]. MixUp linearly interpolates between
77 two images to obtain a new training example, while CutMix masks out a rectangular region of an
78 image with the corresponding region of another image. Besides the aforementioned methods, we
79 also employ FMix due to its irregularly shaped masks sampled from Fourier space, which will play
80 an important role in our analysis. Note that although the masking methods sample the size of the
81 occluding patch from the same distribution, in CutMix part of the rectangle can be outside the image,
82 which leads to less occluded samples overall. We refer to models by the augmentations they were
83 trained with and use “basic” to label the models trained without MSDA.

84 Throughout the paper, we do 5 runs of each experiment with PreAct-ResNet18 [15] as the default
85 architecture. We include results for BagNet [2] and VGG [33] in the Appendices. The main data sets
86 we report results on are CIFAR-10/100 [21], Tiny ImageNet [34], FashionMNIST [38], ImageNet [29].
87 For ImageNet we use pretrained ResNet-101 models made publicly available by Harris et al. [14].
88 Note that the only experiments for which we are unable to run repeats are those on ImageNet, since
89 only one model per augmentation is provided. For full experimental details, see Appendix A. Total
90 emissions of training the models evaluated in this paper are estimated [22] to be 38.07 kgCO₂eq, to
91 which 13.11 kgCO₂eq more are added during the analysis. The hope is that our findings will lead to a
92 better understanding which in the long run would help reduce erroneous research directions without
93 needing to empirically disprove them, lessening the future carbon footprint of the community.

94 2.1 Shape bias

95 For assessing shape bias through sample manipulation, the standard procedure is to choose between
96 dividing the image in 4, 16 or 64 patches to be shuffled. Since FashionMNIST images are smaller,
97 we choose a 2×2 grid, for CIFAR-10/100 and Tiny ImageNet 4×4 , while for ImageNet we use
98 an 8×8 grid. However, similar results are obtained for different grid sizes (Appendix B.2). The
99 large DI index in Table 1 indicates that either basic or MixUp models tend to associate the features
100 artificially introduced by patch-shuffling with a certain class. We take a closer look at the distribution
101 of misclassifications for CIFAR-10 and notice that the basic model tends to wrongly predict the class
102 “Truck” (Figure 2). This is not at all surprising, given that the strong horizontal and vertical edges are
103 highly indicative of this class. Similar observations can be made for other data sets (Appendix B.3).
104 Thus, we believe the grid-shuffling approach is causing models which are not invariant to strong
105 horizontal and vertical edges to appear to rely more heavily on shape information. A model not
106 affected by this transformation could be considered texture-biased if we accept the larger definition
107 of texture as local information. However, there is a question about the extent to which the reciprocal
108 is true; A model can be invariant to the aforementioned edges because it is indeed relying on texture
109 information or simply because it uses different shape-related features.

110 **Is a model necessarily more affected by patch-shuffling if it has a higher shape bias?** To answer
111 this question, we can use another method of determining shape and texture bias to find a counter-

Table 2: Accuracy of augmentation-trained ImageNet and Tiny ImageNet models on the GST data set when the label is taken to be either the shape or texture. There is no clear correlation between masking methods and low texture bias.

	ImageNet		Tiny ImageNet	
	Shape	Texture	Shape	Texture
basic	20.31	53.28	10.56 \pm 0.65	26.04 \pm 1.77
MixUp	24.14	60.31	12.02 \pm 0.33	27.77 \pm 1.56
FMix	21.25	53.43	10.40 \pm 0.39	19.90 \pm 2.12
CutMix	—	—	10.54 \pm 0.38	23.72 \pm 2.42

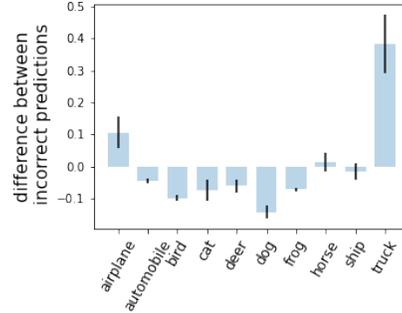


Figure 2: Difference in incorrect predictions on CIFAR-10 for the basic model.

112 example. We analyse the ImageNet models on the Geirhos Style-Transfer (GST) [11] data set. GST
 113 contains artificially generated images where the shape belongs to one class and the texture to another.
 114 There are 16 coarse classes that encompass a number of ImageNet categories to which they are
 115 mapped. The bias of the models is given by the accuracy obtained when the label is set to either the
 116 shape or texture. Using this well-known method of identifying shape bias we want to find models
 117 which have similar biases but different DI indices when patch-shuffling. This would indicate that
 118 sensitivity to shuffling is not necessarily linked to increased shape bias.

119 The results in Table 2 show that the basic model does not have a higher shape bias than masking
 120 methods although it has a significantly higher DI index, as we have seen in Table 1. We repeat the
 121 same experiment on the Tiny ImageNet [34] data set. Geirhos et al. [11] use WordNet [26] to map the
 122 1000 categories to the 16 classes of the GST data set. A number of ImageNet categories that belong
 123 to the 16 higher-level classes of GST are missing. For this reason, a poorer overall performance is
 124 expected and the results could differ slightly given a better fit between the sets. Nonetheless, we find
 125 again no significant correlation between masking augmentation and texture bias. We also include in
 126 Table B.2 the results for BagNet models, which have smaller receptive fields and so are forced to use
 127 more local information. Even in this case, we find a high DI for the basic model and no difference
 128 in texture bias compared to MSDA. Thus, a model which is more affected by patch-shuffling is not
 129 necessary more shaped-bias. In other words, models can appear to have vastly different shape bias
 130 when evaluated on randomly rearranged patches, albeit in reality their bias is similar. The inability of
 131 this method to account for artefacts makes it unfair and unreliable.

132 2.2 Occlusion measurement

133 We want to determine whether the same issue identified in the case of shape bias evaluation applies
 134 to occlusion robustness measures. We focus on CutOcclusion, where a rectangular black patch is
 135 superimposed on test images and the robustness is given by the resulting accuracy. We perform
 136 the same experiment as before, where the *DI* index is now measured when testing on rectangle-
 137 occluded images. There is no standardised distortion when measuring CutOcclusion, with the size
 138 and positioning of the obstructing patch varying between studies. Most often in prior art a lack of
 139 robustness is noted for large occluders [e.g. 5, 42]. For this reason, we sample the size of the patch
 140 from a Beta(2,1) distribution, allowing the occluding patch to lie outside the image (as it is done for
 141 augmenting with CutMix and CutOut [7]). This allows us to capture both the cases in which either
 142 the centre or the border area is masked out but requires a non-uniform distribution to counter for the
 143 patches existing outside the image. We also sample from a uniform distribution where the occluder is
 144 restricted to be positioned within the image boundaries and obtain similar results (See Appendix C.1).

145 Table 3 shows that a significant gap in the DI index can be identified for each of the data sets. This
 146 indicates that some models will again be disadvantaged. We additionally find that data interference is
 147 present for different architectures, when overlapping patches from external images or using differently
 148 shaped masks (Appendix B.4). Thus, the result of CutOcclusion and its variants is highly dependent
 149 on the problem at hand. That is, whether the artefacts introduced by the artificial occlusion are
 150 salient features of the model depends on what features are naturally distinctive for the model. Just as
 151 for randomly shuffling tiles, by occluding images using a particularly shaped patch, one implicitly

Table 3: DI index (%) when occluding with black patches. The highest results are given in italic and the lowest in bold. For each data set, there exists a non-negligible gap in the DI index.

	basic	MixUp	FMix	CutMix
CIFAR-10	<i>5.48</i> ± 1.31	0.75 ± 0.55	0.87 ± 0.92	2.87 ± 3.04
CIFAR-100	<i>3.25</i> ± 1.31	0.80 ± 0.46	1.28 ± 2.48	1.36 ± 0.63
FashionMNIST	0.44 ± 0.34	<i>2.59</i> ± 1.05	0.92 ± 1.79	0.97 ± 1.82
Tiny	<i>2.40</i> ± 0.82	1.88 ± 0.61	0.47 ± 0.41	4.62 ± 4.93
Imagenet	1.28	<i>4.50</i>	1.02	—

152 measures a model’s affinity to certain features, albeit those features might be discriminative. This
 153 deems such methods inappropriate for fairly assessing robustness and texture bias.

154 A related observation was made by Hooker et al. [17] who note the pitfalls of manipulating data to
 155 determine feature importance. They point out that when simply superimposing uniform patches over
 156 image features, it is difficult to assess how much of the reduction in accuracy is caused by the absence
 157 of those features and how much is due to images becoming out of distribution. To address this, the
 158 most important features identified by an estimator are masked out both on train and test data, closing
 159 the gap between the two sets. Hooker et al. then train and evaluate models on the newly generated
 160 images. Unlike for interpretability methods, the subject of occlusion robustness studies is the model
 161 itself, which makes training with a modified version of the data an inviable option. In the following
 162 section we explore ways of overcoming this bias when measuring occlusion robustness.

163 3 What are fairer alternatives?

164 We propose a simple, more carefully defined measure that aims to decouple the machine’s edge bias
 165 from the occlusion robustness, which we refer to as “interplay occlusion” (iOcclusion). Interplay
 166 occlusion reflects the change in the interplay between performance on seen and unseen data. Formally,

$$iOcclusion_i = \left| \frac{\mathcal{A}(\mathcal{D}_{train}^i) - \mathcal{A}(\mathcal{D}_{test}^i)}{\mathcal{A}(\mathcal{D}_{train}) - \mathcal{A}(\mathcal{D}_{test})} \right|, \quad (2)$$

167 where $\mathcal{A}(\mathcal{D})$ denotes the accuracy on a given data set \mathcal{D} , and \mathcal{D}^i is the data set resulting from
 168 removing $i\%$ pixels of each image. The intuition is that on train data robust models are less sensitive
 169 to the artefacts of the occlusion policy for small levels of occlusion, resulting in a large difference in
 170 accuracy from that on unseen data. The performance of both train and test gets close to random as
 171 the percentage of occluded data approaches 90% and we expect the gap to fall off quicker for less
 172 robust models. This change in interplay is taken with respect to the generalisation gap of the model,
 173 such that the quality of the model fit in itself does not interfere with the robustness measure.

174 Although iOcclusion reduces data interference, other factors have to also be considered when choosing
 175 a masking method for computing \mathcal{D}^i , such as the number of contiguous components or the amount of
 176 salient information masked out. In this paper we choose to generate masks using Grad-CAM [30],
 177 such that the area with most salient $i\%$ pixels is covered. It must be noted that this method implicitly
 178 assumes there could be multiple occluders and has the downside of incurring a higher environmental
 179 cost. For this, we also experiment with using rectangular or Fourier-sampled masks and conclude
 180 that although random masking makes the process noisier, the exact choice of masking method is of
 181 secondary importance as long as the occluder’s granularity is accounted for. Appendix C.4 provides
 182 discussion and results on these alternative instances of iOcclusion, as well as differences in their
 183 carbon footprint. For a fair comparison, throughout this section we do not allow the obstructing patch
 184 when measuring CutOcclusion to lie outside the image such that the fraction removed is exact.

185 Assessing the correctness of such a measure is difficult in the absence of a baseline. For the remainder
 186 of this section we will build varied experiments to attest the validity of our method. We focus on the
 187 key results, but include additional ones and further experiments in Appendix C. Since occlusion in
 188 real-life scenarios could be caused by non-uniformly coloured objects, an appropriate measure must
 189 generalise across colour patterns. When computing iOcclusion and CutOcclusion, we superimpose
 190 patches from images belonging to a different data set and compare the results to those obtained
 191 when occluding with black patches only. For visual clarity, Figure 3 presents the results for the

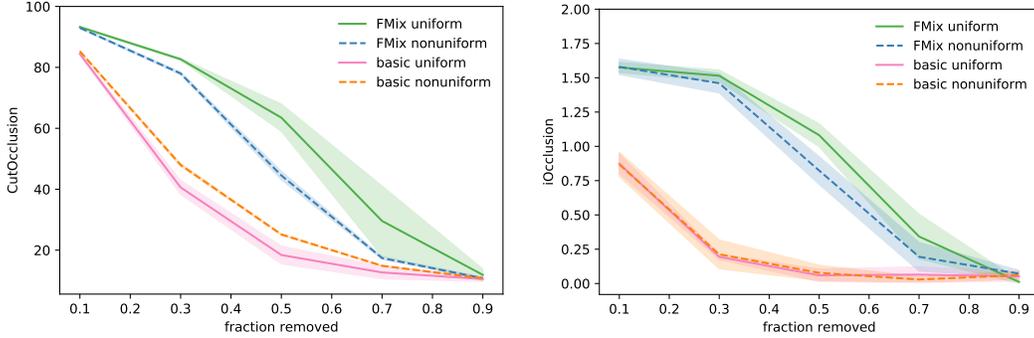


Figure 3: CutOclusion (left) and iOclusion (right) when occluding with black patches (uniform) and patches taken from other images (nonuniform). iOclusion gives more consistent results.

Table 4: DI index (%) and occlusion robustness for models trained on CIFAR-10 when obstructing 30% of the image pixels with non-uniform patches. When measuring the robustness with CutOclusion, RM appears significantly less robust than CutMix due to its sensitivity to patching with rectangles, while iOclusion highlights the robustness specific to training with FMix-like masks. Given in bold is the closest result to that of RM for each evaluation.

	basic	MixUp	CutMix	FMix	RM
DI index	1.67 ± 0.27	1.00 ± 0.31	0.17 ± 0.16	0.15 ± 0.03	0.39 ± 0.07
CutOclusion	47.97 ± 0.52	58.65 ± 1.01	76.56 ± 6.36	78.00 ± 0.45	60.79 ± 5.03
iOclusion	0.21 ± 0.10	0.57 ± 0.18	1.09 ± 0.17	1.46 ± 0.07	1.20 ± 0.23

192 least and most robust models (see Appendix C.2 for a full comparison). For iOclusion, using
 193 uniform occluders gives similar results to its non-uniform version, whereas the CutOclusion measure
 194 provides an inconsistent model evaluation.

195 As we have argued, in addition to not being sensitive to the colour pattern of the patch, a fair measure
 196 must also be invariant to the shape of the patch. To empirically confirm iOclusion reduces the
 197 importance of edge information, we aim to obtain a model that is robust to occlusion, but at the same
 198 time has a high DI index (it is sensitive to edge information). To this end, we create a variation of
 199 FMix, Random Masks (RM), where at the beginning of the training process three masks are randomly
 200 sampled from Fourier space. For each batch, one of the three is chosen uniformly at random. While
 201 the RM models are not sensitive to black-patch occlusion, when masking with patterned patches they
 202 have a higher DI index than FMix, as desired. Table 4 gives results for a fraction of 0.3 pixels covered
 203 by a non-uniform occluder. Our measure reflects the robustness of training with RM, situating it close
 204 to other masking methods. On the other hand, because CutOclusion implicitly penalises models
 205 with high DI index, according to this measure RM appears almost as sensitive to occlusion as MixUp.
 206 Figure B.4 shows results for a wider range of fractions.

207 Another problem that occurs when purely looking at post-masking accuracy is weaker models would
 208 erroneously appear less robust. We show this by reversing the problem: we evaluate the same model
 209 on two different subsets of the CIFAR-100 data set: typical and tail images as categorised by Feldman
 210 and Zhang [10]. They consider a train-test sample pair to belong to the tail of the data distribution if
 211 the test sample is correctly classified when a model is trained with the train sample and incorrectly
 212 without it. CutOclusion would indicate that models are significantly more robust to occluding typical
 213 examples. However, a closer analysis makes us doubt this conclusion. The raw accuracy on both
 214 train and test data for tail examples is lower than for the typical ones. In fact, the performance when
 215 occluding images decreases at the same rate for the two subsets. By way of definition, iOclusion
 216 allows a fair comparison of robustness regardless of the overall performance of a model (Figure 4).

217 As we evidenced through controlled experiments, there are many cases that CutOclusion does not
 218 properly address. From a model analysis perspective, correctly assessing the occlusion robustness
 219 could lead to better understanding and development of models and training procedures. Equally
 220 important, it has applicability for real-world deployments where no prior knowledge exists about the

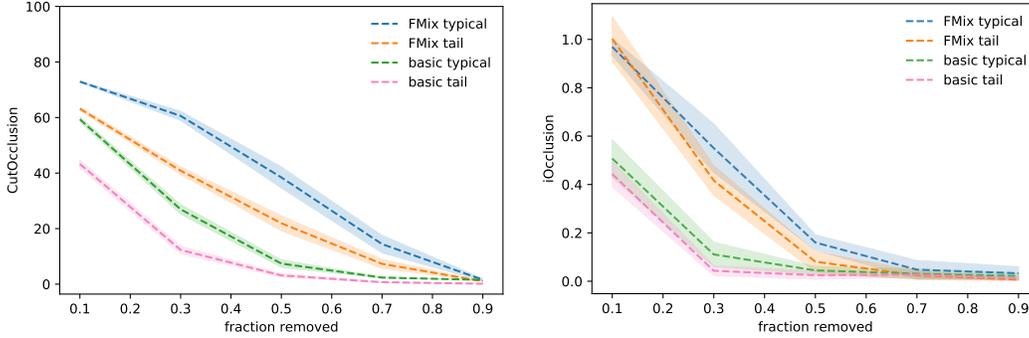


Figure 4: CutOcclusion (left) and iOcclusion (right) for the basic and FMix models on two subsets of the same data set: tail and typical. Evaluating the models with iOcclusion on the two types of samples leads to mostly overlapping robustness levels. That is, they do not differ outside the margin of error. On the contrary, CutOcclusion incorrectly finds the models to be less robust on tail data.

Table 5: Augmentation comparison on CIFAR-10. We consider two variants when calculating diversity. One is computing the cross-entropy loss using the label of the majority class (Diversity), as for mixing in [19]. The alternative, MixDiversity, takes a linear combination of the two cross-entropy losses.

	Affinity	Diversity	MixDiversity
MixUp	-12.58 ± 0.14	0.41 ± 0.01	0.84 ± 0.00
FMix	-25.55 ± 0.26	0.34 ± 0.01	0.65 ± 0.00

Table 6: DI index (%) measured for non-uniform occlusion when training without the class with highest increase in incorrect predictions. Again, a gap can be noted, supporting the idea that data interference is not specific to peculiar cases.

	CIFAR-10	CIFAR-100
MixUp	0.39 ± 0.15	1.22 ± 0.19
FMix	0.08 ± 0.06	0.50 ± 0.21

221 possible shapes of the obstructions. While this aspect of generality is the strength of our approach,
 222 it must be stressed that when there exists a limited set of known possible occluders, evaluating
 223 robustness specifically to them could be safer. Incorrectly assessing robustness can have severe effects
 224 especially when applied to autonomous vehicles or medical imaging. We do not propose a universal
 225 solution, but rather suggest an alternative to the biased approach for the common scenario in which
 226 the environment is not controlled and little is known about all the potential occluders. However, even
 227 in this case, our metric should be taken as a guide when analysing models. Although iOcclusion aims
 228 to address data interference, since a ground-truth does not exist, it cannot be guaranteed that this
 229 method provides fair results in the absolute.

230 The strength of the bias will depend on the data in question and some applications will be more heavily
 231 affected than others. We have seen that for natural images this bias does exist. To confirm that we have
 232 not just identified isolated cases, we remove the class that has the highest increase in mispredictions
 233 and retrain the models on the remaining classes. We find that the bias is again present (Table 6),
 234 but with respect to another class. For example, in the case of CIFAR-10, after removing the “Truck”
 235 class, models mispredict rectangle-occluded images as “Boat” (Appendix D). Thus, the edge artefacts
 236 are very likely to interfere with learnt representations since they are such fundamental features. From
 237 an evaluation perspective, as we have seen, this impacts assessment methods and must be accounted
 238 for. From the training perspective, such a widespread data interference of masking distortions would
 239 indicate a large perceptual shift in the data when performing MSDA. In the following section we
 240 investigate the importance of the artefacts in this case and their implications.

241 4 Is the magnitude of the distribution shift important?

242 Traditionally it was believed that a good augmentation should have minimal distribution shift. Most
 243 recently, it has been argued that it is the degree of the perceived shift that determines augmentation
 244 quality [12]. We start with the perceptual gap of training with MSDA, as proposed in Gontijo-Lopes
 245 et al. [12]. Reiterating, this is given by the difference between the performance of the baseline model
 246 when presented with original test data and augmented test data and is termed “affinity”. Subsequently,

247 we address the gap in the wider sense, as is often sought in prior art. We first argue that high affinity
248 and high diversity are not necessarily desirable. Indeed, on CIFAR-10, we find FMix, a better
249 performing augmentation, to have both lower affinity and lower diversity than MixUp (Table 5). For
250 diversity, we compute the cross-entropy loss where the label is taken to be that of the majority class.
251 Similar results are obtained with the MixUp loss, where a weighted average of the true labels is taken.

252 While intuitively for a high level of affinity, high diversity could correspond to better methods,
253 the converse does not hold. We argue this is because affinity is rather an analysis of the learnt
254 representations of the reference model and cannot give an insight into the quality of the augmentation
255 or its effect on learning. The limitations of affinity are intimately linked to those of CutOcclusion.
256 We have seen in Section 3 that the bias of the basic model is present not only when obstructing an
257 image with a uniform patch, but also when mask-mixing. As such, an augmentation will have a
258 lower affinity if it introduces artefacts that could otherwise lead to learning better representations
259 when used in the training process. We believe this issue extends to other approaches that aim to
260 motivate the success of MSDA through reduced distribution shift. Henceforth, we focus on bringing
261 further supporting evidence that the importance lies in the invariance introduced by the shift and its
262 interaction with the given problem rather than its magnitude.

263 **4.1 If it is not the magnitude that matters, is it the direction?**

264 We use empirical evidence to argue against previous assumptions behind the success of MSDA and
265 propose the study of introduced bias as a more informative research direction. Here we use the term
266 “bias” to refer to a drift in the learnt representations introduced by the change in the training procedure.
267 A fundamental difference to classical training is that in the case of augmentation the samples are no
268 longer independent. Mixed-sampling takes this even further. An immediate question is, does the
269 added correlation lead to more meaningful representations? It is claimed that the strength of MixUp
270 lies in causing the model to behave linearly between two images [40] or in pushing the examples
271 towards their mean [3]. Both of these claims rely on the combined images to be generated from the
272 same distribution. We want to verify to which extent this is necessary for a successful augmentation.

273 It has been argued that label mixing has a negligible effect on the final model performance [19, 18,
274 14, 24]. We use the reformulated objective setting [18, 14], where targets are not mixed and the
275 mixing coefficient is drawn from an imbalanced Beta distribution. This allows us to apply MSDA
276 between data sets. Thus, for training a model on a data set, we use an additional one whose targets
277 will be ignored. As an example, a model that is learning to predict CIFAR-10 images will be trained
278 on a combination of CIFAR-10 and CIFAR-100 images, with the target of the former. This scenario
279 breaks the added correlation between training examples. Note that when mixing between data sets
280 we use the same procedure as when performing regular MSDA, without improving the process.

281 Table 7 contains the results of this experiment, showing that an accuracy similar to or better than that
282 of regular MSDA can be obtained by performing inter-dataset MSDA. This invalidates the argument
283 that the power of MixUp resides in causing the model to act linearly between samples. Another
284 observation is that for FMix and MixUp, introducing elements from CIFAR-100 when training
285 models on the CIFAR-10 problem does not harm the learning process. The reciprocal, however, does
286 not hold. Hence, the “distribution shift” is more intimately linked to the problem at hand and aiming
287 to characterise an augmentation based on the distance from the original distribution is a limiting
288 approach, especially when the distance is measured as perceived by a reference model.

289 We believe an explanation is that the artefacts created when putting together images from CIFAR-10
290 with those of CIFAR-100 could introduce information that makes the separation of the 10 classes
291 easier. However, if the information happens to interfere with a feature that is important for separating
292 the CIFAR-100 categories, the performance could degrade on this data set. This singular experiment
293 is not sufficient to draw any general conclusions. However, it does show that shifting two distributions
294 by the same amount can have different effects on the model performance. Thus, the specifics of the
295 bias introduced could be more important than its magnitude. While some level of data similarity has
296 to be preserved when performing MSDA, it is far from being the objective of such data-distorting
297 approaches which, as we will argue further, should be rather seen as forms of regularisation.

298 We have seen that for all considered data sets, artefacts introduced by masking methods seem to
299 overlap with common features. This has led us to believe that MSDA training could help bypass
300 some of the simplicity bias. The simplicity bias refers to the tendency of deep models to find simple

Table 7: Accuracy on CIFAR-10 (left) and CIFAR-100 (right) upon mixing with samples from a different data set. The baseline is the accuracy when training with a single data set using the reformulated objective. In the interest of space, CIFAR-110 is used to refer to mixing with CIFAR-100 when training on the CIFAR-10 problem and vice-versa.

	MixUp	FMix	CutMix	MixUp	FMix	CutMix
baseline	94.18 \pm 0.34	94.36 \pm 0.28	94.67 \pm 0.20	74.68 \pm 0.37	75.75 \pm 0.31	74.19 \pm 0.50
CIFAR-110	94.70 \pm 0.27	94.80 \pm 0.32	94.66 \pm 0.12	72.36 \pm 1.04	74.80 \pm 0.55	74.47 \pm 0.39
Fashion	92.28 \pm 0.28	95.03 \pm 0.10	94.61 \pm 0.19	66.40 \pm 1.86	74.46 \pm 0.57	74.06 \pm 0.28

301 representations and has been used to justify the success of deep models [28, 36]. Recent research
 302 shows that this propensity causes models to ignore complex features that explain the data well in
 303 favour of elementary features, even when they lead to worse performance [31, 16].

304 Although it could seem natural that since MSDAs are not augmentations in the VRM sense, they will
 305 increase the complexity of the problem, we design an experiment to support this claim. Similarly
 306 to Shah et al. [31], we combine CIFAR-10 and MNIST [23] samples. Since they have the same
 307 number of classes, we can easily associate each class of one data set with a corresponding one from
 308 the other. Thus, we stack a padded image from the k th class of MNIST on top of a sample from the
 309 k th class of CIFAR-10, such that a $3 \times 64 \times 32$ image is obtained. We then randomly combine the
 310 test images and separately compute the accuracy with respect to the targets of each data set.

311 The predictions with respect to the CIFAR-10 labels are no better than random (10.04 ± 0.11), while
 312 the accuracy with respect to the MNIST images remains high (99.57 ± 0.72). Thus, models trained
 313 on this combination are mostly relying on MNIST images to make predictions. Similar behaviours
 314 have previously been associated with simplicity bias. Subsequently, when training, we perform FMix
 315 only on MNIST images and observe that this is enough to reverse the results. Evaluating against the
 316 CIFAR-10 label gives an accuracy of 86.60 ± 0.34 , while testing against the MNIST label only gives
 317 11.61 ± 0.30 . We find that this also holds true for the other MSDAs. Thus, performing these distortions
 318 on the simpler data set increases its complexity to the point where it surpasses that of CIFAR-10.

319 Previously, we presented evidence that masking MSDA does not necessarily promote learning neither
 320 more shape nor texture information. In the light of this fact along with the results from this section,
 321 we believe image distortions force the model to learn more complex both shape and texture-specific
 322 features. Thus, in this paper we pointed out that the shift in learnt representations can lead to better
 323 models and simply quantifying the distribution shift can be misleading. An open question remains:
 324 How can we better capture the bias that is introduced and its quality? We believe understanding how
 325 a relatively small change in the data distribution impacts learnt representations could lead the way to
 326 characterising the relationship between data and model generalisation.

327 5 Conclusions

328 Distorting data is such a commonplace procedure, yet little effort has been devoted to investigating its
 329 broader effects. This is particularly problematic when image modifications are applied in analyses. We
 330 show a number of cases in which this leads to *biased results*. For occlusion robustness measurement,
 331 we propose an alternative. The insights we gain from this endeavour point towards the study of data
 332 characteristics as a cornerstone of our understanding and raise a number of questions about mixed
 333 sample data augmentation, on which we subsequently focus. We note that they interfere with features
 334 that are consistently found across a number of data sets and conclude that the methods commonly
 335 used are forms of mixed sample *regularisation* rather than augmentation. A limitation of previous
 336 studies that aim to explain their success is the focus on trying to argue similarity with original data,
 337 rather than explaining the bias introduced by the distortion. Correctly interpreting it is important
 338 not only for making the models trustable but also for injecting more informed prior knowledge in
 339 future applications. Beyond their practical benefits, we believe MSDAs have the potential to help
 340 characterise the interplay between data and learnt representations. Overall, the purpose of our paper
 341 is to encourage better practice when dealing with all forms of data distortions.

References

- 342
- 343 [1] Rocío Alaiz-Rodríguez and Nathalie Japkowicz. Assessing the impact of changing environments
344 on classifier performance. In *Conference of the Canadian Society for Computational Studies of*
345 *Intelligence*, pages 13–24. Springer, 2008.
- 346 [2] Wieland Brendel and Matthias Bethge. Approximating CNNs with bag-of-local-features models
347 works surprisingly well on imagenet. In *International Conference on Learning Representations*,
348 2019. URL <https://openreview.net/forum?id=SkfMWhAqYQ>.
- 349 [3] Luigi Carratino, Moustapha Cissé, Rodolphe Jenatton, and Jean-Philippe Vert. On mixup
350 regularization. *arXiv preprint arXiv:2006.06049*, 2020.
- 351 [4] Olivier Chapelle, Jason Weston, Léon Bottou, and Vladimir Vapnik. Vicinal risk minimization.
352 In *Advances in neural information processing systems*, pages 416–422, 2001.
- 353 [5] Sanghyuk Chun, Seong Joon Oh, Sangdoon Yun, Dongyoon Han, Junsuk Choe, and Youngjoon
354 Yoo. An empirical evaluation on robustness and uncertainty of regularization methods. *arXiv*
355 *preprint arXiv:2003.03879*, 2020.
- 356 [6] David A Cieslak and Nitesh V Chawla. A framework for monitoring classifiers’ performance:
357 when and why failure occurs? *Knowledge and Information Systems*, 18(1):83–108, 2009.
- 358 [7] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural
359 networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- 360 [8] Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry.
361 Exploring the landscape of spatial robustness. In *International Conference on Machine Learning*,
362 pages 1802–1811. PMLR, 2019.
- 363 [9] Alhussein Fawzi and Pascal Frossard. Measuring the effect of nuisance variables on classifiers.
364 In Edwin R. Hancock Richard C. Wilson and William A. P. Smith, editors, *Proceedings of the*
365 *British Machine Vision Conference (BMVC)*, pages 137.1–137.12. BMVA Press, September
366 2016. ISBN 1-901725-59-6. doi: 10.5244/C.30.137. URL [https://dx.doi.org/10.5244/](https://dx.doi.org/10.5244/C.30.137)
367 [C.30.137](https://dx.doi.org/10.5244/C.30.137).
- 368 [10] Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering
369 the long tail via influence estimation. *Advances in Neural Information Processing Systems*, 33,
370 2020.
- 371 [11] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann,
372 and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias
373 improves accuracy and robustness. In *International Conference on Learning Representations*,
374 2019. URL <https://openreview.net/forum?id=Bygh9j09KX>.
- 375 [12] Raphael Gontijo-Lopes, Sylvia J Smullin, Ekin D Cubuk, and Ethan Dyer. Affinity and diversity:
376 Quantifying mechanisms of data augmentation. *arXiv preprint arXiv:2002.08973*, 2020.
- 377 [13] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
378 examples. *arXiv preprint arXiv:1412.6572*, 2014.
- 379 [14] Ethan Harris, Antonia Marcu, Matthew Painter, Mahesan Niranjan, Adam Prügel-Bennett, and
380 Jonathon Hare. Understanding and enhancing mixed sample data augmentation. *arXiv preprint*
381 *arXiv:2002.12047*, 2020.
- 382 [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual
383 networks. In *European conference on computer vision*, pages 630–645. Springer, 2016.
- 384 [16] Katherine Hermann and Andrew Lampinen. What shapes feature representations? exploring
385 datasets, architectures, and training. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan,
386 and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages
387 9995–10006. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/](https://proceedings.neurips.cc/paper/2020/file/71e9c6620d381d60196ebe694840aaaa-Paper.pdf)
388 [paper/2020/file/71e9c6620d381d60196ebe694840aaaa-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/71e9c6620d381d60196ebe694840aaaa-Paper.pdf).

- 389 [17] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for inter-
390 pretrainability methods in deep neural networks. In *Advances in Neural Information Processing*
391 *Systems*, pages 9737–9748, 2019.
- 392 [18] Ferenc Huszár. mixup: Data-dependent data augmentation, 2017. URL [http://www.
393 inference.vc/mixup-data-dependent-data-augmentation/](http://www.inference.vc/mixup-data-dependent-data-augmentation/).
- 394 [19] Hiroshi Inoue. Data augmentation by pairing samples for images classification. *arXiv preprint*
395 *arXiv:1801.02929*, 2018.
- 396 [20] Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan
397 Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, et al. Captum: A
398 unified and generic model interpretability library for pytorch. *arXiv preprint arXiv:2009.07896*,
399 2020.
- 400 [21] Alex Krizhevsky et al. Learning multiple layers of features from tiny images. 2009.
- 401 [22] Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying
402 the carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*, 2019.
- 403 [23] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL [http:
404 //yann.lecun.com/exdb/mnist/](http://yann.lecun.com/exdb/mnist/).
- 405 [24] Daojun Liang, Feng Yang, Tian Zhang, and Peter Yang. Understanding mixup training methods.
406 *IEEE Access*, 6:58774–58783, 2018.
- 407 [25] Tiange Luo, Tianle Cai, Mengxiao Zhang, Siyu Chen, Di He, and Liwei Wang. Defective
408 convolutional layers learn robust CNNs. *arXiv preprint arXiv:1911.08432*, 2019.
- 409 [26] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):
410 39–41, 1995.
- 411 [27] Chaithanya Kumar Mummadi, Ranjitha Subramaniam, Robin Huttmacher, Julien Vitay, Volker
412 Fischer, and Jan Hendrik Metzen. Does enhanced shape bias improve neural network robustness
413 to common corruptions? In *International Conference on Learning Representations*, 2021. URL
414 <https://openreview.net/forum?id=yUxUNaj2S1>.
- 415 [28] Preetum Nakkiran, Gal Kaplun, Dimitris Kalimeris, Tristan Yang, Benjamin L Edelman, Fred
416 Zhang, and Boaz Barak. SGD on neural networks learns functions of increasing complexity.
417 *arXiv preprint arXiv:1905.11604*, 2019.
- 418 [29] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
419 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.
420 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*
421 (*IJCV*), 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- 422 [30] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi
423 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based
424 localization. In *Proceedings of the IEEE international conference on computer vision*, pages
425 618–626, 2017.
- 426 [31] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The
427 pitfalls of simplicity bias in neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M. F.
428 Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,
429 pages 9573–9585. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.
430 cc/paper/2020/file/6cfe0e6127fa25df2a0ef2ae1067d915-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/6cfe0e6127fa25df2a0ef2ae1067d915-Paper.pdf).
- 431 [32] Baifeng Shi, Dinghuai Zhang, Qi Dai, Zhanxing Zhu, Yadong Mu, and Jingdong Wang. Infor-
432 mative dropout for robust representation learning: A shape-bias perspective. In *International*
433 *Conference on Machine Learning*, pages 8828–8839. PMLR, 2020.
- 434 [33] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
435 image recognition. In *International Conference on Learning Representations*, 2015.

- 436 [34] Stanford. Tiny imagenet visual recognition challenge, 2015. URL [https://tiny-imagenet.](https://tiny-imagenet.herokuapp.com/)
437 [herokuapp.com/](https://tiny-imagenet.herokuapp.com/).
- 438 [35] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfel-
439 low, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*,
440 2013.
- 441 [36] Guillermo Valle-Perez, Chico Q. Camargo, and Ard A. Louis. Deep learning generalizes because
442 the parameter-function map is biased towards simple functions. In *International Conference on*
443 *Learning Representations*, 2019. URL <https://openreview.net/forum?id=rye4g3AqFm>.
- 444 [37] Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media,
445 1999.
- 446 [38] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for
447 benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 448 [39] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon
449 Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In
450 *Proceedings of the IEEE International Conference on Computer Vision*, pages 6023–6032,
451 2019.
- 452 [40] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond
453 empirical risk minimization. In *International Conference on Learning Representations*, 2018.
454 URL <https://openreview.net/forum?id=r1Ddp1-Rb>.
- 455 [41] Tianyuan Zhang and Zhanxing Zhu. Interpreting adversarially trained convolutional neural
456 networks. In *International Conference on Machine Learning*, pages 7502–7511. PMLR, 2019.
- 457 [42] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data
458 augmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, number 34(07),
459 pages 13001–13008, 2020.

460 Checklist

- 461 1. For all authors...
- 462 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
463 contributions and scope? [Yes]
- 464 (b) Did you describe the limitations of your work? [Yes] See lines 174–184, 221–229,
465 292–294.
- 466 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See lines
467 221–229 as well as 89–93 and Appendix C.4.
- 468 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
469 them? [Yes]
- 470 2. If you are including theoretical results...
- 471 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 472 (b) Did you include complete proofs of all theoretical results? [N/A]
- 473 3. If you ran experiments...
- 474 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
475 mental results (either in the supplemental material or as a URL)? [Yes]
- 476 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
477 were chosen)? [Yes] See Appendix A.
- 478 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
479 ments multiple times)? [Yes] The only exception is when using publicly available
480 models, as mentioned in Section 2.
- 481 (d) Did you include the total amount of compute and the type of resources used (e.g., type
482 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix A.

- 483 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 484 (a) If your work uses existing assets, did you cite the creators? [Yes] See lines 84–93 and
- 485 Appendix A.
- 486 (b) Did you mention the license of the assets? [N/A]
- 487 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 488
- 489 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 490 using/curating? [N/A]
- 491 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 492 information or offensive content? [N/A]
- 493 5. If you used crowdsourcing or conducted research with human subjects...
- 494 (a) Did you include the full text of instructions given to participants and screenshots, if
- 495 applicable? [N/A]
- 496 (b) Did you describe any potential participant risks, with links to Institutional Review
- 497 Board (IRB) approvals, if applicable? [N/A]
- 498 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 499 spent on participant compensation? [N/A]