000	PATHOLOGIES OF OUT-OF-DISTRIBUTION DETECTION
001	
002	
003	Anonymous authors
004	Paper under double-blind review
005	
006	
007	Abstract
800	
009	There is a proliferation of out-of-distribution (OOD) detection methods in deep
010	learning which aim to detect distribution shifts and improve model safety. These methods often rely on supervised learning to train models with in distribution
011	data and then use the models' predictive uncertainty or features to identify OOD
012	points. In this paper, we critically re-examine this popular family of OOD detection
013	procedures, revealing deep-seated pathologies. In contrast to prior work, we argue
014	that these procedures are <i>fundamentally answering the wrong question</i> for OOD
015	detection, with no easy fix. Uncertainty-based methods incorrectly conflate high
010	uncertainty with being OOD, and feature-based methods incorrectly conflate far
010	feature-space distance with being OOD. Moreover, there is no reason to expect a
010	classifier trained only on in-distribution classes to be able to identify OOD points;
019	shout the label of an airplane, which may share features with a cat that help
020	distinguish cats from dogs despite generally appearing nothing alike. We show
021	how these pathologies manifest as irreducible errors in OOD detection and identify
023	common settings where these methods are ineffective. Additionally, interventions
024	to improve OOD detection such as feature-logit hybrid methods, scaling of model
025	and data size, Bayesian (epistemic) uncertainty representation, and outlier exposure
026	also fail to address the fundamental misspecification.
027	1
028	I INTRODUCTION
029	In the real world distribution shifts are the norm rather than the exception. We almost always have to
030	deploy our predictive models on test points drawn from at least a somewhat different distribution than
031	the training points: images acquired from different machines and hospitals, lane boundary detection
032	in different cities, speech recognition with different accents (Amodei et al., 2016; Jung et al., 2021;
033	Niu et al., 2016; Zhou et al., 2022; Koh et al., 2021). In order to make good predictions under these
034	shifts, we need to build relevant <i>invariances</i> into our models, so that natural transformations such
035	as rotations, translations, or even mild noise corruptions, do not significantly change the predictive
036	distribution (Hendrycks and Dietterich, 2018; Mintun et al., 2021; Benton et al., 2020).
037	However, rather than generalizing to natural distribution shifts, it has become popular to detect
038	out-of-distribution (OOD) data by training a supervised predictive model on in-distribution data, and
039	then examining the model's uncertainty, logits, or features. A proliferation of works develop such
040	procedures for detection improvements on known benchmarks (e.g., Hendrycks and Gimpel, 2016;
041	Let et al., 2010; Ken et al., 2021; Hendrycks et al., 2019a; Liang et al., 2017; Wang et al., 2021) of propose new and more challenging benchmarks (a.g. Hendrycks et al., 2010a; Vang et al., 2024;
042	Wang et al. 2022: Bitterwolf et al. 2023: Yang et al. 2021)
043	
044	While some prior work has considered limitations of OOD detection, the focus has been on issues with
045	benchmarking (e.g., controlling for covariate shift versus semantic shift detection) (Yang et al., 2024),
045	Specific architectural features of generative models (e.g., normalizing flows with coupling layers) (Kirichenko et al. 2020), or minor method deficiencies that can be straightforwardly addressed (e.g.
047	max-logit being more robust to many classes than max softmax) (Hendrycks et al. 2019a)
048	De la complete de la
049	By contrast, we argue that the whole premise of using supervised models trained on in-distribution

data for out-of-distribution detection is fundamentally flawed — a wholly misspecified enterprise,
 with no easy fix. First, the predictive distribution is over class labels, not whether an input comes
 from a different distribution: it is simply answering a different question. By making a predictive
 distribution better for detection, we could be removing invariances that would help with generalization.
 This effort is particularly self-defeating if the reason for wanting to do detection in the first place is



(a) Feature-based pathologies

(b) Logit-based pathologies

Figure 1: There are irreducible errors when using supervised models for OOD detection because the problem is inherently misspecified. Supervised models can only determine if an input leads to atypical representations or uncertain predictions, which is fundamentally different than determining if the input belongs to the training distribution.

068 to defer to a system that could provide reasonable generalization (e.g., thresholded classification). 069 Moreover, there is often a focus on detecting *semantic shifts* (Yang et al., 2024), particularly new unseen classes, where the model in fact should not be expected to generally say anything reasonable. 071 A car-truck classifier could be highly confident that a dog is a truck, if it uses similar features to 072 distinguish trucks from cars. It is not that OOD detection is "fundamentally difficult", but rather that 073 detection is being approached with methods that are fundamentally answering the wrong question a dog on the whole may look much different than a truck, and should be distinguishable from trucks 074 and cars, but not by training a supervised model only to differentiate between trucks and cars. 075

076 In this paper, we explain the fundamental conceptual limitations of popular supervised approaches 077 for OOD detection, visualized in Figure 1, which we exemplify in several settings. We then consider 078 the limitations — and, in some cases, additional pathologies — of standard interventions, such as epistemic (Bayesian) uncertainty representation and ensembling (Lakshminarayanan et al., 2017; 079 Malinin et al., 2019; Tagasovska and Lopez-Paz, 2019; D'Angelo and Fortuin, 2021; Hendrycks et al., 2018; Pearce et al., 2021), introducing new unseen classes in the model predictive distributions 081 (e.g., Fort et al., 2021), and outlier exposure (Hendrycks et al., 2018; Thulasidasan et al., 2021; Roy et al., 2022; Choi et al., 2023). We also critically examine other complementary procedures, such as 083 generative models, which may appear to be more aligned with the question of OOD detection. 084

085 2 PRELIMINARIES

Let $f_{\theta}: \mathcal{X} \to \mathcal{Y}$ be a neural network with parameters θ which maps training data $X^{\text{tr}} \sim p_{\mathcal{X}}(\cdot)$ 087 to a class from $\mathcal{Y} = \{1, \ldots, K\}$. A model's decision function is derived from its predictive 088 distribution $f_{\theta}(x) = \arg \max_{k \in \{1, \dots, K\}} p_{\theta}(y = k | x)$. OOD detection methods, which leverage 089 trained supervised models, propose a scoring function which for a test example x^* assigns a scalar 090 value $s(x^*, f_{\theta}, \mathcal{D}_{tr})$ given a trained model f_{θ} and training data $\mathcal{D}_{tr} = \{X_i^{tr}, Y_i^{tr}\}_{i=1}^N$. The score 091 $s(x^*, f_{\theta}, \mathcal{D}_{tr})$ is compared to a threshold value to determine whether x^* will be detected as OOD or 092 not. These methods are typically evaluated by computing AUROC (area under the receiver operating characteristic curve) scores on distribution shift benchmarks. 094

Two particularly common families of approaches have emerged for such OOD detection. If we 095 view the model as a composition of transformations $p_{\theta}(y = c|x) = \operatorname{softmax}(c_{\theta} \circ e_{\theta}(x))_{c}$ where 096 $e_{\theta}: \mathcal{X} \to \mathcal{F}$ is the penultimate layer feature extractor, and $c_{\theta}: \mathcal{F} \to \mathbb{R}^{K}$ is the classification layer outputting logits, then there are two natural signals to consider — features or logits. 098

Feature-based approaches. These methods compute the OOD score based on the features, typically 099 from the penultimate layer. The most common approach is based on the squared Mahalanobis 100 Distance (Maha) (Lee et al., 2018), where we fit a class-conditional Gaussian Mixture Model (GMM) 101 to our features with $\mu_c = \frac{1}{N_c} \sum_{i:y_i=c} e(x_i), \Sigma = \frac{1}{N} \sum_{c=1}^{K} \sum_{i:y_i=c} (e(x_i) - \mu_c) (e(x_i) - \mu_c)^{\top}$. The Mahalanobis score is then computed from the negative of the squared Mahalanobis distance as 102 103

$$s_{\text{Maha}}(x) = -\min_{c} \|\mu_{c} - e(x)\|_{\Sigma}^{2} = -\min_{c} (x - \mu_{c})\Sigma^{-1}(x - \mu_{c})^{\top},$$

105 106

104

063

064

065

066

067

Ren et al. (2021) extends this work and proposes Relative Mahalanobis Distance, which computes a 107 likelihood ratio between the most likely class-conditional Gaussian and an unconditional Gaussian fit over all train data with $\mu_{\text{train}} = \frac{1}{N} \sum_{i} e(x_i)$ and $\sum_{\text{train}} = \frac{1}{N} \sum_{i} (e(x_i) - \mu_{\text{train}}) (e(x_i) - \mu_{\text{train}})^{\top}$. The Relative Mahalanobis score is $s_{\text{RelMaha}}(x) = -\min_{c} \|\mu_{c} - e(x)\|_{\Sigma}^{2} - \|\mu_{\text{train}} - e(x)\|_{\Sigma_{\text{train}}}^{2}$. Many other feature-based approaches have also been proposed (Sun et al., 2022; Tack et al., 2020; Sehwag et al., 2021).

Logit-based approaches. These methods operate on the logits of a trained supervised model. The most common approach is *Maximum Softmax Probability* (MSP) (Hendrycks and Gimpel, 2016) $s_{msp}(x) = \max_c p_{\theta}(y = c|x)$. Other popular approaches within this family include the entropy of the predictive distribution $p_{\theta}(y|x)$ (Ren et al., 2019), value of the max logit (Jung et al., 2021) and the energy score (Liu et al., 2020).

Despite a proliferation of methods, simple approaches such as MSP tend to provide state-of-the-art results, even on the more sophisticated benchmarks (Hendrycks et al., 2019a; Yang et al., 2024).
For example, in Table 1 of Yang et al. (2024), it is observed that *"the results confirm that MSP still outperforms all modern methods"*. These methods are thus a natural choice to exemplify the broad conceptual issues with OOD detection, since they provide simple, popular, and still highly competitive approaches.

OOD detection task. These OOD detection methods can be used on various types of OOD detection.
 Much work has focused on using supervised models to identify points where the model does not have any chance of a correct label, often referred to as *semantic shift*, as opposed to label-preserving *covariate shift*. This type of semantic shift detection is often further categorized into *near OOD*, where the points are similar, and *far OOD* for more distinct inputs.

3 RELATED WORK

129

130

131 While anomaly and outlier detection has been studied for many decades in statistics, the related but distinct area of *out-of-distribution* (OOD) detection in deep learning is surprisingly new. Amodei et al. 132 (2016) provides a call to action to build methods that are robust to distribution shifts. Shortly after, 133 Hendrycks and Gimpel (2016) proposed using softmax uncertainty as a simple baseline to detect 134 out-of-distribution (OOD) points. A proliferation of methods followed, using the logits, features, 135 or uncertainty of a supervised model trained on in-distribution data to detect out-of-distribution 136 points, achieving better results on benchmark detection tasks (Lee et al., 2018; Liang et al., 2017; 137 Wang et al., 2022; Sun et al., 2022). Other work has focused on introducing new benchmarks with 138 higher resolution images, or test detection, more specifically under semantic shift (e.g., new unseen 139 classes) versus covariate shift (label-preserving transformations) (Yang et al., 2024; Bitterwolf et al., 140 2023; Huang and Li, 2021). There are also many interventions for boosting performance, including 141 Bayesian uncertainty representation (Lakshminarayanan et al., 2017; Malinin et al., 2019; Tagasovska 142 and Lopez-Paz, 2019; D'Angelo and Fortuin, 2021; Rudner et al., 2022), confidence minimization 143 and outlier exposure (Hendrycks et al., 2018; Papadopoulos et al., 2021; Thulasidasan et al., 2021), and pre-training (Fort et al., 2021; Tran et al., 2022; Hendrycks et al., 2019b). 144

145 While there are several works critical in some way of OOD detection, our focus is significantly 146 different. Critiques tend to be targeted at modifications to existing measures (e.g., max-logit has 147 fewer false positives than MSP) (Hendrycks et al., 2019a), improving the benchmark data (e.g., higher resolution data, data with many classes, and more cleanly separating semantic shift from 148 covariate shift) (Hendrycks and Dietterich, 2019; Yang et al., 2024), specific architectural properties 149 of generative models (e.g., coupling layers in normalizing flows) (Kirichenko et al., 2020), or note 150 that detection might need to be more tailored to specific shifts (Tajwar et al., 2021; Farquhar and Gal, 151 2022). By contrast, we examine whether the predominant approach of training supervised models on 152 in-distribution data for OOD detection is *fundamentally misspecified*, answering a different question 153 than "is this point out-of-distribution?" We conceptually elucidate significant pathologies of both 154 feature and logit-based approaches to OOD detection, and then exemplify these pathologies. We also 155 show that interventions such as Bayesian (epistemic) uncertainty have their own pathological behavior, 156 despite being considered the principled approach to OOD detection in prior work. We further show 157 that other interventions, such as confidence minimization, can introduce a trade-off between detection 158 and generalization. Moreover, we consider whether generative models are more directly answering the question "is this point from a different distribution?" and also evaluate the deficiencies of such 159 approaches. We also directly contrast these approaches, and simple statistical baselines, with the 160 supervised methods. Viewing these issues through the lens of misspecification, we finally consider 161 interventions that can help reduce misspecification to improve detection performance.



Figure 2: Feature-based methods have two key failure modes: indistinguishable features and irrelevant features. (Left): Oracle classifier achieves only around 90% AUROC on ImageNet vs ImageNet-OOD, indicating some OOD inputs have *indistinguishable features*. Oracle PCA projection improves Mahalanobis AUROC, showing many features are *irrelevant for OOD detection*. (Middle): Error decomposition into irreducible, suboptimal feature selection, and other components. (Right): Top discriminating features are OOD dataset-specific. PCA on ID and OOD_A improves AUROC for OOD_A (solid line, name in title) but decreases it for OOD_B (dashed lines) for ViT-S/16 features.

4 OOD DETECTION METHODS ANSWER THE WRONG QUESTIONS

Many OOD detection methods rely on the features or logits from supervised models that are only exposed to in-distribution data. Even though these approaches are sometimes able to achieve reasonable results on OOD detection benchmarks, they fundamentally answer the wrong question: instead of determining whether an input belongs to the training distribution or some different distribution, they instead ask if the input leads to atypical model representations or unconfident predictions. In this section, we explore the concrete instances where the answers to these two questions differ, and we demonstrate that feature and logit-based OOD detection methods have irreducible errors as a result.

188 4.1 FEATURE-BASED METHODS

171

172

173

174

175

176 177

178

187

189 Feature-based methods typically use distance metrics to measure how close the features of the test 190 input are to the features of the train inputs, answering the question "does this input lead to features that are far from the features seen during training?". These methods have two fundamental failure modes: 191 1) the learned features do not sufficiently discriminate between OOD and ID inputs, and 2) the optimal 192 distance metric depends on the OOD data, forcing these methods to use suboptimal, heuristic-based 193 distance metrics given only access to ID data. In particular, only the distance information along a 194 small number of feature dimensions is useful for OOD detection, but it is impossible to infer these 195 most discriminating features from irrelevant features without access to OOD data. 196

OOD features can be indistinguishable from ID features. While OOD inputs generally have
 unique characteristics that distinguish them from ID data, a supervised model may not be incentivized
 to learn these features if they are unhelpful for ID classification. If the OOD and ID features are
 indistinguishable, then no feature-based methods can perform well. This failure mode may have
 especially significant impacts for near OOD detection where fine-grained features are required.

202 To demonstrate this lack of separability between ID and OOD features, we study four different models 203 trained on ImageNet-1k: ResNet-18, ResNet-50, ViT-S/16, and ViT-B/16, with the OOD datasets of 204 ImageNet-OOD (Yang et al., 2024), Textures (Cimpoi et al., 2014), and iNaturalist (Van Horn et al., 2018). For each setting, we train an Oracle, a binary linear classifier, to differentiate between ID 205 features and OOD features and report its performance on held-out ID and OOD features. This Oracle 206 serves as a proxy for the best possible performance of any feature-based OOD detection method since 207 it is directly trained on both ID and OOD features, unlike any realistic methods. We see in Figure 2 208 (left) that even with ground-truth OOD information, the Oracle is unable to clearly disambiguate 209 between ID and OOD examples on challenging OOD datasets, obtaining AUROCs as low as 0.86. 210 For each model, (1– Oracle AUROC) represents an irreducible error: no feature-based method can 211 correctly detect these OOD inputs that have indistinguishable features from ID. 212

Irrelevant features hurt performance and are impossible to fully identify. Even if the model
 has learned features that discriminate between OOD and ID data, it is generally impossible to identify
 which features discriminate between OOD and ID data and which features are irrelevant without
 access to OOD data. As a result of the underspecification of OOD data at train-time, feature-based

methods must resort to suboptimal distance metrics to compare OOD features from ID features that
 do not sufficiently up-weight discriminating features or down-weight irrelevant features.

218 We illustrate this failure mode with the features from ResNet-18, ResNet-50, ViT-S/16, and ViT-B/16 219 trained on ImageNet-1k, and the Mahalanobis (Maha) method, which uses a distance metric defined 220 by the empirical covariance matrix Σ of ID features. We compare the performance of Maha before and 221 after an Oracle PCA projection that preserves only the most discriminating dimensions between OOD 222 and ID features, computed by performing PCA on both ID and OOD features and using the number of 223 PCA components among $\{32, 64, 128, 256\}$ that maximizes the resulting Maha AUROC. In Figure 2 224 (left), we show the addition of an Oracle PCA projection significantly improves Maha performance 225 on all models by an average of over 10 percentage points. Moreover, we see in Figure 2 (middle) 226 that performing this PCA projection accounts for nearly all of the reducible error of Maha for the ViT models. In other words, for ViTs, the gap between Maha and the best possible performance 1) is 227 almost entirely explained by the use of irrelevant features in the distance computation, and 2) requires 228 information unavailable to any feature-based method. While methods such as Relative Mahalanobis 229 and ViM (Ren et al., 2021; Wang et al., 2022) use PCA projections or related ideas to attempt to 230 reduce the impact of irrelevant features, they can only use feature covariances computed on ID data 231 alone, and thus do not address this fundamental limitation as we show in Appendix A.1. 232

In Figure 2 (right), we show that the Oracle PCA projection is highly specific to the particular OOD dataset we wish to detect and does not transfer between OOD datasets. For example, as demonstrated in the first panel, using the top 32 PCA components computed on IN and IN-OOD improves Maha AUROC in detecting IN-OOD but significantly degrades the AUROC for detecting Textures and iNaturalist, using features from ViT-S/16. This result shows that, as long as the OOD dataset is not specified at training time, removing the influence of irrelevant features is impossible for any feature-based method, presenting another fundamental bottleneck to its detection performance.

Visual demonstrations. We visualize clear examples of failure modes for feature-based methods 240 in Appendix A.1. To demonstrate feature overlap, we train a ResNet-18 on a subset of CIFAR-10 241 classes: airplane, cats, and trucks. We then use this trained model to detect OOD images of dogs. We 242 see in Figure A.1 (left) that the feature space between cats and dog have very high overlap, since the 243 model did not learn the features necessary to distinguish between these two classes. This pathology is 244 reflected in the poor performance of feature-based methods such as Mahalanobis distance, which only 245 achieves an AUROC of 0.537 and is barely better than random chance. Furthermore, these failures 246 also occur in larger models trained on diverse datasets. Even when using a ResNet-50 trained on 247 ImageNet-1K, Figure A.1 (right) demonstrates that feature-based methods like Mahalanobis distance 248 fail to correctly differentiate ID from OOD examples and assign low distances to OOD inputs. 249

250 4.2 LOGIT-BASED METHODS

251 Due to the many pathologies of feature-based OOD detection methods, it may be tempting to instead 252 focus on logit-based methods, which gauge a model's uncertainty over an input's predicted labels via 253 its logits. However, the previous limitations are still applicable. For instance, in the scenario where OOD and ID features overlap, logit-based methods would also fail to detect OOD inputs since the 254 logits are a function of the penultimate-layer features. Furthermore, logit-based methods have their 255 own suite of failure modes which arise from the conflation of *label uncertainty*, the uncertainty over 256 the correct ID label, with OOD uncertainty, the uncertainty over whether the sample is ID or OOD. 257 Logit-based methods heuristically assume that higher label uncertainty is equivalent to higher OOD 258 uncertainty, but these are fundamentally different quantities. As a result, there are two distinct failure 259 modes where logit-based methods make the incorrect prediction: instances where ID data naturally 260 has high label uncertainty, and instances where OOD data has low label uncertainty. 261

ID examples often have high uncertainty. To show the misalignment between label and OOD uncertainty, we demonstrate instances where models predict high label uncertainty over in-distribution samples. One example of this failure mode can be found in ImageNet-1K, where it is known that many of the images within the dataset contain concepts from multiple classes (Stock and Cisse, 2018; Shankar et al., 2020). We would expect these multi-label images to have high label uncertainty since there may be multiple correct answers. For our experiments, we used the human annotations from Beyer et al. (2020) as the ground truth for the number of labels corresponding to each image.

We explore the behaviors of ResNet-18, ResNet-50, ViT-S/16, ViT-B/16, all trained on ImageNet-1k, as well as ViT-G/14 DINOv2 pretrained on internet-scale data, for uncertainty-based OOD detection.





271

272

273

274

275

(a) Label uncertainty is high but input is ID

(b) Label uncertainty is low but input is OOD

ID

OOD "Striped"

Figure 3: Logit-based methods incorrectly conflate label uncertainty with OOD uncertainty.
(Left): ID images with multiple correct labels should have high label uncertainty. Each connected line shows the decrease in OOD detection for models listed in the right panel when focusing on this subset of high-uncertainty ID data. (Right): All methods assign low uncertainty to the OOD class
'Striped' from Textures and perform similarly to random chance.

285

When we apply uncertainty-based metrics to these samples where multiple labels may apply, we find 286 in Figure 3 that the average uncertainty of these multi-label images, denoted with Os, is significantly 287 higher than corresponding in-distribution samples, denoted with the connected Xs, across a variety of 288 methods. However, these images are clearly ID, since they are sampled from the same distribution 289 that the model was trained on. Furthermore, we can also see that logit-based methods are not able to 290 distinguish between ID inputs with high natural label uncertainty and OOD inputs; for example, the 291 AUROC for multi-label images (ID with high label uncertainty) vs ImageNet-OOD is only around 292 0.6. These results reveal that uncertainty-based methods are insufficient for OOD detection. 293

OOD examples often have low uncertainty. In Figure 3, we consistently find that logit-based approaches are unable to distinguish between ID and the "Striped" class from Textures across many settings. Furthermore, in Table A.1, we benchmark 14 different models including ResNets, ViTs, and ConvNext V2 in the setting where ImageNet-1K is ID. We record the FPR@95, which indicates how many OOD examples are incorrectly classified as ID due to their low uncertainty (false positive), at a threshold where 95% of ID examples are correctly classified. For logit-based methods such as MSP, max-logit, energy score, and entropy, the average FPR@95 across all settings is over 60%; thus, a majority of OOD examples are misclassified due to their low uncertainty.

We provide visual examples of these failure modes of uncertainty-based methods in Appendix A.2, where the predictive uncertainties of ID inputs are indistinguishable from the uncertainties of OOD inputs. In Figure A.3, we note how the uncertainties of an ID and OOD class entirely overlap for a LeNet-5 trained on a subset of CIFAR-10. We also visualize the feature space of a ResNet-50 trained on ImageNet-1k in Figure A.4 and find that the OOD class is often far from the decision boundary and has high model confidence, even though the examples are not from the input distribution.

Our experiments demonstrate that the difference between label uncertainty and OOD uncertainty, although easy to miss, is a fundamental limitation of logit-based OOD detection methods. This inherit misalignment of goals means no logit-based methods can overcome this pathology.

310 311

5 BUT WHAT ABOUT ...?

Given the prevalence of failure modes when using only feature or logit-based OOD methods, numerous strategies have been proposed to enhance OOD performance. In this section, we examine popular interventions such as combining feature and logit-based approaches, pre-training on larger datasets, modeling epistemic uncertainty, and exposing the model to outliers. For these methods, we analyze their limitations, and demonstrate how they fail to address the fundamental pathologies outlined in Section 4. We also address the limitations of explicitly including an OOD class during training and using unsupervised generative models.

320 5.1 SCALING MODEL AND DATA SIZE

Increasing model size and pre-training on large datasets have been shown to reliably improve OOD detection benchmarks as models tend to learn more diverse and higher-quality features (Fort et al., 2021; Dehghani et al., 2023; Miyai et al., 2023). When models see more diverse data and as the model capacity increases, they can learn more features that help distinguish OOD and ID data.



Figure 4: Scaling model size and training data does not address the fundamental limitations. Scaling from ResNet-18 trained on ImageNet (left-most point) to ViT-G/14 DINOv2 pre-trained on internet scale data (right-most), the Oracle AUROC still shows significant irreducible error for IN-OOD. Furthermore, a large fraction of the gap between the best performing method (described in text) and the Oracle can also be recovered by selecting features through Oracle PCA, indicating that the influence of irrelevant features is not addressed by scale.

However, as we show in Figure 4, scaling alone does not fully address the limitations of OOD 343 detection methods. We benchmark twelve different models of varying sizes and pretraining methods, 344 enumerated in Appendix B.3. First, in challenging near-OOD detection problems such as ImageNet 345 vs. ImageNet-OOD, models learn additional discriminating features between ID and OOD data at an 346 extremely slow rate, such that even the largest ViT-G/14 DINOv2 still has over 5% irreducible error 347 due to indistinguishable features. As we have argued in Section 4, this error can not be decreased 348 regardless of what OOD detection method we use. Indeed, the AUROC achieved by the best method 349 (Best) among Maha, Rel Maha, MSP, Max Logit, Energy, and ViM is consistently below the Oracle, 350 a binary classifier trained on ID and OOD features, by a wide margin. Second, while the error due 351 to indistinguishable features may be decreased (slowly) with scale and can become negligible on far-OOD detection problems such as ImageNet vs iNaturalist, there is still a large gap between the 352 best existing method and the Oracle as we scale the model. Much of this gap can be recovered by the 353 gain from optimally selecting features for Maha, represented by +Maha_{OraclePCA} (the gain is zero 354 if Best already outperforms Maha_{OraclePCA}), suggesting that while scaling allows the model to learn 355 features which almost perfectly discriminate ID and OOD data, the presence of irrelevant features 356 continues limit the performance. We provide additional empirical results in Appendix A.4 which 357 demonstrate the scaling behaviors of logit-based methods using 54 models over nine architectures 358 and six pre-training setups. These results demonstrate the fundamental limitations of existing OOD 359 detection methods even with increasing model and data size.

360 361

5.2 COMBINING FEATURE AND LOGIT-BASED METHODS

Hybrid approaches which combine model features and logits have been proposed for OOD detection (Sun et al., 2021; Wang et al., 2022), and methods like Virtual-logit Matching (ViM) (Wang et al., 2022) have achieved state-of-the-art results for certain OOD benchmarks. To understand the success of these methods, we test a simple hybrid method which sums the normalized scores of a feature-based method (Mahalanobis) with a logit-based method (MSP), to which we refer as "Hybrid-Add". We find in Appendix A.3 that for some models, "Hybrid-Add" improves OOD detection compared to using MSP or Mahalanobis alone on the Textures dataset, indicating feature and logit-based methods can have distinct failure modes.

370 However, hybrid methods do not address the fundamental pathologies caused by the model mis-371 specification. In the many cases where ID and OOD features are indistinguishable, as described in 372 Section 4.1, hybrid methods are equally unable to differentiate OOD examples because both features 373 and logits will overlap. Furthermore, the usefulness of hybrid methods is largely dataset-dependent. 374 For instance, we find that ViM usually outperforms both MSP and Mahalanobis for OOD detection 375 on Textures (Figure A.6) but does not offer a consistent advantage on IN-OOD (Figure 5). Our simple "hybrid-add" method does not offer a clear and consistent advantage over MSP or Mahalanobis on 376 either IN-OOD or iNaturalist. In practice, we find that feature and logit-based methods do not always 377 share distinct failure modes, and so hybrid methods may not be beneficial.



Figure 5: Hybrid OOD methods may not be beneficial. We compare two hybrid methods (ViM, Hybrid-Add) against MSP and Maha.



Figure 6: Training a ResNet-18 with outlier exposure hurts OOD generalization for covariate shifts compared to standard training.

5.3 EXPOSING TO OUTLIERS

Another popular approach to improve OOD detection is outlier exposure, which incorporates OOD examples when training the model (Hendrycks et al., 2018; Choi et al., 2023). In this setting, we explicitly optimize the model to have high uncertainty on the outlier dataset:

$$\mathcal{L}_{\rm CE} + \alpha \mathcal{L}_{\rm OE} = \mathbb{E}_{(x,y)\sim\mathcal{D}_{\rm in}} \ell_{\rm CE}(f(x), y) + \alpha \mathbb{E}_{x'\sim\mathcal{D}_{\rm out}} \ell_{\rm CE}(f(x'), y_u)$$

where y_u is uniform distribution over all K classes. Outlier exposure relies on the diverse dataset D_{out} in order to encourage the model to generally have high predictive uncertainty away from the training data and improve detection with predictive-space methods like MSP. However, even if the model is exposed to OOD data during training, the final model is still misspecified because it only contains ID classes as possible categorizations. As previously discussed in Section 4.1, OOD datasets are quite diverse, and the features necessary to distinguish ID from one OOD dataset often do not generalize to other types of OOD.

Furthermore, *outlier exposure may significantly hurt OOD generalization* because the model is explicitly trained to have high label uncertainty over a large set of inputs; this degradation in performance is especially problematic because OOD generalization is essential for model robustness and reliability. To demonstrate this behavior, we compare two ResNet-18 models trained on CIFAR-10, one with the standard training regime and the other with outlier exposure using TIN-597 as D_{out} following Zhang et al. (2023) (see Appendix B.1 for setup details).

410 In Appendix A.6, we show that outlier exposure does improve OOD detection for most of the semantic 411 shift OOD benchmarks. However, outlier exposure does not improve performance on MNIST, likely 412 because this dataset differs significantly from D_{out} and other natural image benchmarks. This 413 decreased performance highlights the sensitivity of outlier exposure to the choice of OOD data and 414 reiterates that the features which distinguish ID and OOD are not consistent across diverse OOD 415 datasets. Furthermore, we find that while the ID accuracy of the outlier exposed model is negatively impacted, the impacts of outlier exposure on OOD generalization is significantly worse. In Figure 6, 416 on inputs with covariate shifts, outlier exposure hurts the model's accuracy by over 10% across all of 417 our benchmarked datasets. Thus, by explicitly encouraging high uncertainty on the diverse outlier 418 dataset, we sacrifice the generalizability of our model. 419

420 421

378

379

380

381

382

383

384

385

386

387

389 390

391 392

393

394 395 396

5.4 MODELING EPISTEMIC UNCERTAINTY

Predictive uncertainty can be separated into *aleatoric uncertainty*, which is considered irreducible
and stems from inherent data variability, and *epistemic uncertainty*, which is uncertainty over which
solution is correct given the limited data. It has been posited that focusing on epistemic uncertainty
for predictive models is the principled approach to OOD detection because the uncertainty increases
as we move away from the data, and indeed there is a proliferation of methods approximating
epistemic uncertainty for improved performance on OOD detection benchmarks (e.g., Band et al.,
2021; D'Angelo and Fortuin, 2021; Lakshminarayanan et al., 2017; Malinin et al., 2019; Rudner
et al., 2022; Tagasovska and Lopez-Paz, 2019; Tran et al., 2022).

429

Epistemic uncertainty is typically represented through a distribution over the model parameters. For a model f with stochastic parameters Θ , distributed according to $q(\theta)$, we can express the model's predictive uncertainty as a combination of aleatotric and epistemic uncertainty,



Figure 7: Epistemic uncertainty becomes less useful for OOD detection as more in-distribution data is observed. We consider three CIFAR-10 classes as ID, and CIFAR-100 as OOD. We use a ResNet-18 with a last-layer Laplace approximation to measure epistemic uncertainty. As we increase the in-distribution training examples, the posterior collapses and epistemic uncertainty diminishes.

$$\underbrace{\mathcal{H}\left(\mathbb{E}_{q\Theta}[p(y \mid \mathbf{x}, \Theta)]\right)}_{\text{Total Uncertainty}} = \underbrace{\mathbb{E}_{q\Theta}[\mathcal{H}(p(y \mid \mathbf{x}, \Theta))]}_{\text{Aleatoric Uncertainty}} + \underbrace{\mathcal{I}(Y; \Theta)}_{\text{Epistemic Uncertainty}}, \tag{1}$$

where $\mathcal{H}(\cdot)$ is the entropy functional and $\mathcal{I}(Y;\Theta)$ is the mutual information. The predictive distribution, through the rules of probability, is then $p(y = c | \mathbf{x}, \mathcal{D}) = \int \operatorname{softmax}(f_{\theta}(\mathbf{x}))_c \cdot p(\theta | \mathcal{D}) d\theta$. We note that *deep ensembling* procedures (Lakshminarayanan et al., 2017), particularly popular for OOD detection, are a prominent example of epistemic uncertainty representation; by marginalizing over modes in a posterior they often provide a relatively accurate representation of the posterior predictive distribution (Wilson and Izmailov, 2020; Izmailov et al., 2021b).

However, as we have previously discussed, the predictive uncertainty is not over whether a point 456 is OOD, but rather over class labels. Epistemic uncertainty does not address this fundamental 457 misspecification. For a clear demonstration of the conceptual difference between epistemic uncertainty 458 and OOD detection, consider how epistemic uncertainty changes as a function of data size. In the 459 infinite ID-data limit, the epistemic uncertainty of a model approaches zero and the model becomes 460 extremely confident in its parameters. If measuring epistemic uncertainty were the correct approach 461 to OOD detection, then having such low epistemic uncertainty implies that OOD points do not exist 462 in this setting. Therefore, because perfectly capturing epistemic uncertainty is not enough to solve 463 OOD detection, they must answer fundamentally different questions. In fact, as the model sees more in-distribution data during training, its ability to detect OOD inputs may become worse! Given the 464 growing availability of large datasets, this behavior becomes increasingly problematic. 465

466 To illustrate this phenomenon, we consider a last-layer Bayesian approximation (Kristiadi et al., 2020) 467 and train a linear layer f_{θ} over features extracted from a ResNet-18 trained on IN-1K to classify three 468 classes: airplane, dog, and truck. We place a prior over parameters θ , and we approximate the predic-469 tive distribution through a Laplace approximation that uses a Gaussian distribution to approximate the posterior distribution of the model parameters, allowing for the estimation of epistemic uncertainty 470 (MacKay, 2003). In Figure 7, we visualize the learned decision boundaries by applying PCA to 471 reduce the three-dimensional logit space to two dimensions and plot the the maximum softmax 472 probabilities over this projection. As the size of the training data increases, the posterior noticeably 473 contracts, and the performance of the Laplace model decreases to approach the performance of the 474 deterministic model with maximum a posteriori (MAP) parameter estimates. 475

In other work, Izmailov et al. (2021a) studies pathologies in Bayesian methods for OOD generalization
(rather than detection), and D'Angelo and Henning (2021) note the sensitivity to prior specification
in using BNNs for OOD detection.

479 480

481

442

443

444

445 446 447

448

5.5 INTRODUCING AN UNSEEN CLASS

Since standard classification models trained over K classes are fundamentally misspecified for the task of detecting OOD classes in both feature-space and logit-space, it may be tempting to correct the specification problem by adding a (K + 1)-th class corresponding to the OOD category (Thulasidasan et al., 2020). During training, we can then expose models to OOD examples, and use this additional class for OOD detection on test samples.

498



Figure 8: Adding an OOD class is only 496 effective if the train OOD examples are 497 similar to test OOD examples.



Figure 9: Generative models are not optimal OOD detectors: for generative model $\mathcal{N}(\mu, 1)$, the optimal μ is 0 to model the ID data but $-\infty$ for OOD detection.

499 However, we find this method is only effective when the examples that the model is exposed during 500 training are very similar to the OOD examples during test-time, which is often unrealistic. To 501 demonstrate, we train a ResNet-18 model on two CIFAR-100 classes: keyboard and porcupine, and 502 use samples from cup and skyscraper for the OOD class. We then measure the performance of OOD detection over the remaining CIFAR-100 classes. In Figure 8, we use BERT embeddings (Devlin 504 et al., 2018) to compute the cosine similarity of the test-time OOD classes to the train-time ID and 505 OOD classes. We see that the OOD class was effective in capturing test-time examples of bottle and can, since they are similar to the train OOD examples. However, the model is unable to accurately 506 categorize examples like hamster and mouse, which are more closely related to ID classes. 507

508 5.6 USING GENERATIVE MODELS

509 Unlike previously mentioned methods utilize a supervised classification problem to make predictions, 510 unsupervised generative models trained on the in-distribution dataset attempt to directly answer the 511 question of how likely it is that sample x belongs to the training distribution. Generative models, 512 therefore, may appear to be a principled and natural solution to OOD detection. 513

However, better generative models are not always better OOD detectors. Since p(x) answers a 514 fundamentally different question than p(OOD|x), there is generally a conflict between creating a 515 better model for p(x) and the ability to use the likelihood from that model to detect OOD points. 516 We illustrate this phenomenon with a simple 1D example in Figure 9, where the ID data is drawn 517 from $\mathcal{N}(0,1)$ and the OOD data is drawn from $\mathcal{N}(2,1)$. Suppose we model the ID data with 518 $x \sim p_{\mu}(x) = \mathcal{N}(x|\mu, 1)$, where μ is the parameter of our model. Choosing $\mu = 0$ will exactly 519 model the true distribution and achieve the highest likelihood. However, as shown in Figure 9 (right), 520 the optimal choice of μ for OOD detection is $-\infty$, achieving a maximum AUROC but infinite 521 KL divergence from the true ID distribution. We further demonstrate this misalignment between 522 likelihood on the ID data vs OOD detection in Appendix A.5, discussing additional limitations and illustrating the failures of generative approaches such as diffusion models for OOD detection. 523

6 DISCUSSION

526 Fundamentally, we have shown that popular OOD detection procedures, both supervised and genera-527 tive, are often answering a different question than is this unlabeled point from a different distribution? 528 Moreover, interventions like outlier exposure hurt the ability for a model to generalize on covariate 529 shifts which begs the question: why are we doing OOD detection in the first place? If it is really to 530 detect OOD points, then the procedures we are using are severely misspecified and have fundamental limitations. If the goal of OOD detection is to be able to make more reasonable predictions under 531 covariate shifts (e.g., by deferring examples to another model through confidence thresholding), which 532 is arguably a more typical real-world use case than semantic shift, the interventions for detection can 533 be actively harmful. 534

Going forward, it will be important to identify real-world problems where one of generalization or detection is the focus — and in the cases where the ultimate objective really is detection, we should build approaches specifically designed to answer that question. 537

538

524

540 ETHICS STATEMENT

The authors have read and they acknowledge the ICLR Code of Ethics. The authors will strictlyadhere to the ICLR Code of Ethics.

545 REPRODUCIBILITY STATEMENT

Key details to reproduce experiments are provided in the main text and appendix. In addition, weintend to release the code publicly to reproduce all data and figures presented in the paper.

594 REFERENCES

596 597	Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. <i>arXiv preprint arXiv:1606.06565</i> , 2016.
598 599 600	Neil Band, Tim G. J. Rudner, Qixuan Feng, Angelos Filos, Zachary Nado, Michael W. Dusenberry, Ghassen Jerfel, Dustin Tran, and Yarin Gal. Benchmarking Bayesian Deep Learning on Diabetic Retinopathy Detection Tasks. In <i>Advances in Neural Information Processing Systems</i> 34, 2021.
601 602 603 604	Gregory Benton, Marc Finzi, Pavel Izmailov, and Andrew G Wilson. Learning invariances in neural networks from training data. <i>Advances in neural information processing systems</i> , 33:17605–17616, 2020.
605 606 607	James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. <i>Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf</i> , 2(3):8, 2023.
608 609 610	Lucas Beyer, Olivier J Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord. Are we done with imagenet? <i>arXiv preprint arXiv:2006.07159</i> , 2020.
611 612	Julian Bitterwolf, Maximilian Mueller, and Matthias Hein. In or out? fixing imagenet out-of- distribution detection evaluation. <i>arXiv preprint arXiv:2306.00826</i> , 2023.
613 614 615 616	Caroline Choi, Fahim Tajwar, Yoonho Lee, Huaxiu Yao, Ananya Kumar, and Chelsea Finn. Con- servative prediction via data-driven confidence minimization. <i>arXiv preprint arXiv:2306.04974</i> , 2023.
617 618 619	Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describ- ing textures in the wild. In <i>Proceedings of the IEEE conference on computer vision and pattern</i> <i>recognition</i> , pages 3606–3613, 2014.
620 621 622 623	Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In <i>Proceedings of the fourteenth international conference on artificial intelligence and statistics</i> , pages 215–223. JMLR Workshop and Conference Proceedings, 2011.
624 625	Francesco D'Angelo and Vincent Fortuin. Repulsive deep ensembles are bayesian. Advances in Neural Information Processing Systems, 34:3451–3465, 2021.
626 627 628	Francesco D'Angelo and Christian Henning. On out-of-distribution detection with bayesian neural networks. <i>arXiv preprint arXiv:2110.06020</i> , 2021.
629 630	Luke N Darlow, Elliot J Crowley, Antreas Antoniou, and Amos J Storkey. Cinic-10 is not imagenet or cifar-10. <i>arXiv preprint arXiv:1810.03505</i> , 2018.
631 632 633 634 635	Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. In <i>International Conference on Machine Learning</i> , pages 7480–7512. PMLR, 2023.
636 637	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. <i>arXiv preprint arXiv:1810.04805</i> , 2018.
638 639 640	Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. arXiv preprint arXiv:1605.08803, 2016.
641 642	Sebastian Farquhar and Yarin Gal. What'out-of-distribution'is and is not. In <i>NeurIPS ML Safety Workshop</i> , 2022.
643 644 645	Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. Exploring the limits of out-of-distribution detection. <i>Advances in Neural Information Processing Systems</i> , 34:7068–7081, 2021.
646 647	Priya Goyal, Mathilde Caron, Benjamin Lefaudeux, Min Xu, Pengchao Wang, Vivek Pai, Mannat Singh, Vitaliy Liptchinsky, Ishan Misra, Armand Joulin, et al. Self-supervised pretraining of visual features in the wild. <i>arXiv preprint arXiv:2103.01988</i> , 2021.

648 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common 649 corruptions and perturbations. arXiv preprint arXiv:1903.12261, 2019. 650 Dan Hendrycks and Thomas G Dietterich. Benchmarking neural network robustness to common 651 corruptions and surface variations. arXiv preprint arXiv:1807.01697, 2018. 652 653 Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution 654 examples in neural networks. arXiv preprint arXiv:1610.02136, 2016. 655 Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier 656 exposure. arXiv preprint arXiv:1812.04606, 2018. 657 658 Dan Hendrycks, Steven Basart, Mantas Mazeika, Andy Zou, Joe Kwon, Mohammadreza Mostajabi, 659 Jacob Steinhardt, and Dawn Song. Scaling out-of-distribution detection for real-world settings. 660 *arXiv preprint arXiv:1911.11132*, 2019a. 661 Dan Hendrycks, Kimin Lee, and Mantas Mazeika. Using pre-training can improve model robustness 662 and uncertainty. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, Proceedings of the 36th 663 International Conference on Machine Learning, volume 97 of Proceedings of Machine Learning 664 Research, pages 2712–2721. PMLR, 09–15 Jun 2019b. URL https://proceedings.mlr. 665 press/v97/hendrycks19a.html. 666 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in 667 neural information processing systems, 33:6840–6851, 2020. 668 669 Rui Huang and Yixuan Li. MOS: towards scaling out-of-distribution detection for large semantic 670 space. CoRR, abs/2105.01879, 2021. URL https://arxiv.org/abs/2105.01879. 671 Pavel Izmailov, Patrick Nicholson, Sanae Lotfi, and Andrew G Wilson. Dangers of bayesian model 672 averaging under covariate shift. Advances in Neural Information Processing Systems, 34:3309-673 3322, 2021a. 674 675 Pavel Izmailov, Sharad Vikram, Matthew D Hoffman, and Andrew Gordon Gordon Wilson. What are 676 bayesian neural network posteriors really like? In International conference on machine learning, 677 pages 4629-4640, 2021b. 678 Sanghun Jung, Jungsoo Lee, Daehoon Gwak, Sungha Choi, and Jaegul Choo. Standardized max 679 logits: A simple yet effective approach for identifying unexpected road obstacles in urban-scene 680 segmentation. In Proceedings of the IEEE/CVF International Conference on Computer Vision, 681 pages 15425-15434, 2021. 682 Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. Advances 683 in neural information processing systems, 34:21696–21707, 2021. 684 685 Polina Kirichenko, Pavel Izmailov, and Andrew G Wilson. Why normalizing flows fail to detect 686 out-of-distribution data. Advances in neural information processing systems, 33:20578–20589, 687 2020. 688 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-689 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. Wilds: A 690 benchmark of in-the-wild distribution shifts. In International conference on machine learning, 691 pages 5637-5664. PMLR, 2021. 692 693 Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being bayesian, even just a bit, fixes 694 overconfidence in relu networks. In International conference on machine learning, pages 5436-5446. PMLR, 2020. 696 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive 697 uncertainty estimation using deep ensembles. Advances in neural information processing systems, 30, 2017. 699 Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting 700 out-of-distribution samples and adversarial attacks. Advances in neural information processing systems, 31, 2018.

702 703 704	Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. <i>arXiv preprint arXiv:1706.02690</i> , 2017.
705 706	Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. <i>Advances in neural information processing systems</i> , 33:21464–21475, 2020.
707 708	David JC MacKay. <i>Information theory, inference and learning algorithms</i> . Cambridge university press, 2003.
709 710 711	Andrey Malinin, Bruno Mlodozeniec, and Mark Gales. Ensemble distribution distillation. <i>arXiv</i> preprint arXiv:1905.00076, 2019.
712 713 714 715	John P Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal Shankar, Percy Liang, Yair Carmon, and Ludwig Schmidt. Accuracy on the line: on the strong correlation between out-of-distribution and in-distribution generalization. In <i>International conference on machine learning</i> , pages 7721–7735. PMLR, 2021.
716 717 718	Eric Mintun, Alexander Kirillov, and Saining Xie. On interaction between augmentations and corruptions in natural corruption robustness. <i>Advances in Neural Information Processing Systems</i> , 34:3571–3583, 2021.
719 720 721	Atsuyuki Miyai, Qing Yu, Go Irie, and Kiyoharu Aizawa. Can pre-trained networks detect familiar out-of-distribution data? <i>arXiv preprint arXiv:2310.00847</i> , 2023.
722 723	Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Do deep generative models know what they don't know? <i>arXiv preprint arXiv:1810.09136</i> , 2018.
724 725 726 727	Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, and Balaji Lakshminarayanan. Detecting out- ofdistribution inputs to deep generative models using typicality. <i>arXiv preprint arXiv:1906.02994</i> , 2020.
728 729	Jianwei Niu, Jie Lu, Mingliang Xu, Pei Lv, and Xiaoke Zhao. Robust lane detection using two-stage feature extraction with curve fitting. <i>Pattern Recognition</i> , 59:225–233, 2016.
730 731 732	Aristotelis-Angelos Papadopoulos, Mohammad Reza Rajati, Nazim Shaikh, and Jiamian Wang. Outlier exposure with confidence control for out-of-distribution detection. <i>Neurocomputing</i> , 441: 138–150, 2021.
734 735	Tim Pearce, Alexandra Brintrup, and Jun Zhu. Understanding softmax confidence and uncertainty. arXiv preprint arXiv:2106.04972, 2021.
736 737	William Peebles and Saining Xie. Scalable diffusion models with transformers. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 4195–4205, 2023.
738 739 740 741 742	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR, 2021.
743 744 745	Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. <i>Advances in neural information processing systems</i> , 32, 2019.
746 747 748	Jie Ren, Stanislav Fort, Jeremiah Liu, Abhijit Guha Roy, Shreyas Padhy, and Balaji Lakshmi- narayanan. A simple fix to mahalanobis distance for improving near-ood detection. <i>arXiv preprint</i> <i>arXiv:2106.09022</i> , 2021.
749 750 751 752 753	Abhijit Guha Roy, Jie Ren, Shekoofeh Azizi, Aaron Loh, Vivek Natarajan, Basil Mustafa, Nick Pawlowski, Jan Freyberg, Yuan Liu, Zach Beaver, et al. Does your dermatology classifier know what it doesn't know? detecting the long-tail of unseen conditions. <i>Medical Image Analysis</i> , 75: 102274, 2022.
754 755	Tim G. J. Rudner, Zonghao Chen, Yee Whye Teh, and Yarin Gal. Tractable Function-Space Variational Inference in Bayesian Neural Networks. In <i>Advances in Neural Information Processing Systems 35</i> , 2022.

756 757 758 759	Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. <i>Advances in neural information processing systems</i> , 35:36479–36494, 2022.
760 761 762	Vikash Sehwag, Mung Chiang, and Prateek Mittal. Ssd: A unified framework for self-supervised outlier detection. <i>arXiv preprint arXiv:2103.12051</i> , 2021.
763 764 765	Vaishaal Shankar, Rebecca Roelofs, Horia Mania, Alex Fang, Benjamin Recht, and Ludwig Schmidt. Evaluating machine accuracy on imagenet. In <i>International Conference on Machine Learning</i> , pages 8634–8644. PMLR, 2020.
766 767 768 769	Pierre Stock and Moustapha Cisse. Convnets and imagenet beyond accuracy: Understanding mistakes and uncovering biases. In <i>Proceedings of the European conference on computer vision (ECCV)</i> , pages 498–512, 2018.
770 771	Yiyou Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations. Advances in Neural Information Processing Systems, 34:144–157, 2021.
772 773 774	Yiyou Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors. In <i>International Conference on Machine Learning</i> , pages 20827–20840. PMLR, 2022.
775 776 777	Jihoon Tack, Sangwoo Mo, Jongheon Jeong, and Jinwoo Shin. Csi: Novelty detection via contrastive learning on distributionally shifted instances. <i>Advances in neural information processing systems</i> , 33:11839–11852, 2020.
778 779 780	Natasa Tagasovska and David Lopez-Paz. Single-model uncertainties for deep learning. Advances in neural information processing systems, 32, 2019.
781 782	Fahim Tajwar, Ananya Kumar, Sang Michael Xie, and Percy Liang. No true state-of-the-art? ood detection methods are inconsistent across datasets. <i>arXiv preprint arXiv:2109.05554</i> , 2021.
783 784 785 786	Sunil Thulasidasan, Sushil Thapa, Sayera Dhaubhadel, Gopinath Chennupati, Tanmoy Bhattacharya, and Jeff Bilmes. A simple and effective baseline for out-of-distribution detection using abstention. 2020.
787 788 789 790	Sunil Thulasidasan, Sushil Thapa, Sayera Dhaubhadel, Gopinath Chennupati, Tanmoy Bhattacharya, and Jeff Bilmes. An effective baseline for robustness to distributional shift. In 2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA), pages 278–285. IEEE, 2021.
791 792 793 794 795 796	Dustin Tran, Jeremiah Liu, Michael W. Dusenberry, Du Phan, Mark Collier, Jie Ren, Kehang Han, Zi Wang, Zelda Mariet, Huiyi Hu, Neil Band, Tim G. J. Rudner, Karan Singhal, Zachary Nado, Joost van Amersfoort, Andreas Kirsch, Rodolphe Jenatton, Nithum Thain, Honglin Yuan, Kelly Buchanan, Kevin Murphy, D. Sculley, Yarin Gal, Zoubin Ghahramani, Jasper Snoek, and Balaji Lakshminarayanan. Plex: Towards Reliability Using Pretrained Large Model Extensions. In <i>ICML</i> <i>Workshop on Pre-training: Perspectives, Pitfalls, and Paths Forward</i> , 2022.
797 798 799 800 801	Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 8769–8778, 2018.
802 803 804	Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual- logit matching. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern</i> <i>recognition</i> , pages 4921–4930, 2022.
805 806 807	Haoran Wang, Weitang Liu, Alex Bocchieri, and Yixuan Li. Can multi-label classification networks know what they don't know? <i>Advances in Neural Information Processing Systems</i> , 34:29074–29087, 2021.
808 809	Andrew G Wilson and Pavel Izmailov. Bayesian deep learning and a probabilistic perspective of generalization. <i>Advances in neural information processing systems</i> , 33:4697–4708, 2020.

810 811 812	Jingkang Yang, Haoqi Wang, Litong Feng, Xiaopeng Yan, Huabin Zheng, Wayne Zhang, and Ziwei Liu. Semantically coherent out-of-distribution detection. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 8301–8309, 2021.
813 814 815 816 817	Jingkang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi Wang, Guangyao Chen, Bo Li, Yiyou Sun, et al. Openood: Benchmarking generalized out-of-distribution detection. <i>Advances in Neural Information Processing Systems</i> , 35:32598–32611, 2022.
818 819 820	William Yang, Byron Zhang, and Olga Russakovsky. Imagenet-ood: Deciphering modern out- of-distribution detection algorithms. In <i>International Conference on Learning Representations</i> , 2024.
821 822 823 824	Jingyang Zhang, Jingkang Yang, Pengyun Wang, Haoqi Wang, Yueqian Lin, Haoran Zhang, Yiyou Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, et al. Openood v1. 5: Enhanced benchmark for out-of-distribution detection. <i>arXiv preprint arXiv:2306.09301</i> , 2023.
825 826 827	Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning deep features for scene recognition using places database. <i>Advances in neural information processing systems</i> , 27, 2014.
828 829 830	Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A survey. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 45(4):4396–4415, 2022.
831	
832	
833	
834	
835	
830	
837	
030	
039 940	
04U 9/1	
8/12	
843	
844	
845	
846	
847	
848	
849	
850	
851	
852	
853	
854	
855	
856	
857	
858	
859	
860	
861	
862	
863	

TABLE OF CONTENTS

A	A Additional Empirical Results						
	A.1	Feature-based methods	17				
	A.2	Logit-based methods	18				
	A.3	Hybrid Methods	19				
	A.4	Effect of Pre-training	19				
	A.5	Generative models	21				
	A.6	Outlier Exposure	24				
В	Imp	lementation details	25				
	B .1	Outlier exposure experiment.	25				
	B .2	Evaluating pre-trained models.	25				
	B.3	Scaling Experiments	25				

A ADDITIONAL EMPIRICAL RESULTS

A.1 FEATURE-BASED METHODS

We provide visualizations of empirical examples of feature-based failures in Figure A.1, and we demonstrate that advanced methods like relative Mahalanobis and ViM are subject to the same failure modes in Figure A.2.



(a) ResNet-18 on CIFAR-10

(b) ResNet-50 on ImageNet-1k

Figure A.1: Visualizations of failure modes for feature-based OOD detection. (Left): We train a ResNet-18 on a subset of CIFAR-10, and find that the feature space between an ID class and OOD class have significant overlap. (**Right**): Feature-based methods also fail for larger models like ResNet-50 trained on ImageNet 1K, where OOD classes have low Mahalanobis distance.



Figure A.2: Relative Mahalanobis and ViM do not fully address the issue of irrelevant features on ImageNet vs ImageNet-OOD, especially with the more performant ViT models. In all cases, Mahalanobis with an Oracle PCA performs the best. Except for ResNets, Relative Mahalanobis and ViM offer negligible or negative improvement relative to Mahalanobis. The gap between Maha + Oracle PCA and the best-performing feature-based method is especially large for ViTs.

A.2 LOGIT-BASED METHODS

Logit-based methods fail when the uncertainty of ID data looks similar to the uncertainty of OOD data. We see an example in Figure A.3, where we find that the model very confidently classifies OOD dogs as ID trucks. In Table A.1, we find that well over half of OOD examples are misclassified as ID even with powerful pre-trained models. Figure A.4 visualizes a failure mode for a ResNet-50 trained on ImageNet, where 'Stripes' is often miscategorized as ID. In Figure A.5, we find that this failure mode is prevelant across a diverse set of models and logit-based OOD-detection methods.



(a) Confidence of ID vs OOD inputs

(b) Example inputs and model confidence

Figure A.3: **The predictive uncertainty of OOD points may be indistinguishable from ID points.** We train a LeNet5 to classify CIFAR10 automobiles and trucks, and we test the OOD dog class. We see that the model confidence for OOD dogs entirely overlaps with the ID truck class. In this setting, because the uncertainties are identical, no uncertainty-based method would be able to successfully differentiate ID from OOD.

	OOD Dataset	Model	MSP	Max Logit	Entropy	Energy Score
	IN-OOD	ResNet-50	0.774	0.804	0.820	0.778
	IN-OOD	ResNet-50 DINO	0.804	0.830	0.847	0.823
	IN-OOD	ResNet-34	0.807	0.824	0.838	0.809
	IN-OOD	ResNet-18	0.832	0.846	0.855	0.845
	IN-OOD	ViT-S/16	0.797	0.803	0.818	0.798
	IN-OOD	ViT-S/16 DINO	0.761	0.790	0.811	0.768
	IN-OOD	ViT-B/16	0.740	0.733	0.771	0.726
	IN-OOD	ViT-B/16 DINO	0.741	0.764	0.784	0.749
	IN-OOD	ViT-B/16 CLIP	0.764	0.776	0.805	0.726
	IN-OOD	ViT-B/14 DINOv2	0.658	0.621	0.638	0.610
	IN-OOD	ViT-G/14 DINOv2	0.562	0.448	0.450	0.469
j	IN-OOD	ViT-L/14 CLIP	0.686	0.685	0.723	0.631
	IN-OOD	ConvNeXt V2-B	0.701	0.708	0.773	0.673
	IN-OOD	ConvNeXt V2-L	0.696	0.710	0.787	0.663
	Textures	ResNet-50	0.662	0.544	0.522	0.594
	Textures	ResNet-50 DINO	0.681	0.624	0.612	0.637
	Textures	ResNet-34	0.690	0.562	0.533	0.620
	Textures	ResNet-18	0.710	0.571	0.527	0.643
	Textures	ViT-S/16	0.672	0.579	0.506	0.593
	Textures	ViT-S/16 DINO	0.612	0.400	0.363	0.521
	Textures	ViT-B/16	0.586	0.544	0.573	0.521
	Textures	ViT-B/16 DINO	0.531	0.351	0.307	0.437
	Textures	ViT-B/16 CLIP	0.657	0.530	0.538	0.564
	Textures	ViT-B/14 DINOv2	0.535	0.409	0.401	0.451
	Textures	ViT-G/14 DINOv2	0.480	0.344	0.332	0.389
	Textures	ViT-L/14 CLIP	0.543	0.441	0.446	0.462
	Textures	ConvNeXt V2-B	0.530	0.480	0.490	0.441
	Textures	ConvNeXt V2-L	0.551	0.468	0.480	0.440
-	iNaturalist	ResNet-50	0.703	0.700	0.716	0.684
	iNaturalist	ResNet-50 DINO	0.644	0.594	0.598	0.619
	iNaturalist	ResNet-34	0.745	0.721	0.726	0.728
	iNaturalist	ResNet-18	0.739	0.727	0.734	0.727
	iNaturalist	ViT-S/16	0.726	0.683	0.668	0.692
	iNaturalist	ViT-S/16 DINO	0.726	0.660	0.658	0.699
	iNaturalist	ViT-B/16	0.692	0.711	0.791	0.674
	iNaturalist	ViT-B/16 DINO	0.682	0.617	0.613	0.648
	iNaturalist	ViT-B/16 CLIP	0.698	0.655	0.683	0.634
	iNaturalist	ViT-B/14 DINOv2	0.519	0.429	0.426	0.455
	iNaturalist	ViT-G/14 DINOv2	0.448	0.355	0.351	0.384
	iNaturalist	ViT-L/14 CLIP	0.593	0.547	0.566	0.522
	iNaturalist	ConvNeXt V2-B	0.634	0.638	0.712	0.589
	iNaturalist	ConvNeXt V2-L	0.627	0.624	0.704	0.563

Table A.1: FPR@95 for OOD detection remains high with popular models. We record the FPR@95 for the MSP method for 14 models including ResNets, ViTs, and ConvNext V2 models on ImageNet-1K as ID, and Textures, iNaturalist, and ImageNet-OOD as OOD. FPR@95 records how many OOD examples are classified as ID (low uncertainty, false positive) at a threshold where 95% of ID examples are correctly classified (true positive). The average FPR@95 over all models and OOD datasets is 66.5%, thus well over half of OOD examples are classified as ID due to having low uncertainty, and other methods such as max logit, energy score, and entropy all have similar FPR@95s of over 60%.

- 1018
- 1019 A.3 HYBRID METHODS

We find hybrid methods like ViM and Hybrid-Add work well on far-OOD datasets like Textures, where we see noticeable improvement across many models in Figure A.6.

- 1023 A.4 EFFECT OF PRE-TRAINING
- Miller et al. (2021) showed that there is a strong linear relationship between ID accuracy and OOD generalization on OOD data with covariate shifts, suggesting it is sufficient to focus on improving



Figure A.4: For a ResNet-50 trained on ImageNet-1k, we see that the model has very high confidence for the OOD class 'Striped', highlighting the difference between label uncertainty and OOD uncertainty.



Figure A.5: We plot the distribution of OOD scores on ImageNet-1K ID and Describable Textures striped' class OOD data obtained from different OOD detection methods. We discover a systematic failure mode of all methods that utilize logits stemming from the model being overconfident about its predictions on the OOD data. Even though different OOD detection methods have different AUROC numbers, the score distribution plots reveal it is difficult to cleanly separate ID and OOD scores by picking a threshold. We use a ResNet-50 pretrained on ImageNet-1K and use a ViT-B/16 pretrained on ImageNet-1K.

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1067 ID accuracy for better robustness. Similarly, we explore the connection between the test accuracy 1068 and OOD detection performance. We use ImageNet-1K (IN-1K) as ID data and ImageNet-OOD 1069 (IN-OOD) (Yang et al., 2024) and Textures (Cimpoi et al., 2014) as OOD data. We evaluate 54 1070 models covering a wide range of architectures and pretraining methods. In Figure A.7 we plot 1071 AUROC of MSP against ImageNet test accuracy. Generally, ID accuracy and AUROC have close to a linear relationship for models with low- to medium-range performance on ImageNet. However, on 1072 ImageNet-OOD for models performing around or better than 75% ID accuracy, we observe higher 1073 variability in AUROC: for larger-scale highly performant models pre-training data impacts the OOD 1074 detection more significantly. 1075

When models are exposed to a diverse set of data during pre-training, they are likely to learn a wide
range of features, making it possible for them to differentiate between ID and semantically new
classes. There is an edge case when the OOD data is included in pre-training dataset: in Figure A.7
for the best performing ViT and ConvNext models pre-training includes ImageNet-21K, after which
they are fine-tuned on ImageNet-1K. Since ImageNet-OOD consists of images from ImageNet-21K



Figure A.7: The impact of the model architecture, pre-training data and objective on OOD detection
performance. AUROC of MSP on ImageNet vs ImageNet-OOD (left) and ImageNet vs Textures
(right) against ImageNet test accuracy. We observe that improving ImageNet accuracy generally
leads to better OOD detection.

which do not semantically overlap with ImageNet-1K classes, we observe a rapid jump in AUROC for
these models with negligible variability in ID accuracy. Pre-training on diverse data which includes
similar examples to OOD points softens the misspecification of the MSP approach and leads to strong
performance.

A.5 GENERATIVE MODELS



1123Figure A.8: Better generative model of ID data can lead to worse OOD detection. Left: RealNVP1124models achieving better likelihoods on ID CelebA images do not consistently achieve better AUROC1125for detecting CIFAR-10. Right: The Gaussian Mixture Model (GMM) on ResNet-50 features1126achieves best likelihoods with the empirical covariance of ImageNet features, but achieves best1127AUROCs for detecting ImageNet-OOD with the identity covariance ($\alpha = 1$). α represents the linear1128interpolation coefficient towards identity covariance.

Conceptual limitation of generative models for OOD detection. Estimating p(x) is different 1131 from estimating whether x is more likely to be drawn from some different distribution. Conceptually, 1132 for the latter, we would like to compute p(OOD|x), which by Bayes' rule, is p(x|OOD)/p(x) up to 1133 an x-independent constant. In general, knowing p(x) tells us nothing about the value of this ratio. p(x|OOD)/p(x) is also invariant to any coordinate transformation on x, whereas p(x) is not.

- 1134 Measuring typicality rather than density is an alternative method for OOD detection. Rather than 1135 asking whether a point has a high density, typicality asks whether a point belongs to a region with high 1136 probability mass. However, typicality has very similar pathologies compared to density. Consider a common motivating example for typicality: points drawn from a high dimensional Gaussian $\mathcal{N}(0, I)$ 1137 1138 in \mathbb{R}^d will have norms $||x|| = \sqrt{\sum_{i=1}^d x_i^2}$ concentrating around \sqrt{d} by the Law of Large Numbers 1139 (LLN). A point at the origin will be considered highly OOD based on typicality, since it has zero 1140 norm, yet it has the highest density and will thus be considered highly ID based on the density. 1141 But there is no reason why we should judge typicality based on norm rather than other quantities. 1142 Consider the average value of x over the dimensions, $\frac{1}{d} \sum_{i=1}^{d} x_i$, this quantity concentrates around 0 by LLN. Based on this quantity, a point at the origin looks perfectly typical, while a point on a sphere 1143 1144 of radius \sqrt{d} looks highly atypical. Therefore, exactly similar to the density, notions of typicality will 1145 tend to depend on a subjective choice of how to coarse-grain the input space based on quantities that 1146 are most relevant for distinguishing between ID and OOD data. Finally, measures of typicality can 1147 also depend on an arbitrary choice of coordinates. For example, (ϵ, N) -typical set (Nalisnick et al., 1148 2020) relies on the differential entropy, which is not invariant to coordinate transformations. 1149
- **OOD Detection Requires Coarse-Grained Representations.** In general, every test input we encounter will differ from the ID inputs we have previously seen. However, not all test inputs are considered OOD because we are only concerned with differences in certain essential aspects. When learning a generative model p(x) of the ID data, our goal is not necessarily to capture the distribution of x in its finest details. Instead, for the purpose of OOD detection, it is more appropriate to model the distribution over a coarse-grained representation h(x), which captures the attributes necessary for distinguishing OOD from ID data and ideally nothing more.
- Consider an ID dataset consisting of 1000 breeds of dogs and 10 breeds of cats. If our generative 1157 model captures the frequency of each individual breed, any dog input we observe will typically be 1158 considered 100 times more OOD-like than any cat input based on the likelihood of the generative 1159 model. However, if our goal is to detect other animal species and non-animal objects, the likelihood 1160 of this model is clearly not aligned with the objective of OOD detection. In this case, it would 1161 be more suitable to model only the frequency over the dog and cat categories, which serves as an 1162 appropriately coarse-grained representation of the individual breeds. Since the definition of OOD is 1163 ultimately user-defined, the correct coarse-grained representation depends on both the dataset and the 1164 intended definition of OOD, and it might be challenging to accurately specify even when a definition 1165 is known.
- In Figure A.8a, we show the test log-likelihoods (normalized by dimension) of RealNVP (Dinh et al., 2016) normalizing flow models of various sizes trained on CelebA images and their AUROCs for detecting CIFAR-10. While models with the lowest test likelihoods on ID data perform poorly for OOD detection, their OOD detection performance does not improve monotonically with their test likelihoods. In fact, the AUROC eventually decreases with improvements in likelihood.
- 1171 In Figure A.8b, we demonstrate the same phenomenon for a feature-space generative model. We 1172 construct a Gaussian Mixture Model (GMM) model of the features produced by a ResNet-50 pre-1173 trained on ImageNet-1K, the ID dataset. To optimize for the likelihood on ID data, we choose the 1174 cluster means to be the class-conditional means and use the empirical covariance of all features 1175 centered by their respective class means as the covariance of the clusters. This GMM model is 1176 precisely the generative model used by the Mahalanobis method (Lee et al., 2018). As we interpolate 1177 between the empirical covariance and a trivial identity covariance, the ID test likelihood of this GMM 1178 model decreases, yet the AUROC for detecting ImageNet-OOD improves monotonically.
- 1179
- The Impact of Inductive Biases. How a generative model assigns density to data unseen during training is highly dependent on their inductive biases. Despite being highly flexible density models, normalizing flows are known to be poor OOD detectors when trained as a generative model over raw images because their inductive biases encourage the model to focus on low-level pixel correlations rather than high-level semantic properties (Kirichenko et al., 2020; Nalisnick et al., 2018). Here, we demonstrate that the same failure mode applies to diffusion models, a distinct class of generative models achieving state-of-the-art image generation (Betker et al., 2023; Saharia et al., 2022).
- 1187 We use the Diffusion Transformer (DiT), a 256x256-resolution latent diffusion model trained on ImageNet-1K (Peebles and Xie, 2023). We score images based on the diffusion loss, a variance-



Figure A.9: Diffusion Models can fail catastrophically at OOD. (a): Using the diffusion loss, the Diffusion Transformer (DiT) (Peebles and Xie, 2023) fails catastrophically at detecting OOD inputs from the Describable Textures dataset. (b): the DiT model does decently at detecting OOD inputs from iNaturalist.



Figure A.10: Visualization of DiT Reconstruction Error. A DiT trained on ImageNet-1K often accurately reconstructs noised images from Describable Textures despite never having trained on them. Left: Reconstructions of noised Describable Textures images compared to middle: iNaturalist images and right: ImageNet-1K images.

reduced approximation of the variational lower bound (Kingma et al., 2021; Ho et al., 2020). In
 Figure A.9, we show the DiT fails catastrophically in detecting OOD data from Describable Textures
 but achieves decent performance in detecting OOD data from iNaturalist.

In Figure A.10, we qualitatively show that a 256x256 DiT trained on ImageNet-1K often accurately reconstructs noised images from Describable Textures despite never having trained on them. We add noise to the inputs corresponding to the diffusion timesteps at 49, 98, 147 out of 249, where higher timesteps are more noisy.

A.6 OUTLIER EXPOSURE

Training a model with outlier exposure is effective for improving OOD detection, and we see that performance is improved for most OOD problems with semantic shifts in Figure A.11



Figure A.11: Training a ResNet-18 with outlier exposure improves OOD detection for semantic shift datasets but hurts OOD generalization over covariate shifts.







1296 B IMPLEMENTATION DETAILS

1298 B.1 OUTLIER EXPOSURE EXPERIMENT.

On Figure 6, we compare OE model to the baseline ERM training in OOD detection (left panel) and OOD generalization (right panel). For semantic shift detection, we use CIFAR-100, Tiny ImageNet, MNIST, SVHN, Textures (Cimpoi et al., 2014), and Places365 (Zhou et al., 2014). For OOD generalization we evaluate on STL-10 (Coates et al., 2011), CINIC-10 (Darlow et al., 2018) and CIFAR-10-C (Hendrycks and Dietterich, 2019).

1304 We adapt OpenOOD codebase (Zhang et al., 2023; Yang et al., 2022) to train ResNet-18 with baseline 1305 ERM training and Outlier Exposure (Hendrycks et al., 2018) and evaluate models on OOD detection. We train models for 100 epochs with batch size 128 for ID data and batch size 256 for the outlier 1306 dataset, SGD with momentum and initial learning rate 0.1 and weight decay 5×10^{-4} , and we set 1307 the coefficient before the OE loss to alpha = 0.5 (overall, we use standard training hyper-parameters 1308 as in Zhang et al. (2023)). For Figure 6, we run both methods with 3 random seeds and report the 1309 average performance. To evaluate the model on STL-10, we only use the 9 classes which overlap 1310 with CIFAR-10 classes and drop the class "monkey" not present in CIFAR-10 (thus, the evaluation is 1311 marked as STL-9 in Figure 6). For CIFAR-C, we report the average accuracy across 15 corruptions 1312 (Gaussian Noise, Shot Noise, Impulse Noise, Defocus Blur, Glass Blur, Motion Blur, Zoom Blur, 1313 Snow, Frost, Fog, Brightness, Contrast, Elastic transform, Pixelate, JPEG Compression). 1314

1315 B.2 EVALUATING PRE-TRAINED MODELS.

We evaluate 54 models from the timm and torchvision libraries, including 9 different arheitecture types: ResNet, TinyNet, VGG, MobileNet, ConvNeXt, RegNetY, ReXNet, MLP-Mixer, and ViT; and 6 different pre-training data setups: training on IN-1K from scratch, pre-training on IN-21K and fine-tuning on ON-1K, pre-training on IN-12K (a subset of IN-21K) and fine-tuning on IN-1K, CLIP (Radford et al., 2021) pre-training on LAION and fine-tuning on IN-1K, CLIP pre-training on IN-21K and then IN-1K, and Instagram-1B pre-training and further IN-1K fine-tuning of SEER models (Goyal et al., 2021).

1323 B.3 SCALING EXPERIMENTS

We benchmark the following models to demonstrate impact of scale in Figure 4:

- 1326 1. ResNet-18 trained on ImageNet-1k
- 1327 2. ResNet-34 trained on ImageNet-1k
- 1328 3. ResNet-50 trained on ImageNet-1k
- 1330 4. ViT-S/16 trained on ImageNet-1k
- 5. ViT-B/16 trained on ImageNet-1k
- 1333 6. ViT-S/16 trained on ImageNet-1k with DINO
- 1334 7. ViT-B/16 trained on ImageNet-1k with DINO
- 1335 1336 8. ViT-B/16 trained with CLIP
- 1337 9. ViT-L/14 trained with CLIP
- 1338 1339 10. ViT-B/16 pretrained on CLIP, finetuned on ImageNet-1k
- 1340 11. ViT-B/14 trained on 142M images with DINOv2
- 1341 12. ViT-G/14 trained on 142M images with DINOv2

1343

- 1344
- 1345
- 1347

1348