

ON THE UNREASONABLE EFFECTIVENESS OF LAST-LAYER RETRAINING

Anonymous authors

Paper under double-blind review

ABSTRACT

Last-layer retraining (LLR) methods — wherein the last layer of a neural network is reinitialized and retrained on a held-out set following ERM training — have recently garnered interest as an efficient approach to rectify dependence on spurious correlations and improve performance on minority groups. Surprisingly, LLR has recently been found to improve worst-group accuracy even when the held-out set is an imbalanced subset of the training set. We initially hypothesize that this “unreasonable effectiveness” of LLR is explained by its ability to mitigate neural collapse through the held-out set, resulting in the implicit bias of gradient descent benefiting robustness. Our empirical investigation *does not* support this hypothesis. Instead, we present strong evidence for an alternative hypothesis: that the success of LLR is primarily due to better group balance in the held-out set. We conclude by showing how the recent algorithms CB-LLR and AFR perform implicit group-balancing to elicit a robustness improvement.

1 INTRODUCTION

The standard neural network training procedure of empirical risk minimization (ERM) (Vapnik, 1998), which minimizes the average classification loss, is well-known to overfit to spurious correlations in the training set (Geirhos et al., 2020). The conjunction of target labels and spurious features form *minority groups*, upon which ERM performance may be no better than random guessing (Shah et al., 2020). Due to the relevance of spurious correlations in high-consequence applications, *e.g.*, medicine (Zech et al., 2018) and criminal justice (Chouldechova, 2016), significant work has focused on algorithms which maximize worst-group accuracy (WGA) — also called group robustness (Sagawa et al., 2020).

A promising family of group robustness methods are based on last-layer retraining (LLR), wherein the last layer is reinitialized and retrained on a held-out set following ERM training. The original LLR method, called deep feature reweighting (DFR), requisites the held-out set to comprise equal data from each group (Kirichenko et al., 2023). This limits its practical application, as the groups are often unknown ahead of time or difficult to annotate.¹ While DFR still performs the best, group information was recently found to be unnecessary to observe WGA improvement from LLR (Qiu et al., 2023; LaBonte et al., 2023). This surprising observation has led LLR to be termed a “free lunch” for group robustness (LaBonte et al., 2023).

Contributions. In this paper, we take a “scientific method” approach to investigate why LLR on an imbalanced held-out subset of the training set can perform so well. Our contributions include:

- We propose an initial hypothesis (visualized in Figure 1) that neural collapse on the training set results in a biased ERM classifier since the class means are dominated by majority group data. We initially hypothesize that during LLR, the features are not collapsed on the held-out set, and so the implicit bias of gradient descent elicits a maximum-margin classifier that could enjoy better robustness guarantees.
- We put forth evidence which *does not* support our initial hypothesis. First, neural collapse either does not occur or occurs after the standard number of epochs for our benchmark datasets. Moreover, convergence of the LLR classifier to the maximum-margin solution is extremely slow, and the *average* margin over all training data is much more correlated with group accuracy than the (standard) *minimum* margin over training data.

¹However, DFR compares favorably in this respect to methods which require group annotations for the entire training set, such as group distributionally robust optimization (DRO) (Sagawa et al., 2020).

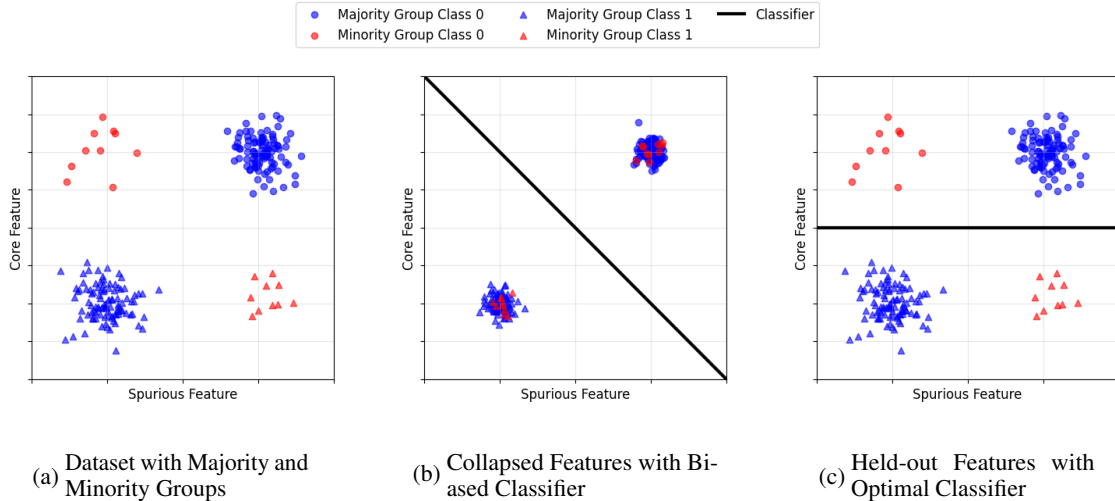


Figure 1: **Visualization of our initial hypothesis.** In Figure 1(a), we visualize the feature space of a training set with two classes and two features, one core and one spurious. In Figure 1(b), we visualize our initial hypothesis that neural collapse on the training set causes the features to collapse to their class means — dominated by majority group data — resulting in a biased ERM classifier. In Figure 1(c), we visualize our initial hypothesis that during LLR, the features are not collapsed on the held-out set, and so the implicit bias of gradient descent elicits a maximum-margin classifier which is invariant to the spurious feature. Importantly, our empirical investigation *does not* support our initial hypothesis. Instead, we find that the success of LLR is primarily explained by better group balance in the held-out set.

- We present strong evidence for an alternative hypothesis: that the success of LLR is primarily due to *better group balance* in the held-out set, often achieved implicitly through class balancing. Our experiments indicate that LLR does not improve over ERM when the held-out set has the same group balance as the training set, with Pearson correlation coefficient $r > 0.99$ on CelebA and CivilComments. On the other hand, LLR with a better group balance enjoys drastically improved WGA, and vice versa.
- Ultimately, we reevaluate the “free lunch” interpretation of LLR by showing that class balanced LLR (CB-LLR) (LaBonte et al., 2023) and automatic feature reweighting (AFR) (Qiu et al., 2023) owe their improved WGA primarily to implicit group-balancing on the held-out set. On a positive note, we show that LLR can recover the WGA of an optimally class-balanced model even when ERM was not optimally class-balanced. More broadly, LLR remains an effective method to achieve group robustness using group annotations only on the held-out set (or proxies thereof, such as in SELF (LaBonte et al., 2023) and AFR (Qiu et al., 2023)).

Related work. Reliance on spurious correlations is a widely-studied phenomenon in machine learning known to exacerbate bias and hinder generalization, *e.g.*, Beery et al. (2018); Geirhos et al. (2020); Sagawa et al. (2020). Last-layer retraining, proposed by Kirichenko et al. (2023), has recently garnered interest as an efficient approach to improve group robustness in a variety of settings (Qiu et al., 2023; LaBonte et al., 2023; Stromberg et al., 2024; Park et al., 2025). Our work aims to further a fundamental understanding of LLR methods, following broadly in the spirit of Izmailov et al. (2022); Ye et al. (2023); Chen et al. (2023); Welfert et al. (2024); Qiu et al. (2024). Our initial hypothesis primarily concerns neural collapse, a phenomenon introduced by Pappas et al. (2020); Han et al. (2022) which has been applied to the study of fairness and spurious correlations in the context of ERM by Lu et al. (2024); Wang et al. (2024); Chen et al. (2024); Xu et al. (2025), but never before to LLR. We also consider the implicit bias of gradient descent in logistic regression towards the maximum ℓ_2 -margin classifier, studied in the separable case by Soudry et al. (2018) and the non-separable case by Ji & Telgarsky (2019) (though a rich literature exists on this topic; see Vardi (2023) for a recent survey). Finally, in Section 3 we evaluate the impact of group balance on LLR performance via a case study on class-balancing, a popular lightweight method to improve group robustness studied by Idrissi et al. (2022); Chaudhuri et al. (2023); Schwartz-Ziv et al. (2023); LaBonte et al. (2024).

Setting. We consider the setting of classification tasks with input domain \mathcal{X} and target classes \mathcal{Y} . We assume the existence of a set of *spurious features* \mathcal{S} such that each data point $x \in \mathcal{X}$ is associated with a single spurious feature

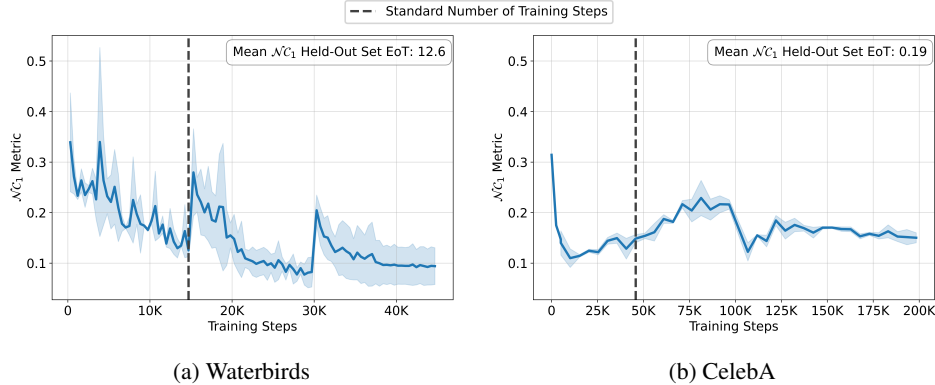


Figure 2: **Collapsing of class feature variability occurs after standard ERM training, if at all.** We plot the empirical metric of neural collapse $\mathcal{N}C_1 := \frac{1}{|\mathcal{Y}|} \text{tr}(\Sigma_A \Sigma_R^\dagger)$ throughout training a ResNet-18 on Waterbirds and CelebA. Each plot displays the mean and standard deviation for $\mathcal{N}C_1$ computed across 3 experimental seeds. We also display the mean $\mathcal{N}C_1$ metric computed on the features of the held-out set at the end of training (EoT). We believe the cyclic spikes in the Waterbirds plot are due to the cosine learning rate scheduler used when training the model.

$s(x) \in \mathcal{S}$. The dataset is partitioned into *groups* \mathcal{G} by the Cartesian product of classes and spurious features: $\mathcal{G} := \mathcal{Y} \times \mathcal{S}$. Let $\Omega_g \subseteq \{1, \dots, m\}$ denote the indices of training points belonging to group $g \in \mathcal{G}$. Similarly, let $\Omega_y \subseteq \{1, \dots, m\}$ denote the indices of training points belonging to class $y \in \mathcal{Y}$. We define the *majority group(s)* to be the groups which maximize $|\Omega_g|$. All other groups are considered to be *minority groups*. Importantly, we also define the *worst group(s)* to be the group(s) with lowest test accuracy. Our goal is to find a model which, despite group imbalance in the training set, performs uniformly well across groups. To evaluate this, we use *worst-group accuracy* (WGA), or the minimum accuracy achieved among all groups (Sagawa et al., 2020).

Datasets and models. We employ four benchmark classification datasets in our experiments: Waterbirds (Welinder et al., 2010; Wah et al., 2011; Sagawa et al., 2020), CelebA (Liu et al., 2015; Sagawa et al., 2020), CivilComments (Borkan et al., 2019; Koh et al., 2021), and MultiNLI (Williams et al., 2018; Sagawa et al., 2020). We discuss each dataset at length in Appendix B.1. We use ResNet-18 and ResNet-50 (He et al., 2016) models pretrained on ImageNet-1K (Russakovsky et al., 2015) for the Waterbirds and CelebA datasets, and a BERT-Base (Devlin et al., 2019) model pretrained on Book Corpus (Zhu et al., 2015) and English Wikipedia for CivilComments and MultiNLI. Our implementation of LLR uses SGD with no learning rate schedule or regularization. Following standard practice, we use 20% of the training set for the held-out set on Waterbirds and half the validation set as the held-out set in all other datasets. Further training details are located in Appendix B.2.

2 NEURAL COLLAPSE AND IMPLICIT BIAS DO NOT EXPLAIN LLR

The interpretation of LLR methods as logistic regression on top of convergent ERM features lends itself to a possible understanding combining a feature learning phenomenon (*neural collapse*) and a logistic regression phenomenon (*implicit bias*). Neural collapse was first identified by Pappas et al. (2020); Han et al. (2022) and implies that the penultimate layer features collapse to their class means over the course of ERM training. On the other hand, implicit bias results state that linear classifiers trained via gradient descent with the unregularized logistic loss converge in direction to the maximum-margin SVM solution (Soudry et al., 2018; Ji & Telgarsky, 2019).

Combining these two ideas, we developed the following initial hypothesis (visualized in Figure 1). We hypothesized that the model undergoes neural collapse during ERM training, causing the features to collapse to their class means. These class means, however, are dominated by the largest groups in each class — resulting in ERM learning a biased classifier. We hypothesized that during LLR, the features are not collapsed as the held-out set was not seen during ERM; hence, the implicit bias of gradient descent elicits a robust maximum-margin classifier.

In this section, we present evidence which *does not* support our initial hypothesis. Specifically, we find that neural collapse either does not occur or occurs after the standard number of epochs of ERM training. Moreover, convergence of

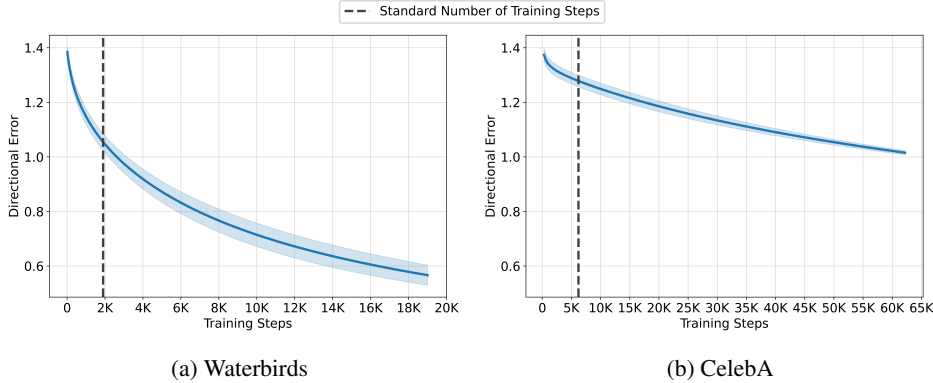


Figure 3: **Convergence of LLR to the maximum-margin SVM solution is extremely slow.** We plot the mean and standard deviation over 3 experimental seeds of the directional error \widehat{Err} between the last layer weights of a ResNet-18 last layer and an SVM (both trained on the features of the held-out set). Here, $\widehat{Err} := \left\| \frac{\theta_{\text{NN}}}{\|\theta_{\text{NN}}\|_2} - \frac{\theta_{\text{SVM}}}{\|\theta_{\text{SVM}}\|_2} \right\|_2$, where θ_{NN} denotes the last layer weights and θ_{SVM} denotes the weights of an SVM trained on the held-out set features.

the LLR classifier to the maximum-margin solution is extremely slow, and the *average* margin is much more correlated with group accuracy than the (standard) *minimum* margin.

2.1 NEURAL COLLAPSE MAY NOT OCCUR DURING STANDARD ERM TRAINING

We will measure the collapse of variability among feature representations via a metric called \mathcal{NC}_1 . For each class $y \in \mathcal{Y}$ and each training example $i \in \Omega_y$ in that class, we denote the penultimate layer feature vector of i as $f_{y,i}$. Neural collapse posits that features collapse to their respective class mean $\mu_y := \frac{1}{|\Omega_y|} \sum_{i \in \Omega_y} f_{y,i}$. We compute an empirical metric of this collapse using the intra-class covariance matrix $\Sigma_A := \frac{1}{m} \sum_{y \in \mathcal{Y}} \sum_{i \in \Omega_y} (f_{y,i} - \mu_y)(f_{y,i} - \mu_y)^\top$, the global feature mean $\mu_G := \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \mu_y$, and the inter-class covariance matrix $\Sigma_R := \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} (\mu_y - \mu_G)(\mu_y - \mu_G)^\top$. We then define $\mathcal{NC}_1 := \frac{1}{|\mathcal{Y}|} \text{tr}(\Sigma_A \Sigma_R^\dagger)$, where Σ_R^\dagger denotes the pseudo-inverse of Σ_R . Neural collapse is formalized by $\mathcal{NC}_1 \rightarrow 0$ (Han et al., 2022; Kothapalli, 2023).²

If neural collapse plays a significant role in ERM learning a biased classifier prior to LLR, then a collapse in feature variability should be observed during a standard number of ERM training steps. We display \mathcal{NC}_1 for a ResNet-18 on Waterbirds and CelebA in Figure 2. We train on both datasets for much longer than normal and denote the standard number of training steps in the spurious correlations literature (Sagawa et al., 2020; Kirichenko et al., 2023).

In Figure 2, we see evidence of feature collapse on Waterbirds. \mathcal{NC}_1 is two orders of magnitude larger for the held-out set than for the training set by the end of training. However, complete collapse is not achieved by the standard number of epochs. The CelebA features do not appear to collapse within any reasonable number of training steps as \mathcal{NC}_1 remains comparable between the training and held-out sets even after more than $3 \times$ the standard number of training steps. Also, \mathcal{NC}_1 measured on the CelebA training set experiences no significant decrease after the first epoch. Thus, it is unlikely that neural collapse significantly affects the classifier learned by ERM prior to LLR in real-world datasets.

2.2 THE MINIMUM MARGIN OF LLR IS NOT PREDICTIVE OF ROBUST GENERALIZATION

Let $f_\theta : \mathcal{X} \rightarrow \{-1, 1\}$ be a binary linear classifier defined by $f_\theta(x) = x^\top \theta + b$. Let $\mathcal{D} := \{(x_i, y_i)\}_{i=1}^n \subseteq \mathcal{X} \times \{-1, 1\}$ be the training set and $\mathcal{C} := \{x_i : y_i f_\theta(x_i) > 0\}$ be the set of points correctly classified by f_θ . Then, we define the minimum ℓ_2 -margin of f_θ to be $\min_{x_i \in \mathcal{C}} y_i f_\theta(x_i)$ and the average ℓ_2 -margin of f_θ to be $\frac{1}{n} \sum_{i=1}^n y_i f_\theta(x_i)$. The hard margin SVM is a binary linear classifier formulated through the following optimization problem: $\min_{\theta, b} \frac{1}{2} \|\theta\|_2^2$ such that $y_i(x_i^\top \theta + b) \geq 1$ for all $(x_i, y_i) \in \mathcal{D}$. The classifier θ_{SVM} learned by a hard margin SVM is the separating hyperplane which maximizes the minimum ℓ_2 -margin. Soudry et al. (2018) show that linear classifiers trained via

²A slightly less precise formulation $\Sigma_A \rightarrow 0$ was studied by Papyan et al. (2020).

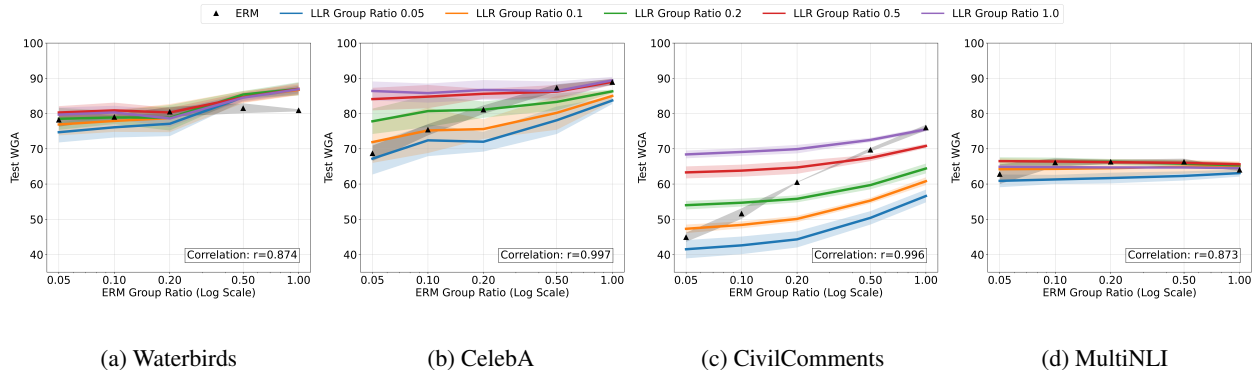


Figure 4: **LLR performance is determined by held-out set group balance.** We compare the test WGA of ERM and LLR models while controlling the group balance of the training and held-out sets. We train a ResNet-50 for the vision datasets and a BERT-Base for the language datasets, and we plot the mean and standard deviation over 3 experimental seeds. A group ratio of 1.00 corresponds to the majority and minority groups having equal size, whereas a group ratio of 0.05 implies that the minority groups are 5% the size of the majority groups. We compute the Pearson correlation coefficient between the test WGA of ERM and the test WGA of LLR for each dataset while keeping group ratios fixed in the comparison. We find that the correlation coefficient is close to 1 for all datasets. In particular, LLR tends to improve over ERM *only if the held-out set has better group balance.*

GD with the unregularized logistic loss converge in direction to θ_{SVM} . In particular, $\|\frac{\theta_{\text{GD}}}{\|\theta_{\text{GD}}\|_2} - \frac{\theta_{\text{SVM}}}{\|\theta_{\text{SVM}}\|_2}\|_2 \rightarrow 0$ at a worst-case rate of $O(\frac{\log \log t}{\log t})$ where t is the number of gradient steps. This result can also be extended to the nonseparable case, though, as one might expect, the implicit bias is more complex (Ji & Telgarsky, 2019).

If LLR maximizes the minimum training margin and this results in a debiased classifier, then we should see a high correlation between the size of group minimum margin and group test accuracy (*i.e.*, minority group points lie near the decision boundary). However, we find that the minimum training margin for a group is not correlated with the test accuracy for that group. Interestingly, we find that *average* group training margin actually is correlated with group test accuracy. We display the Pearson correlation coefficients in Table 1 (deferred to Appendix A).

Moreover, convergence towards θ_{SVM} is slow in practice. We train a ResNet-18 last layer on the held-out set via SGD with learning rate 0.001 and unregularized logistic loss, and we plot its directional error with θ_{SVM} in Figure 3. In line with the slow $O(\frac{\log \log t}{\log t})$ rate, we find that even after $10\times$ the number of standard training steps, the LLR classifier has yet to converge to θ_{SVM} . As convergence to θ_{SVM} is slow and minimum ℓ_2 -margin is uncorrelated with group test accuracy, we conclude that SVM convergence and minimum margin bounds are unlikely to explain LLR.

3 GROUP BALANCE CORRELATES STRONGLY WITH LLR PERFORMANCE

The negative results of Section 2 prompted us to reconsider our initial hypothesis. We now present strong evidence for an alternative hypothesis: that the success of LLR is primarily due to *better group balance* in the held-out set than the training set. We also show how recent LLR algorithms owe their success to implicit group-balancing on the held-out set.

3.1 LLR PERFORMANCE IS DETERMINED BY HELD-OUT SET GROUP BALANCE

To isolate the effect of group balance on LLR, we perform an ablation study, shown in Figure 4. We directly control the group ratios — defined as the ratio between the number of minority group points and the number of majority group points — by removing data until the desired group ratio is achieved for each class. We vary the group ratios for both the ERM training set and the LLR held-out set between 0.05 and 1.00, keeping the total data in each stage constant. For a fair comparison, the ERMs in Figure 4 are trained with the held-out set added in. For example, on Waterbirds we compare LLR with 20% of data following an ERM with 80% of data (colored line) to an ERM with 100% of data (gray line). Note that LLR with a group ratio of 1.00 corresponds to DFR.

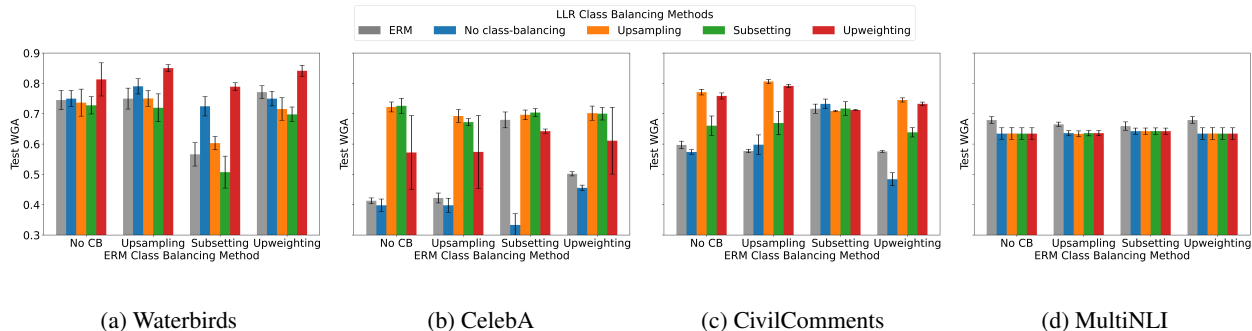


Figure 5: **LLR does not significantly improve over optimally class-balanced ERM.** We plot the test WGA of ERM (trained on 100% of available data) and LLR (trained using a training/held-out set split) using no class-balancing, upsampling, upweighting, and subsetting. We train a ResNet-50 for the vision datasets and a BERT-Base for the language datasets and plot the mean and standard deviation over 3 experimental seeds. We find ERM matches the performance of LLR when ERM is trained with optimal class-balancing. Note that MultiNLI is class balanced *a priori*.

We find that LLR performance is almost completely determined by the change in group balance between the training and held-out datasets. LLR only shows significant gains over ERM when it is trained on a held-out set with a *higher group ratio* than the ERM training set, and vice versa. In Figure 4, we detail the Pearson correlation coefficients between ERM and LLR test WGA across all group ratios and observe a strong trend, with $r > 0.99$ on CelebA and CivilComments. Moreover, the test WGA of LLR trained with a fixed group ratio is nearly identical to the test WGA of ERM trained with that same group ratio (observed most impressively on CelebA).

These results indicate that the two-stage training procedure of LLR methods is not fundamentally oriented towards learning a debiased classifier as we initially hypothesized. Rather, some form of group-balancing is necessary to capture any of the robustness benefits achieved by DFR. Moreover, LLR methods are limited by how well they improve the group ratio, and thus DFR likely upper bounds the WGA of any LLR method which does not utilize group annotations.

3.2 RECENT LLR METHODS SUCCEED VIA IMPLICIT GROUP-BALANCING

In this section, we show how the recent LLR methods of class-balanced LLR (CB-LLR) (LaBonte et al., 2023) and automatic feature reweighting (AFR) (Qiu et al., 2023) perform *implicit group-balancing* to elicit a robustness improvement without group annotations.³ Oftentimes, this is achieved by *class-balancing*.

We utilize three standard techniques for class-balancing: *subsetting*, wherein the size of each class is reduced to match the size of the smallest class; *upsampling*, wherein sampling probabilities for each class are adjusted so that the mini-batches are balanced in expectation; and *upweighting*, wherein minority class samples are assigned larger weights in the loss function (Idrissi et al., 2022; LaBonte et al., 2024). Since groups are defined with respect to classes, each of these class-balancing methods partially group-balance the dataset without requiring group annotations.

In Figure 5, we compare the test WGA of class-balanced ERM and class-balanced LLR, and we find that the improvement of LLR over ERM is dependent on the class-balancing methods used. We find optimally class-balanced ERM achieves nearly the same test WGA as optimally class-balanced LLR on all datasets! The CB-LLR results of LaBonte et al. (2023) use upsampling for ERM (as well as LLR), which on CelebA and CivilComments is suboptimal and explains why CB-LLR seems to work so well. On a positive note, this means that LLR can recover the WGA of an optimally class-balanced model even when ERM experiences catastrophic collapse.

Compared to standard LLR, which utilizes the cross-entropy loss only, AFR incorporates a weighted loss function which prioritizes points upon which the ERM model performs poorly. Specifically, AFR introduces a weight for each held-out example i proportional to $\exp(-\gamma \hat{p}_i)$, where \hat{p}_i is the probability for the correct class y_i and $\gamma \geq 0$ is a tunable temperature parameter.⁴ Therefore, the AFR held-out set effectively has better group balance than the training set! Notably, the ablations of Qiu et al. (2023) show that setting $\gamma = 0$ reduces to CB-LLR with upweighting, thus explaining its improvement over an ERM with no class-balancing (see Figure 5).

³LaBonte et al. (2023) also propose SELF: unlike AFR and CB-LLR, it *explicitly* constructs a more group-balanced held-out set.

⁴Loss-based adjustments are common among group robustness methods not using group annotations, e.g., Liu et al. (2021).

312 REFERENCES

- 313 Sara Beery, Grant van Horn, and Pietro Perona. Recognition in terra incognita. In *European Conference on Computer*
314 *Vision (ECCV)*, 2018.
- 316 Daniel Borkan, Lucas Dixon, Jeffrey Sorenson, Nithium Thain, and Lucy Vasserman. Nuanced metrics for measuring
317 unintended bias with real data for text classification. In *World Wide Web (WWW)*, 2019.
- 318 Kamalika Chaudhuri, Kartik Ahuja, Martin Arjovsky, and David Lopez-Paz. Why does throwing away data improve
319 worst-group error? In *International Conference on Machine Learning (ICML)*, 2023.
- 321 Yongqiang Chen, Wei Huang, Kaiwen Zhou, Yatao Bian, Bo Han, and James Cheng. Understanding and improving
322 feature learning for out-of-distribution generalization. In *Conference on Neural Information Processing Systems*
323 *(NeurIPS)*, 2023.
- 324 Zhikang Chen, Min Zhang, Sen Cui, Haoxuan Li, Gang Niu, Mingming Gong, Changshui Zhang, and Kun Zhang.
325 Neural collapse inspired feature alignment for out-of-distribution generalization. In *Conference on Neural Information*
326 *Processing Systems (NeurIPS)*, 2024.
- 328 Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments.
329 In *Conference on Fairness, Accountability, and Transparency in Machine Learning (FATML)*, 2016.
- 330 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional
331 transformers for language understanding. In *North American Association for Computational Linguistics (NAACL)*,
332 2019.
- 334 William Falcon and the PyTorch Lightning maintainers and contributors. PyTorch Lightning. *GitHub*, 2019.
- 335 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and
336 Felix A. Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2:665–673, 2020.
- 338 X. Y. Han, Vardan Papyan, and David Donoho. Neural collapse under MSE loss: Proximity to and dynamics on the
339 central path. In *International Conference on Learning Representations (ICLR)*, 2022.
- 340 Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric
341 Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van
342 Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-
343 Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E.
344 Oliphant. Array programming with NumPy. *Nature*, 585(1):357–362, 2020.
- 346 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference*
347 *on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 348 John D. Hunter. Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007.
- 350 Badr Youbi Idrissi, Martín Arjovsky, Mohammad Pezeshki, and David Lopez-Paz. Simple data balancing achieves
351 competitive worst-group-accuracy. In *Conference on Causal Learning and Reasoning (CLeAR)*, 2022.
- 352 Pavel Izmailov, Polina Kirichenko, Nate Gruver, and Andrew Gordon Wilson. On feature learning in the presence of
353 spurious correlations. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- 355 Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning*
356 *Theory (COLT)*, 2019.
- 357 Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient for robustness to
358 spurious correlations. In *International Conference on Learning Representations (ICLR)*, 2023.
- 360 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu,
361 Michihiro Yasunaga, Richard Lanus Phillips, Sara Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey
362 Levine, Chelsea Finn, and Percy Liang. WILDS: A benchmark of in-the-wild distribution shifts. In *International*
363 *Conference on Machine Learning (ICML)*, 2021.

364 Vignesh Kothapalli. Neural collapse: A review on modelling principles and generalization. *Transactions on Machine*
365 *Learning Research (TMLR)*, 2023.

366 Tyler LaBonte. Milkshake: Quick and extendable experimentation with classification models. <http://github.com/tmlabonte/milkshake>, 2023.

367 Tyler LaBonte, Vidya Muthukumar, and Abhishek Kumar. Towards last-layer retraining for group robustness with
368 fewer annotations. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2023.

369 Tyler LaBonte, John C. Hill, Xinchun Zhang, Vidya Muthukumar, and Abhishek Kumar. The group robustness is in the
370 details: Revisiting finetuning under spurious correlations. In *Conference on Neural Information Processing Systems*
371 *(NeurIPS)*, 2024.

372 Evan Zheran Liu, Behzad Haghgoo, Annie S. Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang,
373 and Chelsea Finn. Just train twice: Improving group robustness without training group information. In *International*
374 *Conference on Machine Learning (ICML)*, 2021.

375 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *International*
376 *Conference on Computer Vision (ICCV)*, 2015.

377 Shenyu Lu, Junyi Chai, and Xiaoqian Wang. Neural collapse inspired debiased representation learning for min-max
378 fairness. In *Conference on Knowledge Discovery and Data Mining (KDD)*, 2024.

381 Vardan Papyan, X.Y. Han, and David Donoho. Prevalence of neural collapse during the terminal phase of deep learning
382 training. *Proceedings of the National Academy of Sciences (PNAS)*, 117(40):24652–24663, 2020.

383 Juhyeon Park, Seokhyeon Jeong, and Taesup Moon. TLDR: text based last-layer retraining for debiasing image
384 classifiers. In *Winter Conference on Applications of Computer Vision (WACV)*, 2025.

385 Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban
386 Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *Conference on Neural*
387 *Information Processing Systems (NeurIPS) Workshop on Automatic Differentiation*, 2017.

388 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming
389 Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zach DeVito, Martin
390 Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch:
391 An imperative style, high-performance deep learning library. In *Conference on Neural Information Processing*
392 *Systems (NeurIPS)*, 2019.

393 GuanWen Qiu, Da Kuang, and Surbhi Goel. Complexity matters: Feature learning in the presence of spurious
394 correlations. In *International Conference on Machine Learning (ICML)*, 2024.

395 Shikai Qiu, Andres Potapczynski, Pavel Izmailov, and Andrew Gordon Wilson. Simple and fast group robustness by
396 automatic feature reweighting. In *International Conference on Machine Learning (ICML)*, 2023.

400 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy,
401 Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet large scale visual recognition
402 challenge. *International Journal of Computer Vision (IJCV)*, 115(1):211–252, 2015.

403 Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks for
404 group shifts: On the importance of regularization for worst-case generalization. In *International Conference on*
405 *Learning Representations (ICLR)*, 2020.

406 Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of simplicity
407 bias in neural networks. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.

408 Ravid Shwartz-Ziv, Micah Goldblum, Yucen Lily Li, C Bayan Bruss, and Andrew Gordon Wilson. Simplifying neural
409 network training under class imbalance. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2023.

410 Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient
411 descent on separable data. *Journal of Machine Learning Research (JMLR)*, 19:1–57, 2018.

416 Nathan Stromberg, Rohan Ayyagari, Monica Welfert, Sanmi Koyejo, and Lalitha Sankar. Robustness to subpopulation
417 shift with domain label noise via regularized annotation of domains. *Transactions on Machine Learning Research*
418 (*TMLR*), 2024.

419 Saeid Asgari Taghanaki, Aliasghar Khani, Fereshte Khani, Ali Gholami, Linh Trana, Ali Mahdavi-Amiri, and Ghassan
420 Hamarneh. MaskTune: mitigating spurious correlations by forcing to explore. In *Conference on Neural Information*
421 *Processing Systems (NeurIPS)*, 2022.

422 TorchVision maintainers and contributors. TorchVision: PyTorch’s computer vision library. *GitHub*, 2016.

423 Vladimir Vapnik. *Statistical Learning Theory*. Wiley, 1998.

424 Gal Vardi. On the implicit bias in deep-learning algorithms. *Communications of the ACM*, 66(6):86–93, 2023.

425 Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The Caltech-UCSD birds-200-2011
426 dataset. Technical report, California Institute of Technology, 2011.

427 Yining Wang, Junjie Sun, Chenyue Wang, Mi Zhang, and Min Yang. Navigate beyond shortcuts: Debaised learning
428 through the lens of neural collapse. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.

429 Monica Welfert, Nathan Stromberg, and Lalitha Sankar. Theoretical guarantees of data augmented last layer retraining
430 methods. In *International Symposium on Information Theory (ISIT)*, 2024.

431 Peter Welinder, Steve Branson, Takeshi Mita, Catherine Wah, Florian Schroff, Serge Belongie, and Pietro Perona.
432 Caltech-UCSD birds 200. Technical report, California Institute of Technology, 2010.

433 Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence understanding
434 through inference. In *North American Association for Computational Linguistics (NAACL)*, 2018.

435 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac,
436 Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine
437 Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander
438 Rush. Transformers: State-of-the-art natural language processing. In *Conference on Empirical Methods in Natural*
439 *Language Processing (EMNLP) System Demonstrations*, 2020.

440 Jingxuan Xu, Wuyang Chen, Linyi Li, Yao Zhao, and Yunchao Wei. Collapsed language models promote fairness. In
441 *International Conference on Learning Representations (ICLR)*, 2025.

442 Haotian Ye, James Zou, and Linjun Zhang. Freeze then train: Towards provable representation learning under spurious
443 correlations and feature noise. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2023.

444 John R. Zech, Marcus A. Badgeley, Manway Liu, Anthony B. Costa, Joseph J. Titano, and Eric Karl Oermann. Variable
445 generalization performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional
446 study. *PLoS Medicine*, 15, 2018.

447 Michael Zhang, Nimit S. Sohoni, Hongyang R. Zhang, Chelsea Finn, and Christopher Ré. Correct-n-contrast: A
448 contrastive approach for improving robustness to spurious correlations. In *International Conference on Machine*
449 *Learning (ICML)*, 2022.

450 Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler.
451 Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In
452 *International Conference on Computer Vision (ICCV)*, 2015.

A ADDITIONAL EXPERIMENTS

In this section, we present additional experiments deferred from the main body of the work. First, in Table 1, we include the Pearson correlation coefficient between the test accuracy on each group and its minimum/average margin, calculated as defined in Section 2.2. We find that the minimum training margin on a group is surprisingly uncorrelated with the test accuracy on that group; instead, the average training margin on a group is strongly correlated with its test accuracy.

Table 1: **Correlation between types of training margins and group test accuracy.** We compute the Pearson correlation coefficient between the test group accuracy of a ResNet-18 and the minimum/average margin of each group across 3 experimental seeds. We see that average training margin is much more correlated with group test accuracy than minimum training margin. We leave an explanation of the predictive power of average margin for future work.

Margin Metric	Waterbirds	CelebA
Minimum Margin	-0.132	0.093
Average Margin	0.668	0.789

Second, we study an extension of the *implicit* group-balancing discussion from Section 3.2 to *explicit* group balancing. As discussed in Section 3.2, we investigate three explicit balancing methods: subsetting, upsampling, and upweighting. It has previously been observed that using upsampling or upweighting for class-balancing can result in drastic decreases to test WGA during training (LaBonte et al., 2024). If the performance improvement from LLR is indeed primarily due to the implicit group-balancing achieved by class-balancing, we would expect similar collapse when we explicitly group-balance using these same techniques.

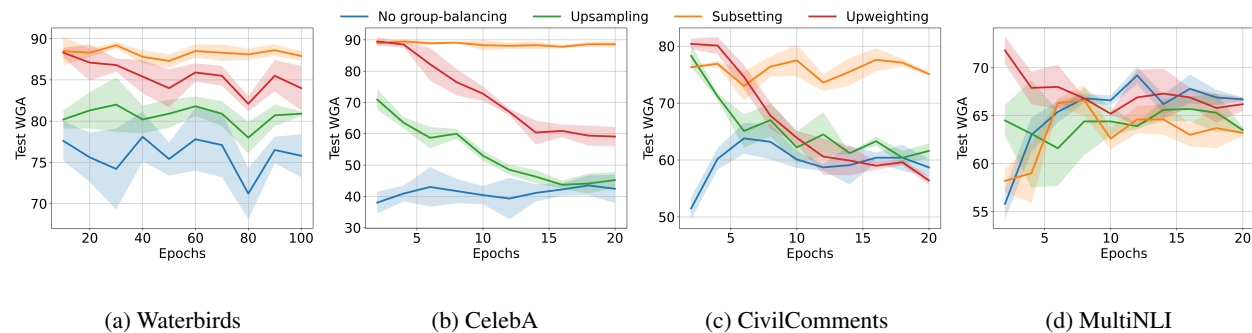


Figure 6: **Group-balancing with upsampling and upweighting can lead to catastrophic collapse.** We compare three group-balancing methods: *subsetting*, *upsampling*, and *upweighting*. We plot the mean and standard deviation over 3 experimental seeds for a ResNet-50 on the vision datasets and a BERT-Base on the language datasets. We note a dramatic decrease in test WGA during training for both upsampling and upweighting on CelebA and CivilComments. We also see a collapse in WGA for upweighting on MultiNLI.

Figure 6 shows the effects that different group-balancing methods have on WGA during training. We find that similarly to class-balancing, upsampling and upweighting lead to catastrophic collapse on CelebA and CivilComments. Additionally, group-balanced upweighting leads to a catastrophic collapse on MultiNLI and a mild decrease in WGA on Waterbirds. Overall, the training dynamics which result from the choice of subsetting, upweighting, and upsampling for group-balancing are quite similar to the dynamics observed by (LaBonte et al., 2024) for class-balancing. The notable exceptions are MultiNLI, which is class balanced *a priori* but not group balanced, and subsetting on Waterbirds. Class balanced subsetting on Waterbirds overly penalizes the minority group in the majority class which is problematic due to Waterbirds’ small size. Group balanced subsetting does not have this issue and hence performs very well on Waterbirds.

B ADDITIONAL DATASET AND TRAINING DETAILS

B.1 DATASET DETAILS

We study four benchmarks for group robustness across vision and language tasks, outlined below and detailed in Table 2. Note that Waterbirds is the only dataset that has a distribution shift and MultiNLI is the only dataset which is class-balanced *a priori*.

- *Waterbirds* (Welinder et al., 2010; Wah et al., 2011; Sagawa et al., 2020) is an image dataset where birds are classified into land species (“landbirds”) or water species (“waterbirds”). The spurious feature is the image background: landbirds are more frequently associated with land backgrounds, and waterbirds are more often found with water backgrounds.⁵
- *CelebA* (Liu et al., 2015; Sagawa et al., 2020) is an image dataset where celebrities are categorized as either blond or non-blond. The spurious feature is gender, with a $16\times$ greater number of blond women than blond men in the training data.
- *CivilComments* (Borkan et al., 2019; Koh et al., 2021) is a language dataset where online comments are classified as either toxic or non-toxic. The spurious feature involves the presence of one of the following identity categories: male, female, LGBT, black, white, Christian, Muslim, or other religion.⁶ Toxic comments tend to contain one of these identity categories more often than non-toxic comments, and vice versa.
- *MultiNLI* (Williams et al., 2018; Sagawa et al., 2020) is a language dataset where pairs of sentences are classified as a contradiction, entailment, or neither. The spurious feature is the presence of a negation in the second sentence — contradictions tend to have negations more often than entailments or neutral pairs.

B.2 TRAINING DETAILS

We utilize ResNet-18 and ResNet-50 (He et al., 2016) models pretrained on ImageNet-1K (Russakovsky et al., 2015) for Waterbirds and CelebA. We also utilize a BERT-Base (Devlin et al., 2019) model pretrained on Book Corpus (Zhu et al., 2015) and English Wikipedia for CivilComments and MultiNLI. These pretrained models are used as the initialization for ERM finetuning under the cross-entropy loss. We use standard ImageNet normalization with standard flip and crop data augmentation for the vision tasks and BERT tokenization for the language tasks (Izmailov et al., 2022). Our implementation uses the following packages: NumPy (Harris et al., 2020), PyTorch (Paszke et al., 2017; 2019), Lightning (Falcon & the PyTorch Lightning maintainers and contributors, 2019), TorchVision (TorchVision maintainers and contributors, 2016), Matplotlib (Hunter, 2007), Transformers (Wolf et al., 2020), and Milkshake (LaBonte, 2023).

To our knowledge, the licenses of Waterbirds and CelebA are unknown. CivilComments is released under the CC0 license, and information about MultiNLI’s license may be found in Williams et al. (2018).

Our experiments were conducted on two local 24GB Nvidia RTX A5000 GPUs. We provide our ERM finetuning hyperparameters in Table 3. Our LLR experiments were run for the same number of epochs as ERM on a held-out dataset (20% of the training set for Waterbirds and half the validation set for the other three datasets); we used SGD with learning rate 0.01 with no weight decay or learning rate schedule.

⁵It is worth noting that the Waterbirds dataset contains incorrect labels (Taghanaki et al., 2022). We report results based on the original, uncorrected version, as is customary in the literature.

⁶This version of CivilComments includes four identity groups, as used in this work and by Sagawa et al. (2020); Idrissi et al. (2022); Izmailov et al. (2022); Kirichenko et al. (2023); LaBonte et al. (2023). There exists another version where identity categories are not merged into one spurious feature; that version is employed by Liu et al. (2021); Zhang et al. (2022); Qiu et al. (2023). Both versions use the WILDS split (Koh et al., 2021).

572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623

Table 2: **Dataset composition.** The class probabilities change dramatically when conditioned on the spurious feature. Probabilities may not sum to 1 due to rounding.

Dataset	Group g		Training distribution \hat{p}			Data quantity		
	Class y	Spurious s	$\hat{p}(y)$	$\hat{p}(g)$	$\hat{p}(y s)$	Train	Val	Test
Waterbirds	landbird	land	.768	.730	.984	3498	467	2225
	landbird	water		.038	.148	184	466	2225
	waterbird	land	.232	.012	.016	56	133	642
	waterbird	water		.220	.852	1057	133	642
CelebA	non-blond	female	.851	.440	.758	71629	8535	9767
	non-blond	male		.411	.980	66874	8276	7535
	blond	female	.149	.141	.242	22880	2874	2480
	blond	male		.009	.020	1387	182	180
CivilComments	neutral	no identity	.887	.551	.921	148186	25159	74780
	neutral	identity		.336	.836	90337	14966	43778
	toxic	no identity	.113	.047	.079	12731	2111	6455
	toxic	identity		.066	.164	17784	2944	8769
MultiNLI	contradiction	no negation	.333	.279	.300	57498	22814	34597
	contradiction	negation		.054	.761	11158	4634	6655
	entailment	no negation	.334	.327	.352	67376	26949	40496
	entailment	negation		.007	.104	1521	613	886
	neither	no negation	.333	.323	.348	66630	26655	39930
	neither	negation		.010	.136	1992	797	1148

Table 3: **Default ERM and LLR hyperparameters.**

Dataset	Optimizer	Initial LR	LR schedule	Batch size	Weight decay	Epochs
Waterbirds	AdamW	1×10^{-5}	Cosine	32	1×10^{-4}	100
CelebA	AdamW	1×10^{-5}	Cosine	32	1×10^{-4}	20
CivilComments	AdamW	1×10^{-5}	Linear	32	1×10^{-4}	20
MultiNLI	AdamW	1×10^{-5}	Linear	32	1×10^{-4}	20