
Bias Amplification in Image Classification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent research suggests that predictions made by machine-learning models can
2 amplify biases present in the training data. Mitigating such bias amplification
3 requires a deep understanding of the mechanics in modern machine learning that
4 give rise to that amplification. We perform the first systematic, controlled study
5 into when and how bias amplification occurs. To enable this study, we design a
6 simple image-classification problem in which we can tightly control (synthetic)
7 biases. Our study of this problem reveals that the strength of bias amplification
8 is correlated to measures such as model accuracy, model capacity, and amount
9 of training data. We also find that bias amplification can vary greatly during
10 training. Finally, we find that bias amplification may depend on the difficulty of
11 the classification task relative to the difficulty of recognizing group membership:
12 bias amplification appears to occur primarily when it is easier to recognize group
13 membership than class membership. Our results suggest best practices for training
14 machine-learning models that we hope will help pave the way for the development
15 of better mitigation strategies.

16 1 Introduction

17 Several recent studies have presented results suggesting that, beyond reproducing biases present
18 in the training data, machine-learning models can *amplify* such biases as well [11, 34, 38]. Bias
19 amplification is concerning as it can foster the proliferation of undesired stereotypes [9, 32, 38, 37]
20 or lead to unjustifiable differences in model accuracy between subgroups of users [5, 8].

21 The existence of bias amplification suggests that machine-learning models are not always doing what
22 we expect them to do: *viz.*, make predictions according to the statistics present in their training data.
23 Although several studies have proposed measures for the severity of bias amplification [11, 25, 34, 38],
24 this question of when and why bias amplification occurs remains largely unanswered.

25 We present a systematic, controlled study of bias amplification. We design a simple image-
26 classification task that facilitates tight control of synthetic biases. In line with prior work [11, 34, 38],
27 we find that models trained for this classification task, indeed, amplify biases present in their training
28 data. We use the ability to control biases to study key research questions (RQs) that increase our
29 understanding of bias amplification:

- 30 • *RQ1*: How does bias amplification vary as the bias in the data varies?
- 31 • *RQ2*: How does bias amplification vary as a function of model capacity?
- 32 • *RQ3*: How does bias amplification vary during model training?
- 33 • *RQ4*: How does bias amplification vary as a function of the relative difficulty of recognizing
34 class membership versus recognizing group membership?

35 We observe that bias amplification tends to increase with bias in the training set in many of our
36 experiments. We find that bias amplification varies with model capacity: models with more parameters

37 and/or less regularization can amplify biases, but models with too few parameters and/or too much
38 regularization can amplify biases even more. Bias amplification also greatly varies with training
39 set size: models trained on very small or very large training sets appear to amplify biases less. We
40 observe that the degree of bias amplification can vary greatly during model training. In many of
41 our experiments, we find that the behavior of bias amplification depends on the difficulty of the
42 classification task relative to the difficulty of group membership recognition.

43 The results of our study provide intuitions for when bias amplification occurs and why. They suggest
44 some best practices that may help reduce bias mitigation in real-world machine-learning models, such
45 as careful cross-validation of hyperparameters related to model capacity, regularization, and training
46 duration to substantially reduce bias amplification of the final model. Collecting more training data
47 may reduce bias amplification as well. We hope that our study helps pave the way for the development
48 and adoption of mitigation strategies for bias amplification in common computer vision tasks.

49 2 Experiments

50 We design an image-classification task in which each image has both a class and a group, and in
51 which we can introduce synthetic biases by altering the group assignment of images.

52 2.1 Experimental setup

53 **Classification task.** We perform image-classification experiments on three image datasets: (1) the
54 Fashion MNIST dataset [36], (2) the CIFAR-10 dataset [21], and (3) the CIFAR-100 dataset [21].
55 Because our analyses are easier to perform with binary classification problems, we convert the datasets
56 to have binary labels by randomly selecting half of the classes to be positive and the other half to
57 be negative. All our models are residual networks [16] that we train according to the procedures
58 discussed in Appendix A.1.1. To mitigate the effect of a particular random class assignment, we
59 average the test accuracy and bias amplification values over 20 random assignments of the original
60 classes to the binary classes and include 95% confidence intervals.

61 **Group membership.** We consider a classification model to be *biased* if it predicts a particular *class*
62 at a disproportionate rate for examples from a particular *group*. However, rather than using real-world
63 groups we choose to establish synthetic groups for our experiments. This reduces the noise in our
64 measurements and allows us to perform additional root-cause investigations of bias amplification that
65 may pose ethical or technical challenges with real-world groups, such as determining the model’s
66 ease of predicting the group in the image.

67 Specifically, we create two groups in our image-classification problems by *inverting* some of the
68 images in a dataset and not inverting others. Using image inversion to create groups has two main
69 advantages over other synthetic methods like color changes or adding random noise: (1) it hardly
70 introduces new visual features into the images that may alter the image-classification problem and (2)
71 it is straightforward for image-recognition models to recognize whether or not an image is inverted.¹
72 This allows us to tightly control the correlation between classes and groups without introducing
73 causal relations between them. Figure 6 shows examples of inverted and non-inverted images.

74 **Controlling dataset bias.** For all images corresponding to a single task in the input dataset, we
75 randomly select positively labeled images with rate $1/2 - \epsilon$ and invert them, and we randomly invert
76 negatively labeled images with rate $1/2 + \epsilon$ (we choose $\epsilon \in [0, 1/2]$). This leads to a bias of strength 2ϵ
77 in the dataset: If $\epsilon = 0$, image inversion (*i.e.*, group membership) carries no information on whether
78 the images has a positive or a negative label (*i.e.*, class membership). By contrast, group membership
79 uniquely defines class membership when $\epsilon = 1/2$. Hence, $\epsilon = 0$ corresponds to an *unbiased* dataset
80 in which group membership does not carry information about class membership, $\epsilon = 1/2$ corresponds
81 the a *fully biased* setting in which group membership uniquely determines class membership, and
82 values of $\epsilon \in (0, 1/2)$ correspond to *partly biased* datasets.

¹In preliminary experiments, we found that the test accuracy of a residual network trained to recognize image inversion is 100% on Fashion MNIST images and 96% on CIFAR-100 images.

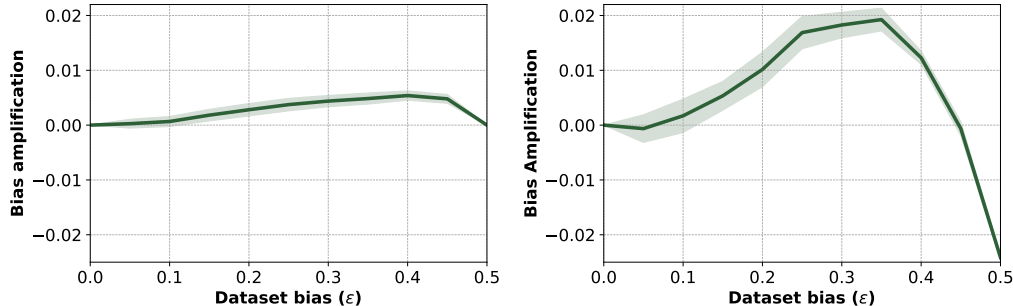


Figure 1: Bias amplification, $\text{BiasAmp}_{A \rightarrow T}$, as a function of the degree of bias, ϵ , for (left) ResNet-18 models trained on the Fashion MNIST dataset and (right) ResNet-110 models trained on the CIFAR-100 dataset. Shaded regions indicate the 95% confidence intervals over 20 independent experiments.

83 **Bias amplification measure.** We adopt the directional bias amplification measure $\text{BiasAmp}_{A \rightarrow T}$
 84 from [34]. This measure disambiguates different types of bias amplification and accounts for varying
 85 base rates of group membership.

86 The value of $\text{BiasAmp}_{A \rightarrow T}$ is 0 if the model predictions are exactly as biased as the labels in the
 87 dataset. If $\text{BiasAmp}_{A \rightarrow T}$ is negative, the model predictions dampen the bias present in the dataset
 88 and a positive $\text{BiasAmp}_{A \rightarrow T}$ value indicates that the model predictions amplify the bias in the dataset.
 89 Refer to Appendix A.1.2 for more details about the metric.

90 2.2 Results

91 We present the results of our experiments organized by research question (RQ). In accompanying
 92 materials, we also discuss the effect of training dataset size on bias amplification (Appendix A.3) and
 93 the relationship between overconfidence and bias amplification (Appendix A.4).

94 **RQ1: How does bias amplification vary as the bias in the data varies?** We perform experiments
 95 in which we vary the amount of bias in the Fashion MNIST dataset by generating training and test
 96 sets with different levels of bias, *i.e.*, by varying ϵ .

97 The results from this experiment (left pane of Figure 1) show that when the training set is unbiased
 98 ($\epsilon = 0$), no bias amplification occurs because no bias is present. No bias amplification occurs when
 99 the training set is fully biased ($\epsilon = 1/2$), as it is impossible to amplify an already maximum bias.
 100 However, for intermediate ϵ values (*i.e.*, in partially biased training sets), the trained models amplify
 101 the bias present in the training data. Bias amplification generally *increases* with the amount of bias in
 102 training data, until the bias in the data is nearly maximized ($\epsilon = 0.5$).

103 We repeat the same experiment on the CIFAR-100 dataset with ResNet-110 models. The results
 104 (in the right pane of Figure 1) show a similar pattern. A notable difference, however, is that bias
 105 amplification is negative when the CIFAR-100 dataset is maximally biased ($\epsilon = 1/2$). We surmise
 106 that this happens because the group membership of CIFAR-100 images cannot always be recognized
 107 correctly by a model. To obtain zero bias amplification at $\epsilon = 1/2$, a model needs to be a perfect
 108 predictor of group membership. Hence, when the model incorrectly recognizes the group membership
 109 of some of the images, a negative bias amplification (*i.e.*, bias dampening) is obtained.

110 **RQ2: How does bias amplification vary as a function of model capacity?** It is well-known that
 111 the capacity of machine-learning models influences their classification performance. To understand
 112 how model capacity impacts bias amplification, we perform experiments in which we measure bias
 113 amplification while adjusting the capacity of our models. We adjust model capacity in three ways: (1)
 114 via the *depth* of the model; (2) via the *width* of the model; and (3) via the *regularization* of the model.

115 We focus on the CIFAR-100 dataset here because CIFAR-100 images are harder to classify than
 116 Fashion MNIST images: this makes it more likely that models with different capacities will produce
 117 substantially different predictions. We use the ResNet-110 model from our RQ1 experiments as our

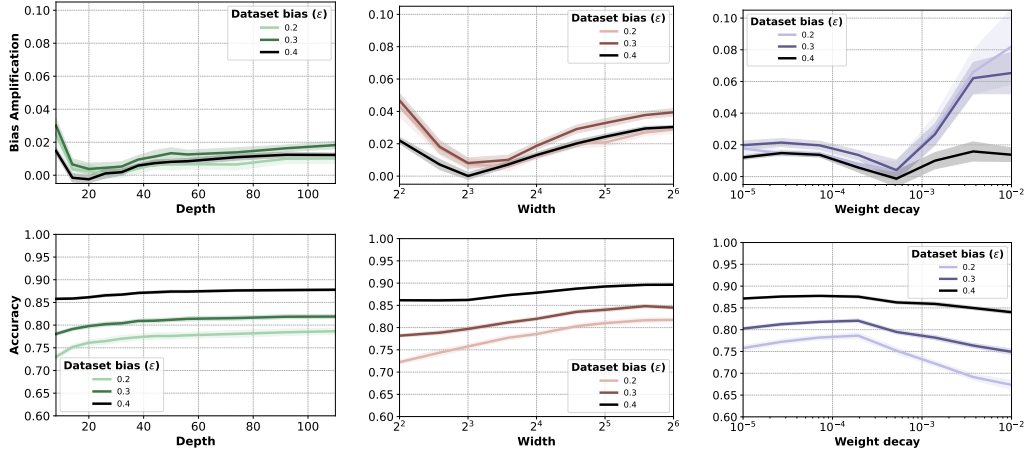


Figure 2: Bias amplification (**top**) and test accuracy (**bottom**) on the CIFAR-100 dataset as a function of three measures of model capacity. Each line represents a different amount of bias (ϵ) in the training set. Shaded regions indicate the 95% confidence intervals across 20 models. **Left:** Results for varying model depths. **Middle:** Results for varying model widths. **Right:** Results for varying weight decays.

118 base model. We experiment with depths that range between 8 and 110, widths ranging between 4 and
 119 64, and logarithmically spaced weight decay values between 10^{-5} to 10^{-2} . As before, we vary the
 120 dataset bias, ϵ , between 0 and $1/2$. The top row of Figure 2 shows the results of these experiments.

121 Irrespective of whether we vary model depth, width, or weight decay, the results suggest that bias
 122 amplification follows a “v-shape”: it increases when model capacity increases beyond a certain level,
 123 but it also increases when model capacity is reduced below a certain level. We surmise there are
 124 different explanations for these two increases. When the capacity of a model is limited, it needs to
 125 rely on features that are easy to extract when making class predictions. When the dataset is biased
 126 ($\epsilon > 0$), the model thus relies on image inversion, which is easy to recognize, in its class predictions.
 127 This explains why bias amplification is relatively large when the model has low capacity.² In contrast,
 128 when the capacity of a model is large, bias amplification may increase because the model has the
 129 capacity to extract both features that indicate class membership and features that indicate group
 130 membership. This allows the model to use group membership features to increase the confidence of
 131 its predictions, which reduces the training loss.³

132 The relation between model capacity and bias amplification resembles the well-known relation
 133 between model capacity and *generalization error*. Models with insufficient capacity have high
 134 generalization error because they cannot model the data distribution well, whereas high-capacity
 135 models may have high generalization error due to *overfitting*. Our results suggests that there exists a
 136 model-capacity “sweet spot” in which bias is minimally amplified, akin to model-capacity sweet spot
 137 that minimizes generalization error (for a given training set).

138 To investigate whether the optimal model capacities for bias amplification and generalization error
 139 coincide, we plot the test accuracy of our models in the bottom row of Figure 2. Test accuracy
 140 increases monotonically with model depth and width, suggesting that the (overall) optimal model is
 141 larger than the range of models we experimented with. However, we do observe that a weight decay
 142 of $1.9 \cdot 10^{-4}$ appears optimal for test accuracy. This weight-decay value is smaller than the value that
 143 minimizes bias amplification ($5.2 \cdot 10^{-4}$), which suggests that model designers may sometimes have
 144 to trade off bias amplification and accuracy when tuning hyperparameters.

145 **RQ3: How does bias amplification vary during model training?** Thus far, we have only mea-
 146 sured bias amplification of models that were trained until convergence for 500 epochs. While it is

²This explanation relies on the assumption that *group membership* features are relatively easy to extract and, hence, that our observations may change had we not used image inversion to construct our synthetic groups. We investigate how the difficulty of recognizing group membership influences bias amplification in RQ5.

³Indeed, we surmise the increase in bias amplification in very high-capacity models is related to the tendency of such models to be overconfident [14]; we investigate this relation further in Appendix A.4.

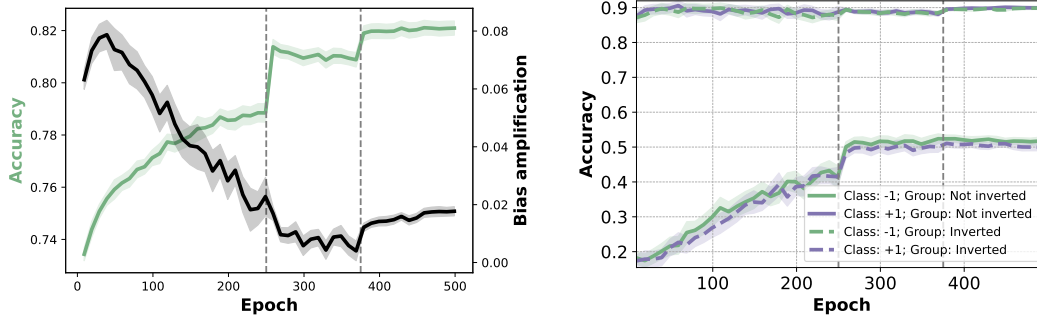


Figure 3: **Left:** Bias amplification and accuracy of ResNet-110 models during training on the CIFAR-100 dataset with a bias of $\epsilon = 0.3$. **Right:** Accuracy of the same models’ per class-group combination. Shaded regions indicate the 95% confidence intervals across 50 models. Vertical dashed lines indicate epochs at which the learning rate of the mini-batch SGD optimizer is decreased.

147 feasible to train models until convergence on small datasets, it may not be practical to do so on very
 148 large training sets. To evaluate how much bias amplification of a model varies during training, we
 149 measure bias amplification during the training of ResNet-110 models on a version of the CIFAR-100
 150 dataset with bias $\epsilon = 0.3$.

151 The left pane of Figure 3 plots bias amplification and accuracy as a function of training epoch in this
 152 setting. The results in the figure show that bias amplification *varies greatly* during training; models
 153 amplify biases much more strongly in the early stages of training. Bias amplification gradually
 154 declines as training proceeds and the recognition accuracy of the model increases. However, the bias
 155 amplification increases again slightly in the final stages of training, in particular, after the learning
 156 rate is decreased to its smallest value. Notably, bias amplification appears to increase slightly every
 157 time the learning rate is decreased.

158 To better understand what drives these changes in bias amplification during training, we disaggregate
 159 the model’s test accuracy into the four group-task combinations in the right pane of Figure 3. The
 160 model very quickly achieves high accuracy on examples for which the class label, $\{-1, +1\}$, matches
 161 the corresponding majority group, $\{\text{inverted}, \text{not inverted}\}$, per the bias in the dataset. By contrast,
 162 the accuracy on examples for which the class label does not match the majority group is very low in
 163 the initial stages of learning and increases much more gradually during training. We surmise this
 164 happens because group membership (image inversion) is easier to recognize than class membership
 165 (CIFAR-100 binary label). In the early stages of training, the model rapidly picks up on the easy-to-
 166 detect group membership signal as it provides the fastest way to reduce the model’s loss. In turn, this
 167 leads to bias amplification because the model makes predictions based on group membership signals
 168 whilst ignoring class membership signals. As training progresses, the group membership signal loses
 169 value because it is not a perfect predictor of class membership (note that $\epsilon = 0.3$). Hence, the model
 170 starts to utilize more class membership signals as training progresses, which results in an increase in
 171 accuracy and a decrease in bias amplification.

172 To test this hypothesis, we perform an experiment in which we swap the role of the group and the
 173 class: *i.e.*, the class label now represents whether or not the image is inverted and the group label
 174 depends on the object depicted in the CIFAR-100 image. Indeed, we find that bias is dampened in the
 175 early stages of training as the model latches onto the easy-to-extract class membership signal first, but
 176 this largely disappears in the later stages of training as the model starts to utilize group membership
 177 signals for recognition as well. See Appendix A.2 for more.

178 **RQ4: How does bias amplification vary as a function of the relative difficulty of recognizing**
 179 **class membership versus recognizing group membership?** Hitherto, we repeatedly observed
 180 that bias amplification may depend on the relative difficulty of recognizing class membership versus
 181 recognizing group membership: as the group signal is easier to extract in our setup, models amplify
 182 bias more in early stages of training and/or when they have lower capacity. We perform a more
 183 detailed study of this relationship.



Figure 4: Example of different amounts of overlay ($\eta = 0.0, 0.2, 0.5, 0.8, 1.0$) for an example belonging to the “airplane” class and “bird” group. Only class information is visible when $\eta = 0.0$ (left); only group information is visible when $\eta = 1.0$ (right).

184 We alter our problem setup such that we can control the relative difficulty of class recognition and
 185 group recognition. We abandon our image-inversion setup and, instead, create datasets that contain a
 186 convex combination of two CIFAR-10 images: a “group image” and a “class image”. By changing the
 187 weight of the convex combination, we can make the group image or the class image more prominent
 188 in the resulting image, thereby altering the difficulty of recognizing the class and the group.

189 We create the two groups, a and b , by randomly choosing two CIFAR-10 classes that we sample
 190 group images from. We also randomly choose two CIFAR-10 classes to form the binary classification
 191 task (*i.e.*, one class is the positive class and the other the negative class). Next, we create an example
 192 by sampling a class image, $\mathbf{I}_{\text{class}}$, from one of the two classes and a corresponding group image,
 193 $\mathbf{I}_{\text{group}}$, from one of the two groups. We linearly mix these two images:

$$\mathbf{I} = \eta \mathbf{I}_{\text{group}} + (1 - \eta) \mathbf{I}_{\text{class}}, \quad (1)$$

194 where $\eta \in [0, 1]$ is a mixing parameter and the final example \mathbf{I} is assigned the label of $\mathbf{I}_{\text{class}}$. Figure 4
 195 shows an example of the resulting examples for different η values. As before, we assign positive
 196 examples to group a with probability $0.5 + \epsilon$ or to group b with probability $0.5 - \epsilon$. Negative examples
 197 are assigned group b with probability $0.5 + \epsilon$, and to group a with probability $0.5 - \epsilon$.

198 When $\eta = 0$, this task reduces to classifying two
 199 classes from the standard CIFAR-10 images as
 200 the model cannot observe the group image at all.
 201 Conversely, directly recognizing class membership
 202 is impossible when $\eta = 1$ but recognizing
 203 group membership is easy in that setting. Hence,
 204 η provides a knob that facilitates varying the
 205 relative difficulty of recognizing group member-
 206 ship versus class membership. Additionally, this
 207 new combinatorial method provides insight into
 208 bias amplification when there are specific visual
 209 features associated with individual sub-groups.

210 Figure 5 presents results of experiments in
 211 which we measure bias amplification as a func-
 212 tion of the trade-off parameter, η , for different
 213 degrees of bias, ϵ . It shows that bias is damp-
 214 ened when it is relatively difficult to recognize
 215 group membership (*i.e.*, when η is low). When
 216 η increases past the point where group infor-
 217 mation is more visible than class information
 218 ($\eta = 0.5$), however, the bias amplification starts
 219 to progressively increase and becomes positive
 220 for larger η . This observation provides additional evidence for the hypothesis that bias amplification
 221 depends heavily on the relative difficulty of recognizing group membership versus class membership.

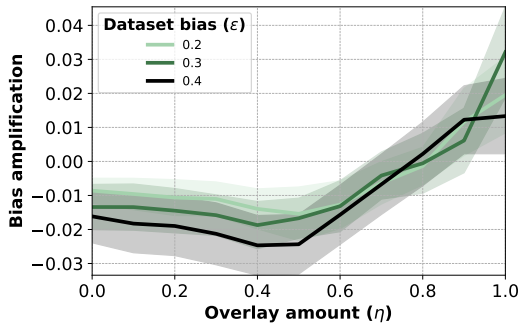


Figure 5: Bias amplification as a function of the relative difficulty of predicting class and group membership, η , for three different levels of bias, ϵ . Recognizing class membership is easier for small η values; recognizing group membership is easier for large η values. Shaded regions indicate 95% confidence intervals across 20 models.

222 3 Related work

223 This study is one of many studies on fairness and bias amplification in machine-learning models.

224 **Fairness.** Prior work has introduced a large number of formulations of fairness, including equalized
225 odds and equalized opportunity [15], fairness through awareness [10] or unawareness [13, 22],
226 treatment equality [1], and demographic parity [10, 22]. Measures associated with these fairness
227 formulations include differences in accuracy [1], differences in true or false positive rate [7, 15], and
228 the average per-class accuracy across subgroups [5]. These measures differ from bias amplification
229 measures in that they focus on correlations in the model predictions, whereas bias amplification
230 focuses on *differences* between the correlations in the training data and those in the model predictions.
231 In other words, bias-amplification measures discern between bias that is adopted from the training
232 data and bias that is amplified by the model; fairness measures make no such distinction.

233 **Bias amplification.** The study of bias amplification is of interest because it allows us to study how
234 design choices in our models, training algorithms, *etc.* contribute to bias in machine-learning models
235 beyond biases in the training data [17]. Prior work has measured bias amplification using generative
236 adversarial networks [6, 18], by considering binary classifications without attributes [24], and by
237 measuring correlations in model predictions [19, 38]. In our work, we use the BiasAmp_{A→T} measure
238 from [34], which addressed shortcomings in prior work [38], to measure bias amplification. Bias
239 amplification has also been studied in the context of causal statistics [2, 26, 29, 30, 35], but that line
240 of work has remained disparate from the study of bias amplification in machine learning. Despite
241 the plethora of prior work on measuring bias amplification, little is known on *when and how* bias
242 amplification arises in machine-learning models supporting vision tasks. Our study is among the first
243 to shed some light on the context under which bias amplification occurs.

244 4 Discussion

245 The results of our experiments shed light on the conditions under which bias amplification can
246 occur in machine-learning models for vision tasks. In particular, we find that bias amplification
247 varies as a function of bias in the dataset, model capacity, training time, and the amount of training
248 data. We also find that bias amplification depends on the relative difficulty of recognizing class
249 membership and recognizing group membership. This creates a predicament as the Bayes error of
250 those two recognition tasks are generally beyond the control of the model developer. Moreover, the
251 model developer may not always be able to measure the difficulty of recognizing group membership
252 empirically as doing so may involve developing a model that predicts sensitive attributes—something
253 that model developers may want to avoid [20, 23, 31, 34].

254 Although our study does not resolve this predicament, it may provide some useful best practices
255 to mitigate bias amplification as much as possible during model development. Our result suggests
256 that there is value in using cross-validation to carefully select a model architecture, regularizer, and
257 training recipe that minimizes bias amplification. Model developers may reduce bias amplification
258 using the same tuning process that they routinely use to minimize classification error. Our study
259 provides intuitions for how key levers available to the model developer can affect bias amplification.
260 However such tuning does require access to sensitive attribute values, *viz.* group-membership
261 information, and our study does not provide a complete overview of how all relevant levers influence
262 bias amplification. We intend to perform a more comprehensive investigation in future work.

263 **Limitations.** While our study provides useful insights and suggests best practices, it also suffers
264 from several key limitations. It is limited to binary classification tasks in the image-recognition
265 domain and uses synthetic indicators of group membership. Further work is needed to understand
266 how our findings apply to real world groups and biases, and how bias amplification manifests in
267 different vision tasks and other modalities. We note that in recommendation tasks especially, bias
268 amplification may arise in more complex ways because such systems generally have a human-in-the-
269 loop influencing the behavior of the system [3].

270 Another limitation of our study is that it only studies bias *amplification*, which requires tradeoffs
271 between other fairness guarantees and performance measures. Eliminating undesired biases altogether
272 and ensuring fair, optimal performance thus requires careful design of the entire pipeline from data
273 collection to model deployment.

274 **References**

- 275 [1] R. Berk, H. Heidari, S. Jabbari, M. Kearns, and A. Roth. Fairness in criminal justice risk
276 assessments: The state of the art. *Sociological Methods and Research*, 50(1):3–44, 2021.
- 277 [2] J. Bhattacharya and W.B. Vogt. Do instrumental variables belong in propensity scores?, 2007.
- 278 [3] Leon Bottou, Jonas Peters, Joaquin Quinonero-Candela, Denis X. Charles, D. Max Chickering,
279 Elon Portugaly, Dipankar Ray, Patrice Simard, and Ed Snelson. Counterfactual reasoning and
280 learning systems: The example of computational advertising. *Journal of Machine Learning*
281 *Research*, 14:3207–3260, 2013.
- 282 [4] S. Boucheron, O. Bousquet, and G. Lugosi. Theory of classification: a survey of recent advances.
283 *ESAIM: Probability and Statistics*, 9:323–375, 2005.
- 284 [5] J. Buolamwini and T. Gebru. Gender shades: Intersectional accuracy disparities in commercial
285 gender classification. In *Fairness, Accountability, and Transparency (FAT)*, 2018.
- 286 [6] K. Choi, A. Grover, T. Singh, R. Shu, and S. Ermon. Fair generative modeling via weak
287 supervision. In *Proceedings of the International Conference on Machine Learning (ICML)*,
288 2020.
- 289 [7] A. Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction
290 instrument. *Big Data*, 5(2):153–163, 2016.
- 291 [8] T. DeVries, I. Misra, C. Wang, and L.J.P. van der Maaten. Does object recognition work for
292 everyone? In *CVPR Workshop on Computer Vision for Global Challenges*, 2019.
- 293 [9] Emily Dinan, Angela Fan, Adina Williams, Jack Urbanek, Douwe Kiela, and Jason Weston.
294 Queens are powerful too: Mitigating gender bias in dialogue generation. In *Proceedings of the*
295 *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020.
- 296 [10] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel. Fairness through awareness. In
297 *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, 2012.
- 298 [11] J. Foulds, R. Islam, K.N. Keya, and S. Pan. An intersectional definition of fairness. In
299 *Proceedings of the International Conference on Data Engineering (ICDE)*, pages 1918–1921,
300 2020.
- 301 [12] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola,
302 Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: training
303 imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- 304 [13] N. Grgic-Hlaca, M.B. Zafar, K. P. Gummadi, and A. Weller. The case for process fairness
305 in learning: Feature selection for fair decision making. In *NeurIPS Symposium on Machine*
306 *Learning and the Law*, 2016.
- 307 [14] C. Guo, G. Pleiss, Y. Sun, and K.Q. Weinberger. On calibration of modern neural networks. In
308 *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1321–1330,
309 2017.
- 310 [15] M. Hardt, E. Price, and N. Srebro. Equality of opportunity in supervised learning. In *Advances*
311 *in Neural Information Processing Systems (NeurIPS)*, 2016.
- 312 [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for im-
313 age recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
314 *Recognition (CVPR)*, 2016.
- 315 [17] S. Hooker. Moving beyond “algorithmic bias is a data problem”. *Patterns*, 2(4), 2021.
- 316 [18] N. Jain, A. Olmo, S. Sengupta, L. Manikonda, and S. Kambhampati. Imperfect imaGANation:
317 Implications of GANs exacerbating biases on facial data augmentation and Snapchat selfie
318 lenses. In *arXiv preprint arXiv:2001.09528*, 2020.

- 319 [19] S. Jia, T. Meng, J. Zhao, and K.-W. Chang. Mitigating gender bias amplification in distribution
320 by posterior regularization. In *Annual Meeting of the Association for Computational Linguistics*
321 (*ACL*), 2020.
- 322 [20] O. Keyes. The misgendering machines: Trans/HCI implications of automatic gender recognition.
323 In *Proceedings of the ACM on Human-Computer Interaction*, 2018.
- 324 [21] Alex Krizhevsky. Learning multiple layers of features from tiny images, 2009.
- 325 [22] M.J. Kusner, J.R. Loftus, C. Russell, and R. Silva. Counterfactual fairness. In *Advances in*
326 *Neural Information Processing Systems (NeurIPS)*, 2017.
- 327 [23] B.N. Larson. Gender as a variable in natural-language processing: Ethical considerations. In
328 *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, 2017.
- 329 [24] K. Leino, E. Black, M. Fredrikson, S. Sen, and A. Datta. Feature-wise bias amplification. In
330 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- 331 [25] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan. A survey on bias and
332 fairness in machine learning. *ACM Computing Surveys*, 54(6):1–35, 2022.
- 333 [26] J.A. Middleton, M.A. Scott, R. Diakow, and J.L. Hill. Bias amplification and bias unmasking.
334 *Political Analysis*, 3:307–323, 2016.
- 335 [27] M.P. Naeini, G.F. Cooper, and M. Hauskrecht. Obtaining well calibrated probabilities using
336 Bayesian binning. In *AAAI*, 2015.
- 337 [28] Y. Nesterov. A method for solving a convex programming problem with convergence rate
338 $o(1/k^2)$. *Soviet Mathematics Doklady*, 27:372–367, 1983.
- 339 [29] J. Pearl. On a class of bias-amplifying variables that endanger effect estimates. In *Uncertainty*
340 *in Artificial Intelligence*, 2010.
- 341 [30] J. Pearl. Invited commentary: Understanding bias amplification. *American Journal of Epidemi-*
342 *ology*, 174, 2011.
- 343 [31] M. K. Scheuerman, K. Wade, C. Lustig, and J. R. Brubaker. How we’ve taught algorithms to see
344 identity: Constructing race and gender in image databases for facial analysis. In *Proceedings of*
345 *the ACM on Human-Computer Interaction*, 2020.
- 346 [32] P. Stock and M. Cisse. Convnets and ImageNet beyond accuracy: Explanations, bias detection,
347 adversarial examples and model criticism. In *arXiv:1711.11443*, 2017.
- 348 [33] V. N. Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer-Verlag, 1982.
- 349 [34] Angelina Wang and Olga Russakovsky. Directional bias amplification. In *Proceedings of the*
350 *International Conference on Machine Learning (ICML)*, 2021.
- 351 [35] J.M. Wooldridge. Should instrumental variables be used as matching variables? *Research in*
352 *Economics*, 70:232–237, 2016.
- 353 [36] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for
354 benchmarking machine learning algorithms, 2017.
- 355 [37] D. Zhao, A. Wang, and O. Russakovsky. Understanding and evaluating racial biases in image
356 captioning. In *arXiv preprint arXiv:2106.08503*, 2021.
- 357 [38] J. Zhao, T. Wang, M. Yatskar, V. Ordonez, and K.-W. Chang. Men also like shopping: Reducing
358 gender bias amplification using corpus-level constraints. In *Proceedings of the Conference on*
359 *Empirical Methods in Natural Language Processing (EMNLP)*, 2017.



Figure 6: **Left:** Two examples of inversions performed on Fashion MNIST images. **Right:** Two examples of inversions performed on CIFAR-100 images. For each pair, the original image is on the left and the inverted image is on the right.

360 A Appendix

361 A.1 Experimental set-up

362 Here we discuss additional details in our experimental set-up.

363 A.1.1 Model training.

364 All our models are residual networks [16] that are trained to minimize the binary cross-entropy loss
 365 between the model prediction and the true (binary) class label. We follow the training procedures
 366 in [16] and train our models using mini-batch stochastic gradient descent (SGD) with a Nesterov
 367 momentum [28] of 0.9 for 500 epochs. The models are trained using weight decay (ℓ_2 -regularization)
 368 with a decay parameter of 10^{-4} . We warm up the training by setting the learning rate to 0.01 for one
 369 epoch as in [12]. Subsequently, the learning rate is set to 0.1 and decayed twice by a factor of 10 after
 370 250 and 375 epochs. We train on a single GPU using a batch size of 128.

371 During training, we adopt the data augmentation procedure of [16] by: randomly cropping training
 372 images, flipping the resulting image horizontally with probability $1/2$, and resizing the crops to size
 373 28×28 pixels (for Fashion MNIST) or 32×32 pixels (for CIFAR-10 and CIFAR-100). No data
 374 augmentation is used at test time. We normalize all images by subtracting a per-channel mean
 375 value and dividing by a per-channel standard deviation. When training models on CIFAR-10 and
 376 CIFAR-100, we follow [16] and pad the images with zeros.

377 A.1.2 Directional bias amplification.

378 We give a concise treatise of the measure here and refer the reader to [34] for further details.

379 Suppose we have a set of *classes*, \mathcal{T} , and a set of *groups*, \mathcal{A} . In our setup, $\mathcal{T} = \{-1, +1\}$ and
 380 $\mathcal{A} = \{\text{inverted, not inverted}\}$, where the binary labels $t \in \mathcal{T}$ were obtained by the random class
 381 assignment described above. The $\text{BiasAmp}_{\mathcal{A} \rightarrow \mathcal{T}}$ measure defines *bias* as a difference in the prevalence
 382 of a class label $t \in \mathcal{T}$ between groups $a \in \mathcal{A}$. For example, bias is present if inverted images are
 383 more likely to be positively labeled. Denote by $Pr(T_t = 1)$ the probability that an example in the
 384 dataset has class label t , and by $Pr(\hat{T}_t = 1)$ the probability that an example in the dataset is labeled
 385 as class t by the model. With these definitions, [34] defines *bias amplification* as the difference in
 386 bias between the labels in the dataset and the labels predicted by the model:

$$\text{BiasAmp}_{\mathcal{A} \rightarrow \mathcal{T}} = \frac{1}{|\mathcal{A}||\mathcal{T}|} \sum_{a \in \mathcal{A}, t \in \mathcal{T}} y_{at} \Delta_{at} - (1 - y_{at}) \Delta_{at}. \quad (2)$$

387 Δ_{at} measures the difference between the bias in the dataset and in the model predictions:

$$\Delta_{at} = Pr(\hat{T}_t = 1 | A_a = 1) - Pr(T_t = 1 | A_a = 1). \quad (3)$$

388 In the definition of $\text{BiasAmp}_{\mathcal{A} \rightarrow \mathcal{T}}$, y_{at} alters the sign of the difference Δ_{at} to correct for the fact that
 389 the bias can have two directions. Specifically, $y_{at} \in \{0, 1\}$ is a binary variable that indicates the
 390 direction of the bias:

$$y_{at} = [Pr(T_t = 1, A_a = 1) > Pr(T_t = 1)Pr(A_a = 1)], \quad (4)$$

391 where $[\dots]$ are Iverson brackets. In all our experiments, we compute $\text{BiasAmp}_{\mathcal{A} \rightarrow \mathcal{T}}$ by measuring
 392 both $Pr(\hat{T}_t = 1)$ and $Pr(T_t = 1)$ on the test set after training the model on the training set. The
 393 train and test datasets come from the same distribution.

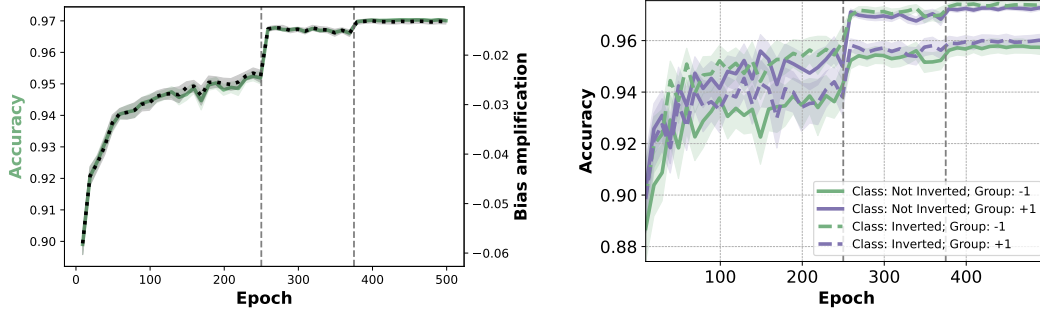


Figure 7: **Left:** Bias amplification and accuracy of ResNet-110 models during training on the CIFAR-100 dataset with a bias of $\epsilon = 0.3$ in which *the role of classes and groups is swapped* compared to the experiment in Figure 3: the class label indicates whether or not an image is inverted, and the group label is determined based on the visual content of the image. **Right:** Accuracy of the models’ per class-group combination. Shaded regions indicate the 95% confidence intervals across 50 models. Vertical dashed lines indicate epochs at which the learning rate of the mini-batch SGD optimizer is decreased.

394 A.2 Swapping group- and task-classes

395 In RQ4, we hypothesize that the early-stage bias amplification is due to group membership being
 396 easier to recognize than class membership in our setup. To test this hypothesis, we perform an
 397 experiment in which we swap the role of the group and the class: *i.e.*, the class label now represents
 398 whether or not the image is inverted and the group label depends on the object depicted in the CIFAR-
 399 100 image. We would expect the differences in accuracy between the majority / minority groups to
 400 disappear and bias amplification to actually be negative early on in training. As before, we measure
 401 bias amplification during training and plot the results in the left pane of Figure 7. The corresponding
 402 disaggregated accuracies are in the right pane of Figure 7. Indeed, we find that bias is dampened in
 403 the early stages of training as the model latches onto the easy-to-extract class membership signal
 404 first, but this largely disappears in the later stages of training as the model starts to utilize group
 405 membership signals for recognition as well.

406 A.3 Effect of training size on bias amplification.

407 It is well-established that the error of machine-learning models can be reduced by increasing the
 408 amount of training data (as it reduces the estimation error [4, 33]). This raises the obvious question
 409 if bias amplification varies with training set size as well. To answer this question, we perform
 410 experiments in which we train ResNet-110 models on stratified subsamples of the CIFAR-100
 411 training set. We vary the size of the subsamples to be a proportion, $p \in [0.1, 1.0]$, of the original
 412 training set. We increase the number of training epochs by a factor of $1/p$ so that each model
 413 performs the same number of parameter updates during training. We do not alter any of the other
 414 hyperparameters.

415 Figure 8 shows the results of our experiments. Whereas model accuracy increases monotonically
 416 with training set size, bias amplification varies in a more complex way. Beyond a certain training
 417 set size, bias amplification decreases with more training data. This is unsurprising: the additional
 418 training examples enable more accurate modeling of the data distribution, reducing bias amplification.
 419 However, bias amplification is also reduced when the training set becomes very small. We surmise
 420 this observation is due to overfitting: when trained on a small dataset, models tend to learn spurious
 421 correlations in that dataset rather than true statistical patterns such as the biases that exist in our
 422 training sets. The model cannot amplify bias if it is unable to capture that bias in the first place.

423 A.4 Overconfidence and bias amplification

424 Our observation that models with higher capacity amplify bias more is reminiscent of observations
 425 that higher-capacity models tend to be more miscalibrated [14]. If high-capacity models are not
 426 explicitly calibrated, they are often overconfident in the sense that the accuracy of predictions that

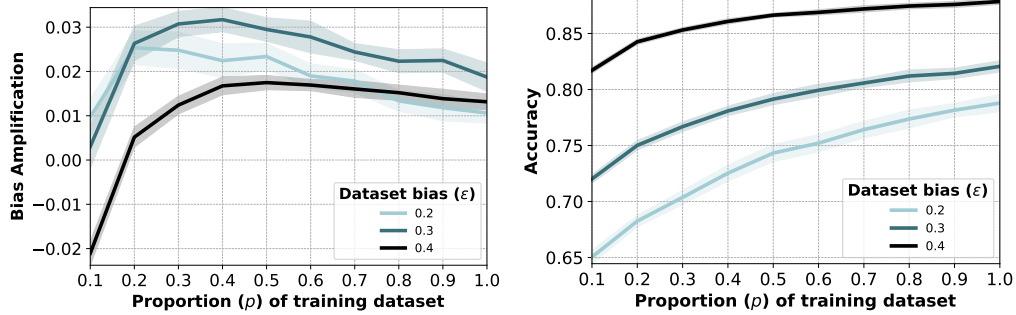


Figure 8: Bias amplification (**left**) and test accuracy (**right**) of ResNet-110 models on the CIFAR-100 dataset as a function of the proportion of the training set used for training the models. The number of epochs for each model is scaled depending on the amount of training data used. Shaded regions indicate the 95% confidence intervals across 20 models.

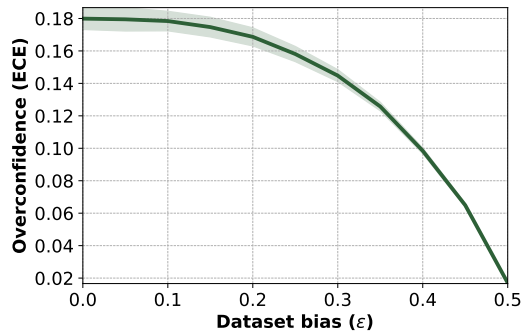


Figure 9: Expected calibration error (ECE) of ResNet-110 models on the CIFAR-100 dataset as a function of dataset bias, ϵ . Shaded regions indicate the 95% confidence intervals across 20 models.

427 they make with, say, 90% confidence is lower than 90%. We perform experiments to investigate if
 428 bias amplification is correlated to such model overconfidence.

429 To do so, we measure the overconfidence of our models in terms of the expected calibration error
 430 (ECE) [27]. The ECE measures the expected value of the (absolute) difference between the model
 431 accuracy and the model confidence:

$$ECE(\hat{P}) = \mathbb{E} \left[|Pr(\hat{Y} = y | \hat{C} = c) - c| \right], \quad (5)$$

432 where \hat{Y} and \hat{C} are random variables indicating the class label of an example and the model-
 433 prediction confidence for that same example, respectively, and the expectation is over all possible
 434 confidence values $c \in [0, 1]$. Because we only have access to a finite number of samples of the
 435 distribution $p(\hat{C})$, we approximate the expected value by binning $p(\hat{C})$ into 15 values and averaging
 436 those values, weighted by the number of examples per bin. A higher ECE value indicates a larger
 437 discrepancy between the prediction confidence values and the corresponding accuracies, *i.e.*, a higher
 438 degree of model overconfidence.

439 Figure 9 shows ECE as a function of the bias in the dataset, ϵ , for ResNet-110 models on CIFAR-100.
 440 We observe that model overconfidence decreases with bias in our experiment, because the task
 441 becomes easier as bias increases: if a task is very easy, a model is generally less overconfident as it
 442 correctly predicts nearly every example.

443 Next, we study the relation between overconfidence and bias amplification by varying the capacity of
 444 the model. Figure 10 shows this relation for three levels of dataset bias, ϵ , and for three model-capacity
 445 measures: depth, width, and weight decay. Darker points in the figure correspond to higher-capacity
 446 models. The results show that bias amplification initially decreases as model overconfidence increases

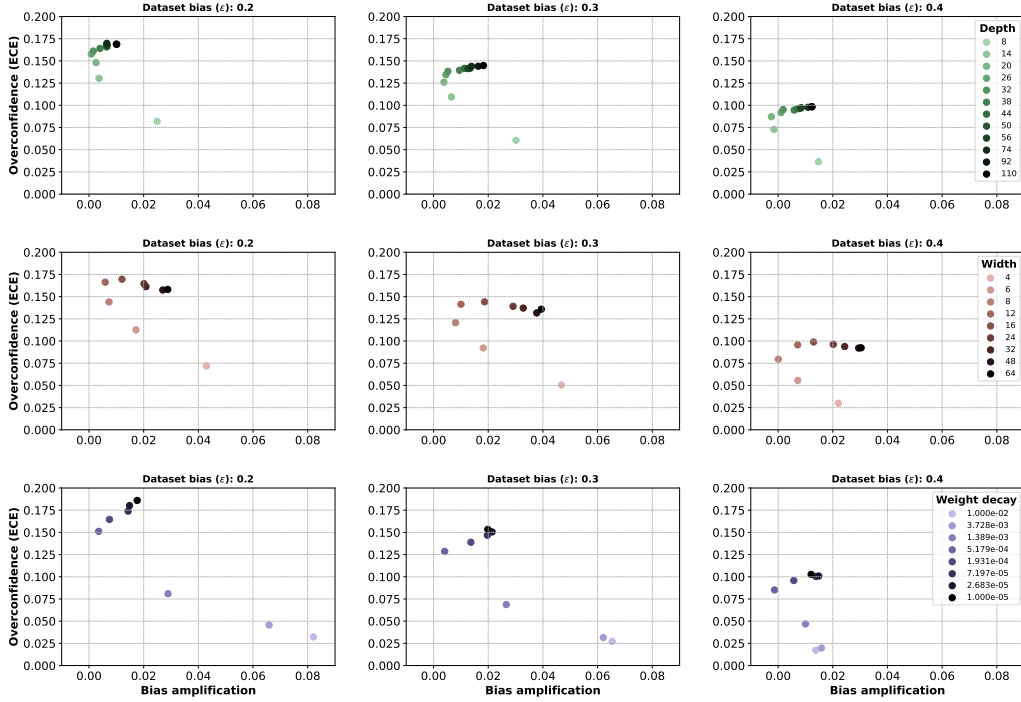


Figure 10: Bias amplification and expected calibration error (ECE) of ResNet models of varying depth (**first row**), width (**second row**), and weight decay (**third row**) on the CIFAR-100 dataset, for three values of the dataset bias, ϵ . Results are averaged over 20 runs.

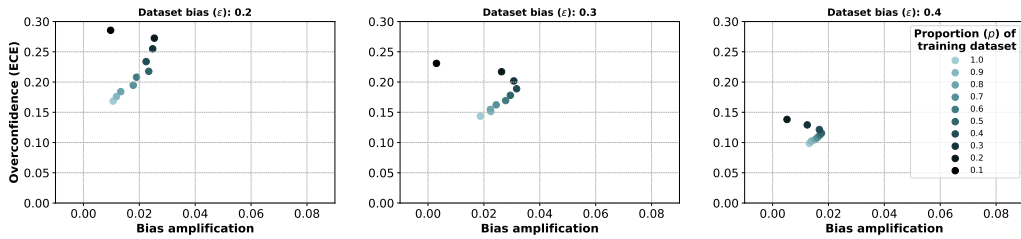


Figure 11: Bias amplification and expected calibration error (ECE) of ResNet models of varying training dataset size on the CIFAR-100 dataset, for three values of the dataset bias, ϵ . Results are averaged over 20 runs.

447 (for low-capacity models), but that bias amplification and overconfidence both increase for higher-
 448 capacity models.

449 Finally, Figure 11 studies the relationship between bias amplification and overconfidence as the
 450 size of the training set changes. Darker points in the figure correspond to smaller training sets. As
 451 expected, reducing the number of training examples increases the model's overconfidence (ECE).
 452 Bias amplification, however, initially increases as the training set size decreases but decreases again
 453 when the training set becomes very small.