

# DECIDER: Leveraging Foundation Model Priors for Improved Model Failure Detection and Explanation

**Rakshith Subramanyam\***

*Axio.ai*

RAKSHITH.SUBRAMANYAM@AXIO.AI

**Kowshik Thopalli\***

*Lawrence Livermore National Laboratory*

THOPALLI1@LLNL.GOV

**Vivek Narayanaswamy\***

*Lawrence Livermore National Laboratory*

NARAYANASWAM1@LLNL.GOV

**Jayaraman J. Thiagarajan**

*Lawrence Livermore National Laboratory*

JJTHIAGARAJAN@GMAIL.COM

**Editor:**

## Abstract

Reliably detecting when a deployed machine learning model is likely to fail on a given input is crucial for ensuring safe operation. In this work, we propose DECIDER (Debiasing Classifiers to Identify Errors Reliably), a novel approach that leverages priors from large language models (LLMs) and vision-language models (VLMs) to detect failures in image classification models. DECIDER utilizes LLMs to specify task-relevant core attributes and constructs a “debiased” version of the classifier by aligning its visual features to these core attributes using a VLM, and detects potential failure by measuring disagreement between the original and debiased models. In addition to proactively identifying samples on which the model would fail, DECIDER also provides human-interpretable explanations for failure through a novel attribute-ablation strategy. Through extensive experiments across diverse benchmarks spanning subpopulation shifts (spurious correlations, class imbalance) and covariate shifts (synthetic corruptions, domain shifts), DECIDER consistently achieves state-of-the-art failure detection performance, significantly outperforming baselines in terms of the overall Matthews correlation coefficient as well as failure and success recall. Our codes can be accessed at <https://github.com/kowshikthopalli/DECIDER/>

**Keywords:** Failure Detection, Vision-Language Models, Large-language Models

## 1 Introduction

A crucial step in ensuring the safety of deployed models is to proactively identify if a model is likely to fail for a given test input. This enables the implementation of appropriate correction mechanisms without impacting the model’s operation, or even deferring to human expertise for decision-making. While failures in vision models can be attributed to a variety of factors, the most significant cause is the violation of data distribution assumptions made during training (Jiang et al., 2019), which is the focus of this work. In general, data comprises both task-relevant *core attributes* and irrelevant *nuisance attributes*, and they are never

---

. \* equal contribution

explicitly annotated. Consequently, models can fail to generalize if (i) the training data contains spurious correlations (to nuisance attributes) that do not appear at test time, (ii) class-conditional distribution of nuisance attributes can arbitrarily change between train and test data (e.g., patient race imbalance in clinical datasets), or (iii) novel attributes emerge only at test time (e.g., style changes). Note that, when the class-conditional distributions of core attributes themselves change between train and test data, it leads to the more challenging scenario of *concept shifts*, and is not considered in this work. Nevertheless, detecting failures across all these scenarios is known to be challenging (Joshi et al., 2022; Yang et al., 2023; Geirhos et al., 2020), and hence there has been a surge in research interest (Hendrycks and Gimpel, 2017; Guillory et al., 2021; Gal and Ghahramani, 2016; Kirsch et al., 2021; Jain et al., 2023).

We begin by acknowledging that it is not only difficult, but also inefficient, to describe such nuisance attribute discrepancies solely using visual features. In this regard, we explore the utility of large language models (LLMs) and vision-language models (VLMs) in characterizing data attributes through a combination of visual and natural language descriptors. Subsequently, one can leverage these descriptors to design powerful failure detectors that systematically discern gaps in model generalization. Based on this idea, we develop **DECIDER** (Debiasing Classifiers to Identify Errors Reliably), a new approach for failure detection in vision models. At its core, **DECIDER** (i) utilizes LLMs (e.g., GPT-3 (Brown et al., 2020)) to specify task-relevant core attributes, (ii) uses a VLM (e.g., CLIP (Radford et al., 2021b)) to construct a “debiased” version of the task model by aligning its visual features to the core attributes, and (iii) detects failure by measuring disagreement between the original and debiased models for any given test input.

Additionally, **DECIDER** can be used to provide explanations for failure cases. This is done by employing an attribute-ablation strategy that adjusts the relative importance of core attributes such that the prediction probabilities of the debiased matches the original model. Our extensive empirical evaluation shows that our method achieves state-of-the-art performance in detecting failures across various datasets and test scenarios. In summary, our work provides early evidence for the utility of large-scale foundation models as priors for designing novel safety mechanisms.

## 2 Related Work

**Failure Detection.** Failure detection in classification involves identifying incorrect predictions made by the model (Hendrycks and Gimpel, 2017; Zhu et al., 2022; Qu et al., 2022). This problem ultimately boils down to identifying an appropriate metric or a *scoring function* that can delineate failed samples from successful ones. Early work involves using simple scores directly derived from the predictions of the model such as Maximum Softmax Probability (MSP) (Hendrycks and Gimpel, 2017), predictive entropy (Kirsch et al., 2021) and energy (Liu et al., 2020) to identify failed samples. More recent work focuses on scores that quantify failure by evaluating the local manifold smoothness (Ng et al., 2022) around a given sample and those that are based on agreement of a sample between different components of an ensemble (Jiang et al., 2022; Trivedi et al., 2023). However, such metrics can become unreliable to characterize failure as the model used to derive them can be potentially mis-calibrated and unreliable (Guo et al., 2017; Minderer et al., 2021). Failure

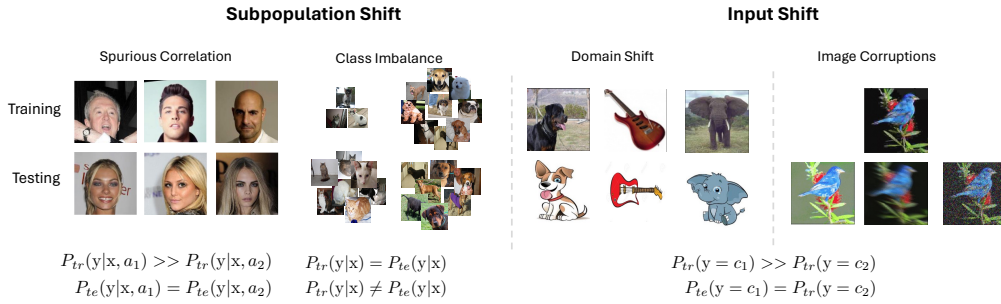


Figure 1: A visual illustration of the different failure scenarios we consider. These include scenarios when the model relies on spurious correlations present in the data i.e., when an attribute is spuriously correlated with the label (e.g., color of hair and gender). Another cause of failure is when the training data has class imbalance, leading to poorer generalization on images from the under-sampled class. Lastly, another important cause of failures are when the distribution of the test data is different from the training data. This can range from natural image corruptions to covariate shifts.

detection has also been studied under the lens of generalization gap estimation (Guillory et al., 2021; Narayanaswamy et al., 2022) where the goal is to predict the accuracy of the model on an unlabeled target distribution using distributional metrics derived from a number of calibration datasets.

**Failure Detection with Vision Language Foundation Models.** Visual-Language Models (VLMs) (Radford et al., 2021a; Li et al., 2022) are pre-trained on a large-corpora of image-text captions using a self-supervised objective. VLMs facilitate flexible adaptation to downstream tasks through zero-shot transfer or fine-tuning, demonstrating enhanced performance in zero-shot classification and OOD detection (Wei et al., 2023; Wortsman et al., 2022; Goyal et al., 2023; Ming et al., 2022; Wang et al., 2023; Michels et al., 2023; Esmailpour et al., 2022). Recently, VLMs have been used as a lens to understand the failure modes and weaknesses of any pre-trained model. For instance, the authors of (Jain et al., 2023) fit a post-hoc failure detector on the latent spaces of the VLM to estimate whether a sample has been correctly identified or not by the pre-trained classifier. The detector is then used to identify the directions of classifier failure modes. However, this approach requires a carefully tailored calibration set to fit the detector which is often unavailable in practice. On the other hand, the authors of (Deng et al., 2023) demonstrate that the latent space agreement between the pre-trained model and the VLM is a potential indicator for failure. In contrast, our paper aims to perform failure detection by first designing an improved classifier leveraging the VLM latent space and assessing the agreement between the classifier and its enhanced version while providing explanations for failure.

### 3 Background

**Preliminaries.** Let  $\mathcal{F}$  denote a multi-class classifier with parameters  $\theta$ , trained on a dataset  $\mathcal{D} = (x_i, y_i)_{i=1}^M$  comprising  $M$  samples. Here,  $x_i \in \mathcal{X}$ , is a 3 channel, input RGB image, and  $y_i \in \mathcal{Y}$  is the corresponding label, where  $\mathcal{Y}$  is the set of class labels i.e.,  $\mathcal{Y} = \{1, 2, \dots, C\}$ .

Here,  $C$  denotes the total number of distinct classes. The classifier  $\mathcal{F}$  operates on the input to produce the logits  $\mathcal{F}(x)$  corresponding to every class which is followed by a `softmax` operation to estimate output probabilities  $p(y = c|x)$  where  $c$  corresponds to the class index.

In this paper, we consider the problem of failure detection in classification models, where the source of failure arises due to the following scenarios (Fig. 1) - (i) Input level shifts where the training and test images share identical conditional output distributions i.e.,  $P_{tr}(y|x) = P_{te}(y|x)$  but different input marginals  $P_{tr}(x) \neq P_{te}(x)$ . Here, the test data can correspond to domain variations or image corruptions. (ii) Sub-population shifts (a) Spurious correlation where the labels are non-causally associated (Yang et al., 2023) with certain input characteristics or attributes in the training data over others leading to learning non-generalizable decision rules. For instance, let  $a_1$  and  $a_2$  correspond to two attributes of an image  $x$  and the training distribution is such that  $P_{tr}(y|x, a_1) \gg P_{tr}(y|x, a_2)$ . This model is susceptible to spurious correlations between the inputs and the targets and can fail during test time when  $P_{te}(y|x, a_1) = P_{te}(y|x, a_2)$ , (b) Class imbalance where the number of examples in a given class can be significantly greater than those present in another i.e.,  $P_{tr}(y = c_1) \gg P_{tr}(y = c_2)$ . This does not allow the classifier to optimally capture the image statistics and semantics of class  $c_2$  leading to sub-optimal generalization performance. **Failure Detector Design.** Failure detection is a binary classification problem of identifying whether an input sample has been correctly predicted or not by the model. We define our failure detector  $\mathcal{G}$  as follows,

$$\mathcal{G}(x; \theta, \tau) = \begin{cases} \text{failure,} & \text{if } s(x; \theta) < \tau, \\ \text{success,} & \text{if } s(x; \theta) \geq \tau. \end{cases} \quad (1)$$

Here,  $s(\cdot)$  is a scoring function derived from the classifier  $\mathcal{F}$  that assigns higher values for correctly identified samples and vice-versa and  $\tau$  is the user-controlled threshold for detection. Following standard practice from the generalization gap literature (Trivedi et al., 2023; Garg et al., 2022), we identify  $\tau$  such that  $\sum_i \mathbb{I}(s(x_i; \theta) \geq \tau)$  approximates the true accuracy of the held-out validation dataset.

**Contrastive Language-Image Pre-training (CLIP).** CLIP (Radford et al., 2021b) is a vision-language model trained on large corpus of image-text pairs with self-supervised learning. It aligns images with natural language descriptions in a shared embedding space, enabling zero-shot learning and fine-tuning for downstream tasks such as image captioning (Subramanyam et al., 2023) and visual question answering (Song et al., 2022; Yu et al., 2024; Guo et al., 2023; Schwenk et al., 2022). CLIP employs image ( $I(\cdot)$ ) and text ( $T(\cdot)$ ) encoders to generate embeddings ( $z_I$  and  $z_T$ ). For zero-shot inference, it computes the cosine similarity (*cos sim*) between image and text embeddings. This similarity yields class-specific logit scores for zero-shot classification, where the prediction probability  $p(y|x)$  is calculated using `softmax`.

## 4 Proposed Approach

### 4.1 Motivation

Typically, a classifier  $\mathcal{F}$  is trained on a dataset  $\mathcal{D}$  to learn the mapping between inputs and target labels. The datasets contain both task-relevant *core attributes* and irrelevant

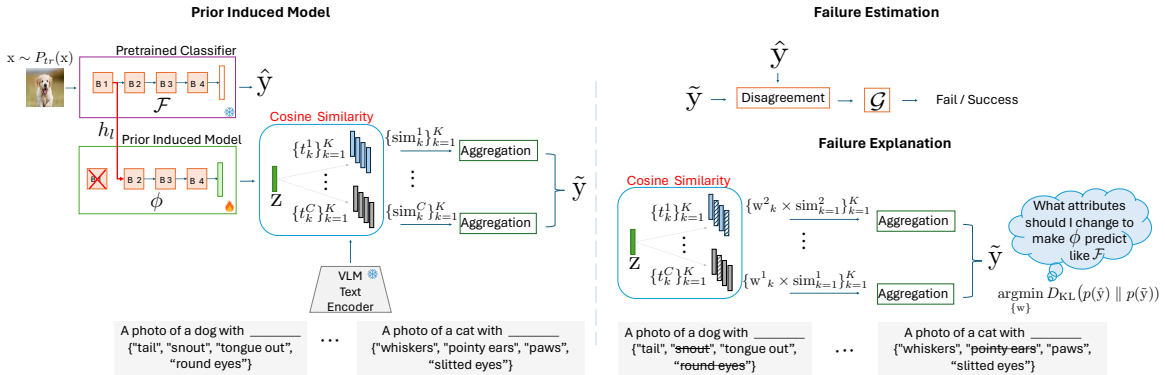


Figure 2: **DECIDER for failure detection.** (Left) DECIDER trains a Prior Induced Model (PIM)  $\phi$ , identical to the architecture of the pre-trained classifier  $\mathcal{F}$ , utilizing priors from a VLM model. (Top Right) The disagreement between the predictions of  $\phi$  and  $\mathcal{F}$  serves as an indicator for failure detection. (Bottom Right) By adjusting attribute level weights, DECIDER offers explanatory insights into failures.

*nuisance attributes*, which are not explicitly annotated. Consequently, the decision rules of the classifier could rely on nuisance attributes leading to poor generalization. For e.g., the model can fail to generalize if the training data contains spurious correlations with nuisance attributes that do not appear during testing. We underscore that this problem of reliance on nuisance attributes arises due to the difficulty in describing them solely using visual features.

To address this, we go beyond using only visual features and propose to leverage a combination of vision and language descriptors through the use of LLMs and VLMs and design failure detectors that discern the gap in model generalization. In this section we describe our novel strategy for failure detection which involves training a classifier referred to as the Prior Induced Model (PIM)  $\phi$  with the aid of LLMs and VLMs. We believe that the prior knowledge induced by VLMs will help PIM associate task-relevant core attributes. We first describe our paradigm that incorporates foundation models in classifier training. We then develop a prediction disagreement based strategy between PIM and the original classifier to conduct failure detection. Finally, we elucidate the capability of our approach in extracting failure explanations in order to support interpretability.

## 4.2 Incorporating Foundation Model Priors

A key challenge in traditional classification models is the direct mapping of images to coarse labels which encapsulate several attributes. For instance, in distinguishing between a dog and a cat, the label “dog” encompasses attributes like “wagging tail” and “snout”, while “cat” includes “whiskers” and “pointy ears”. Without explicit access to such detailed attribute information and due to potential biases in the training data, models are susceptible to rely on overly simplistic decision rules. In contrast, VLMs such as CLIP offer capabilities to encode both image and textual attribute descriptions into a unified latent space that is enriched to support meaningful image-text attribute associations.

To improve the effectiveness of classification model training, we hypothesize that aligning the model’s visual features with the textual descriptions of core attributes related to the class of interest in the VLM latent space can enhance training. This alignment is expected to equip the classifier with the ability to develop decision-making rules that are both more reliable and generalizable, while also reducing the influence of existing biases.

To achieve this, we introduce the PIM model  $\phi$ , which is guided by the LLM and VLM based priors (see Fig. 2 left). The architecture of PIM closely resembles that of its counterpart  $\mathcal{F}$ , with the notable distinction being that its final layer projects onto the VLM latent space. This projection supports the alignment with the textual descriptions of class-level attributes, thereby harnessing the linguistic capabilities of foundational models. PIM is specifically engineered to accept early-stage features from  $\mathcal{F}$ , denoted as  $h_l$ , which are then processed through PIM’s analogous layers to produce the image encoding  $z$  within the VLM latent space. For instance, when both  $\mathcal{F}$  and  $\phi$  are based on the ResNet architecture (He et al., 2016), the output from block 1 of  $\mathcal{F}$  serves as the input for block 2 in  $\phi$ .

It must be noted that the success of our approach relies upon the quality of the fine-grained text attributes extracted for every class. While there exists strategies (Merullo et al., 2022) that are capable of extracting image-level textual descriptions, they usually involve the text decoders in the loop which can be computationally expensive. Therefore, we resort to using Large Language Models (LLMs) to compute task-specific attribute descriptions offline.

### 4.3 Generating Task-specific Core-attribute Descriptions

LLMs (Touvron et al., 2023; Brown et al., 2020) have demonstrated their utility across a range of language tasks (Radford et al., 2019; Wei et al., 2022; Nakano et al., 2021; Pratt et al., 2023) and are particularly adept at contextual understanding, and generating coherent text even with descriptive prompting. To extract the class-specific attribute descriptions, we query GPT-3 (Brown et al., 2020) with the prompts “List visually descriptive attributes of  $\text{!CLASS!}$ .” This allows us to gather a set of  $K$  attributes  $\mathcal{A}^c = \{a_k^c\}_{k=1}^K$  for every class  $c$ .

### 4.4 Training PIM

**(i) Computing Cosine Similarities.** We first compute the cosine similarity scores between the image embedding  $z$  produced by PIM for a given image and the text embeddings associated with attribute  $k$  from each class  $c$ . It is given by,

$$\Omega_{\mathcal{A}^c} = \{\omega_k^c\}_{k=1}^K \text{ where } \omega_k^c = \cos \text{sim}(z, e_k^c) \quad (2)$$

Here, the text embeddings  $E_{\mathcal{A}^c} = \{e_k^c\}_{k=1}^K$  for each attribute of every class are obtained using the CLIP text encoder.

**(ii) Attribute Similarity Aggregation.** Subsequently, we aggregate these attribute similarity scores,  $\Omega_{\mathcal{A}^c}$ , for each class  $c$  to obtain coarse prediction logits corresponding to the class label  $y \in \mathcal{Y}$ . We investigate two aggregation strategies namely - (i) Class-level mean and (ii) Class-level max to consolidate these scores into final class predictions which are eventually normalized using `softmax`. These strategies enable a more refined and attribute-aware determination of classification outcomes.

**(iii) Optimization Objective.** The optimization is primarily guided by the cross-entropy loss which evaluates the discrepancy between the predicted probabilities from PIM and the ground truth label. In addition, we include consistency driven augmentations namely CutMix (Yun et al., 2019) and AugMix (Hendrycks et al., 2020) to improve its robustness. Additionally, we upweight the losses corresponding to the instances where (i) the biased classifier  $\mathcal{F}$  predicts accurately, but  $\phi$  does not and (ii) the biased classifier  $\mathcal{F}$  does not predict accurately, as well as  $\phi$  does not, within a training batch.

#### 4.5 DECIDER: Failure Estimation Using PIM

To assess the failure of the biased classifier  $\mathcal{F}$ , we compute the disagreement between PIM and  $\mathcal{F}$  based on the discrepancy between their predictions. This disagreement score is calculated as the cross-entropy between the sample-level probability distributions between the two models with PIM being the reference distribution given by  $s(\mathbf{x}) = -\sum_{c=1}^C p(y = c|\mathbf{x}) \cdot \log(q(y = c|\mathbf{x}))$  where  $p(\cdot)$  and  $q(\cdot)$  represent the predicted probabilities of  $\mathcal{F}$  and PIM, respectively.

#### 4.6 Extracting Explanations for Failure

Our failure explanation protocol is designed to elucidate the underlying reasons behind the discrepancies between predictions of  $\mathcal{F}$  and  $\phi$ . The primary objective is to identify the optimal subset of attributes necessary for aligning the PIM’s prediction probabilities with those of the task model. To achieve this, we implement an attribute ablation strategy where we iteratively adjust a group of weights corresponding to each attribute across all classes. Our iterative process begins by initially assigning uniform weights to every attribute for each class within a batch. These weights are then optimized by minimizing the Kullback-Leibler (KL) divergence between the probability distributions predicted by  $\mathcal{F}$  and those adjusted by PIM, accounting for the influence of the weighted attributes. As the algorithm converges, the weights will highlight those attributes that have significant impact on the predictions of  $\mathcal{F}$ , providing insights into the features considered by  $\mathcal{F}$  when making decisions. Fig. 2 right illustrates our failure explanation mechanism.

### 5 Empirical Evaluation

We conduct comprehensive evaluations of DECIDER using various classification benchmarks and assess performance under various failure scenarios with different architectures. We employ OpenAI’s CLIP ViT-B-32 model in all experiments (Radford et al., 2021a).

#### 5.1 Experimental Setup

**Datasets.** Our experiments are centered around datasets reflecting four common sources of model failure:

- **Input-Level Shifts:** CIFAR100-C (Hendrycks and Dietterich, 2019), comprising 19 types of corruptions at five severity levels over the CIFAR100 test images across 100 categories.

- **Spurious Correlations:** (1) Waterbirds (Yang et al., 2023) involves classifying images as ‘water bird’ or ‘land bird’. The training data offers biases tied to the background (water/land). (2) CelebA (Liu et al., 2015; Yang et al., 2023) involves classifying if individuals have blond hair or not, with labels spuriously correlated with gender.
- **Class Imbalance:** We modify the Kaggle Cats vs Dogs dataset (Cukierski, 2013), adjusting the distribution to create a training imbalance with 5,989 cat and 19,966 dog images for training, while maintaining balanced test data.
- **Distribution Shifts:** (1) PACS (Li et al., 2017) includes images from four domains (Photo, Art-painting, Cartoon, Sketch), to be classified into seven categories. As two large-scale benchmarks, we consider (2) DomainNet (Peng et al., 2019) which contains images from 345 categories from 6 domains (Real, Painting, Infograph, Quickdraw, Cartoon and Sketch) and (3) ImageNet-Sketch (Wang et al., 2019) benchmark which contains sketch images from 1000 ImageNet (Russakovsky et al., 2015) classes.

**Model Architectures.** We consider the ResNet-50 architecture for CelebA dataset and for all other datasets, we employ ResNet-18 trained on their respective datasets as the original classifier  $\mathcal{F}$ . In the supplementary, we study the performance of DECIDER on more architectures and we provide additional training details.

## 5.2 Baselines

We consider different baselines that use sample-level scores  $s$  for failure estimation:- (i) Maximum Softmax Probability (MSP) (Hendrycks and Gimpel, 2017) which is given by  $s(\mathbf{x}) = \max_j p(y = j|\mathbf{x})$ , (ii) Predictive Entropy (**Ent**) is essentially the entropy among the predictions of a sample and is given by  $s(\mathbf{x}) = -\sum_{j=1}^K p(y = j|\mathbf{x}) \cdot \log(p(y = j|\mathbf{x}))$ , (iii) **Energy** (Liu et al., 2020) score is defined by  $s(\mathbf{x}) = -T \cdot \log \sum_{j=1}^K \exp^{\mathcal{F}_\theta(\mathbf{x}_j)}$ . Following standard practice, we consider  $T = 1$  in all our experiments. (iv) Generalized Model Disagreement (**GDE**) (Jiang et al., 2022; Chen et al., 2021) - Let  $\mathcal{F}_{\theta_1}, \mathcal{F}_{\theta_2} \dots \mathcal{F}_{\theta_r}$  denote  $r$  models trained with different random seeds. Let  $\mathcal{F}_{\theta_1}$  denote the base classifier. Then the score is computed as  $s(\mathbf{x}) = \frac{1}{r} \sum_{i=1}^r \frac{1}{r-1} \sum_{j \neq i}^r \mathbb{I}(\mathcal{F}_{\theta_i} \neq \mathcal{F}_{\theta_j})$ . We set  $r$  to 5.

It must be noted that we utilize negative versions of entropy and energy to reflect the fact the samples that are correctly predicted are associated with higher scores.

## 5.3 Metrics

We consider the following metrics to evaluate failure detection performance: (i) Failure Recall (FR) which corresponds to the fraction of samples that have been correctly identified as failure, (ii) Success Recall (SR) corresponds to the fraction of samples that have been correctly predicted as successful. The trade-off between the two metrics is indicative of how aggressive or conservative the failure detector is. (iii) Matthew’s Correlation Coefficient (MCC) holistically assesses the quality of the binary classification task of failure detection and provides a balanced measure when the class sizes are different. It takes into account both true and false positives and negatives respectively while assessing performance.



Dataset	Method	FR	SR	MCC
CIFAR100	MSP	0.6835	0.809	0.4943
	Energy	0.6776	0.7965	0.4747
	Ent	0.6894	0.8105	0.514
	<b>DECIDER</b>			
	+ mean	0.7949	0.7436	0.5267
	+ max	0.7933	0.7474	<b>0.5292</b>
	CIFAR100-C	MSP	0.7448	0.6345
Energy		0.8145	0.5442	0.3577
Ent		0.7761	0.616	0.3766
<b>DECIDER</b>				
+ mean		0.8507	0.5393	0.4007
+ max		0.8448	0.5506	<b>0.4015</b>
Waterbirds		MSP	0.3166	0.8891
	Energy	0.4803	0.8047	0.2814
	Ent	0.4878	0.8022	0.2827
	<b>DECIDER</b>			
	+ mean	0.5303	0.8310	0.3580
	+ max	0.6063	0.8580	<b>0.4598</b>
	CelebA	MSP	0.4058	0.9653
Energy		0.4292	0.9616	0.3677
Ent		0.4214	0.9631	0.3675
<b>DECIDER</b>				
+ mean		0.5443	0.9701	<b>0.4928</b>
+ max		0.4390	0.9621	0.3738
Cats and Dogs		MSP	0.4076	0.9235
	Energy	0.4303	0.9196	0.3428
	Ent	0.4233	0.9212	0.3402
	<b>DECIDER</b>			
	+ mean	0.5993	0.9468	0.544
	+ max	0.5783	0.9554	<b>0.5532</b>

**(a)**
**(b)**

Figure 3: Results on failure detection across different benchmarks - (a) CIFAR100, and image corruptions on CIFAR-100-C, and (b) subpopulation shifts from spurious correlations on Waterbirds, CelebA datasets, and class imbalance on Cats vs Dogs. DECIDER consistently outperforms baselines in terms of the overall Matthew’s Correlation Coefficient (MCC) as well as achieving higher failure and success recalls.

## 5.4 Findings

**Input Shifts.** Fig. 3(a) showcases the results on the CIFAR100 and CIFAR100-C datasets. On the clean CIFAR100, DECIDER outperforms the baselines with a superior MCC of 0.5292 for the max variant (versus 0.514 for the best baseline), attributed to higher failure recall (0.7933) and success recall (0.7474). On the more challenging CIFAR100-C (severity level 4), DECIDER further highlights its efficacy by achieving an MCC of 0.4015 with max aggregation, exceeding the top baseline (entropy) which has an MCC of 0.3766. This is due to a balanced trade-off between failure recall (0.8448) and success recall (0.5506), distinguishing DECIDER from other baselines that fail to maintain such balance. These findings clearly demonstrate DECIDER as robust in detecting classifier failures amid input-level shifts, surpassing other baselines in performance metrics.

**Subpopulation Shifts.** Our comprehensive evaluation addresses datasets affected by various subpopulation shifts. The summarized results in Fig. 3(b) underline the effectiveness of DECIDER in navigating these challenges:

Waterbirds: DECIDER achieves a high failure recall of 0.6063, outperforming the best baseline (entropy) which has a recall of 0.4878. Importantly, DECIDER maintains a high success recall

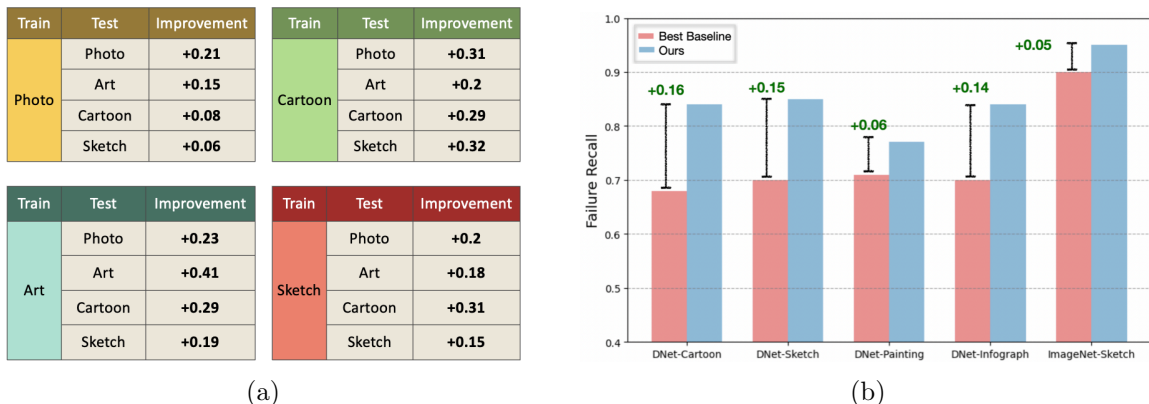


Figure 4: DECIDER produces the best performance on covariate shifts.. (left) Comparison of DECIDER against the best baseline in terms of the difference in MCC on the PACS dataset involving covariate shifts across 4 different visual domains. (Right) Improvement in failure recall performance of the best performing baseline and DECIDER on large-scale covariate shift benchmarks- DomainNet (DNet) and ImageNet-Sketch. The classifiers and PIMs are trained on DomainNet Real and Imagenet train sets respectively and evaluated on the different distribution shift datasets.

(0.858) with minimal compromise compared to MSP (0.8891). The outcome is a leading MCC of 0.4598, attesting to DECIDER’s balanced detection ability in environments with misleading background cues.

CelebA: With mean aggregation, DECIDER delivers the highest MCC of 0.4928, combining a failure recall of 0.5443 with a success recall of 0.9701, showcasing its strength in addressing gender and hair color spurious correlations.

Cats vs Dogs: Exhibiting strong performance in class imbalance, DECIDER (max aggregation) achieves an MCC of 0.5532, significantly surpassing the top baseline (energy) with an MCC of 0.3428, underlining its efficacy in balanced success and failure recall. DECIDER not only demonstrates high failure detection capability but also ensures high success recall rates above 0.94, highlighting its proficiency in class-imbalanced settings.

**Covariate Shifts.** In this section, we evaluate the performance of DECIDER in the challenging setting of identifying failure due to covariate shifts. We first consider the PACS dataset which contains 4 different domains. We train PIM and derive individual thresholds for each of the four domains and evaluate its performance across all domains. While we present detailed results for baselines and metrics in the supplementary, in Fig. 4(a), we report the gain in MCC scores between the best performing baseline and DECIDER. It can be seen that DECIDER outperforms the baselines by a large margin across all the domains. To further validate the effectiveness of DECIDER, we conducted experiments on large-scale covariate shift benchmarks, including DomainNet and ImageNet. In the DomainNet case, we trained the classifier and PIM on images from the real domain and evaluated their performance on four different target domains: Cartoon, Sketch, Painting, and Infograph. For ImageNet, we trained on the ImageNet training dataset and assessed the performance on the challenging

ImageNet-Sketch benchmark. Fig. 4(b) presents the failure recall performance of the best-performing baseline and DECIDER, clearly demonstrating the superiority of our approach even when applied to large-scale datasets.

In summary, these results highlight the importance of leveraging language priors together with priors from the VLM to construct debiased models that reliably help detect failures across different scenarios.

## 6 Failure Explanation

Having empirically demonstrated the superior failure detection capabilities of DECIDER, we now turn our attention to the task of explaining the reasons behind failures. To that end, we consider the max variant of DECIDER and adjust the influence of individual attributes to ensure that the prediction probabilities generated by DECIDER closely mirror those of the original model as explained in Section 4. This manipulation offers evidence of what attributes the task model uses. For e.g., on the top left of Fig. 5, the task is to correctly identify the hair color. Here, the classifier  $\mathcal{F}$  incorrectly classifies the image, while PIM accurately identifies the same. We observe that our optimization process reduces the influence of core attributes such as "Browning Tresses" and "Red Highlights" on PIM's predictions. This manipulation serves as evidence that the biased classifier  $\mathcal{F}$  may not have considered these crucial attributes in its decision-making process. Similarly, in the example shown in Fig. 5,  $\mathcal{F}$  misclassifies a Cat as a Dog (top right) and the proposed optimization shows that the classifier has not focused enough on the important core attributes such as "Thin Whiskers" thus making the erroneous classification.

## 7 Analyses

**Biases or insufficiency of GPT-3 attributes.** The success of DECIDER relies on the quality of the attributes generated by the LLM. To study the impact on failure detection on the quality of text attributes, we consider two practical scenarios: (i) GPT-3 generates irrelevant attributes: In this case, the PIM model has the risk of learning noisy decision rules that the even the classifier might not have; (ii) GPT provides insufficient attributes: With only partial attributes, PIM's predictive performance can be limited. To comprehensively evaluate the impact of both scenarios, we employ the following protocol on the Waterbirds dataset. For scenario (i), we add 5 randomly sampled core attributes from the other class to the attribute set of each class. For case (ii), we remove 5 randomly selected attributes from the attribute set of each class. We train PIM under both these scenarios. As the results in Table 1 show, although there is a noticeable drop in performance due to the severe attribute corruptions, DECIDER still outperforms the best baseline (Ent) method. This demonstrates the robustness of DECIDER to imperfect attribute sets.

**Impact of Layer Selection of  $\mathcal{F}$  on  $\phi$ .** In this study, we explore how the performance of the PIM model  $\phi$  is influenced by the specific layer in  $\mathcal{F}$  from which we extract features. This experiment uses the ResNet-18 architecture, with models trained on the CIFAR100 and Waterbirds datasets. From the results presented in the table in Fig. 6, using features from the early layers (layer 1 and layer 2) of ResNet-18 yields the highest MCC (Matthews Correlation Coefficient) scores. In contrast, leveraging features from the later layers leads

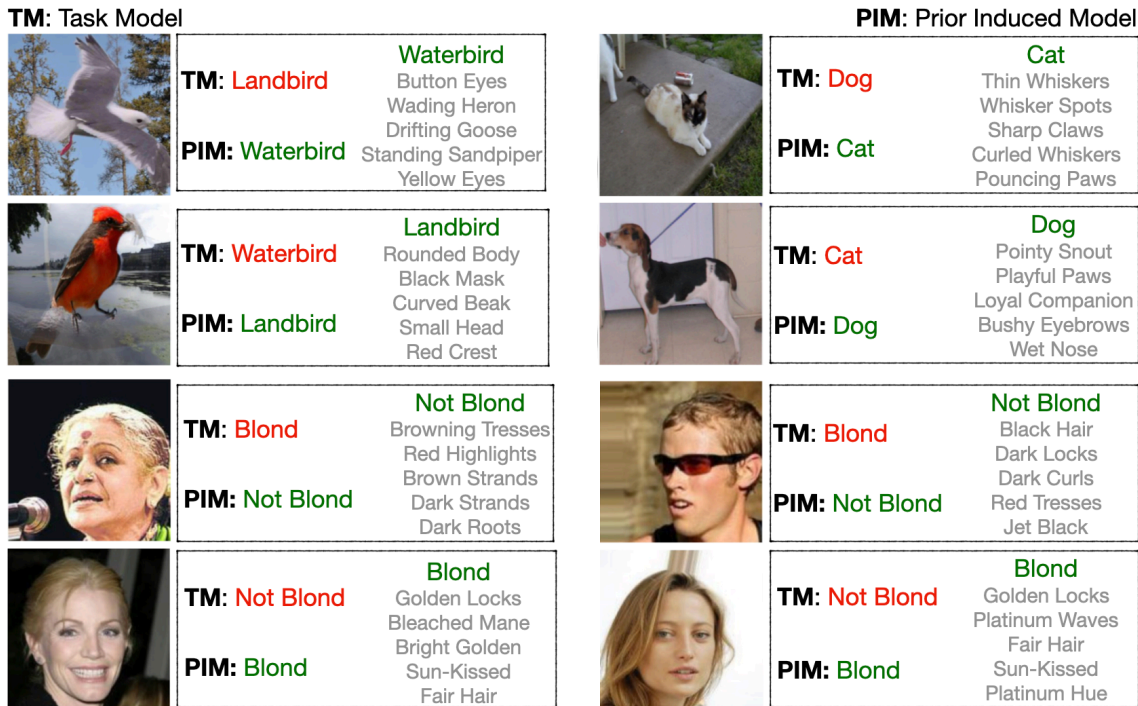


Figure 5: **Failure Explanations.** We explain the failures of the biased classifier  $\mathcal{F}$ , by manipulating the influence of individual attributes in PIM, such that the prediction probabilities of PIM match that of  $\mathcal{F}$ . The knowledge of the attributes whose influence was needed to be reduced provides an indication that  $\mathcal{F}$  has not focused on those attributes to make its decisions. We show qualitative examples on Water birds in top left, Cats vs dogs in top right and from CelebA dataset in bottom.

Table 1: **Impact of attribute quality** – (i) *irrelevant*: add 5 nuisance attributes; (ii) *insufficient*: remove 5 core attributes. Although there is a drop in performance under attribute corruptions, DECIDER still outperforms existing baselines.

Metric	Baseline (Ent)	DECIDER	DECIDER (irrelevant)	DECIDER (insufficient)
Failure Recall	0.48	<b>0.60</b>	0.54	0.49
Success Recall	0.80	<b>0.85</b>	0.81	0.83
MCC	0.28	<b>0.45</b>	0.34	0.33

to a noticeable decline in performance. This observation suggests that the initial layers of the network are less prone to carrying biases than the later ones, supporting the findings from previous research (Lee et al., 2022).

**Model Ensembles for Disagreement Analysis.** It has been shown that the prediction disagreement between different constituent members of a model ensemble can serve as an indicator of failure (Jiang et al., 2022; Trivedi et al., 2023). In this experiment, we compare

the failure estimation performance obtained through the disagreement between PIM and  $\mathcal{F}$  to the performance obtained by the disagreement between an ensemble (GDE). To that end, we trained five different classifiers with different initial seeds on three different datasets: Waterbirds, CelebA, and Cat vs Dogs. Figure 6, evidences the superiority of the proposed approaches compared to GDE.

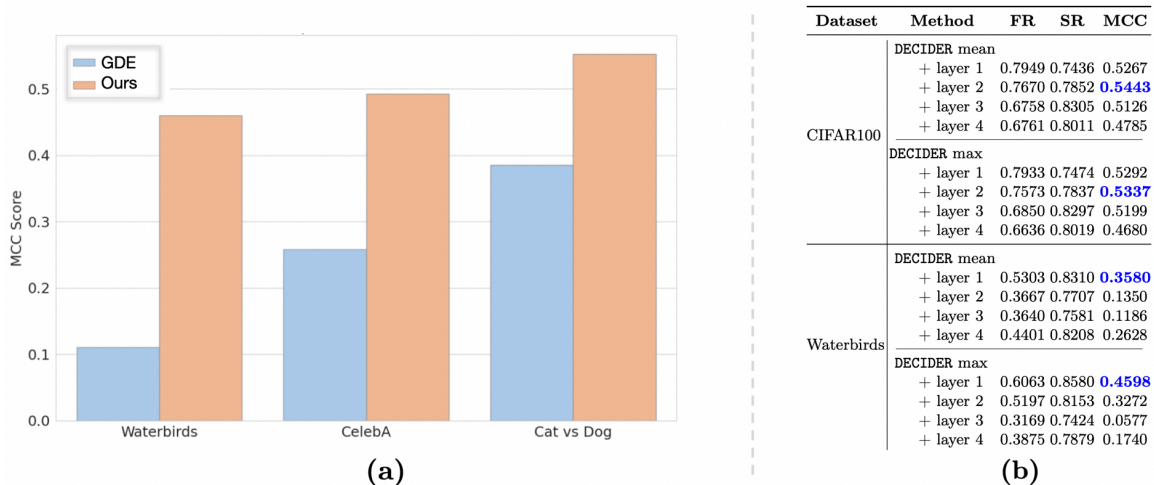


Figure 6: (a) Comparison of DECIDER against the failure detection performance obtained through disagreement between predictions from an ensemble of multiple instances of  $\mathcal{F}$  on Waterbirds, CelebA and Cats vs Dogs datasets respectively. (b) Ablation study analyzing the impact of using features from different layers of the base model  $\mathcal{F}$  as input to the Prior Induced Model (PIM)  $\phi$  on CIFAR-100 and Waterbirds datasets.

**Impact of PIM accuracy on failure detection.** Since we attempt to train a debiased classifier, in this section, we study the impact of its accuracy on failure detection. Table 2 in the appendix reveals that, despite the occasional slight decrease in the predictive performance of the debiased model PIM, the core-nuisance attribute disambiguation, which is crucial for failure detection, is not compromised. Consequently, DECIDER consistently achieves superior failure recall compared to the baselines.

**Replacing PIM with CLIP classifiers.** Given that we propose leveraging the priors from CLIP to obtain a debiased version of the classifier, it is natural to consider utilizing CLIP’s zero-shot classifier directly as PIM. Table 3 in appendix demonstrates that such an approach yields poor failure detection performance when CLIP’s zero-shot classifier is employed as PIM. This is because the visual features and their correlations to the core attributes of CLIP can differ significantly from the task model, thus rendering the model disagreement based failure detection highly ineffective.

## 8 Conclusion

In this work, we introduced DECIDER, a novel approach that leverages LLMs and vision-language foundation models to detect failures in pre-trained image classification models. Our key insight was to train an improved version of the pre-trained classifier, PIM, that

learns robust associations between visual features and class-level attributes by projecting into the shared embedding space of a VLMs such as CLIP. By analyzing the disagreement between PIM’s predictions and the original biased model, **DECIDER** can reliably identify potential failures while offering human-interpretable explanations. Extensive experiments across multiple benchmarks evidences the consistent superiority of **DECIDER** over baselines, achieving substantially higher overall scores and better trade-offs between failure and success recalls. Our work highlights the promise of integrating vision-language priors into model failure analysis pipelines to facilitate more reliable and trustworthy deployment of vision models in safety-critical applications. Extending **DECIDER** to other vision-language models and exploring its application to other failure modes such as adversarial attacks constitute our future work.

## Acknowledgments and Disclosure of Funding

This work was performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under Contract No. DE-AC52-07NA27344, Lawrence Livermore National Security, LLC. Supported by LDRD project 24-FS-002. LLNL-CONF-862086.

## References

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Jiefeng Chen, Frederick Liu, Besim Avci, Xi Wu, Yingyu Liang, and Somesh Jha. Detecting errors and estimating accuracy on unlabeled data with self-training ensembles. *Advances in Neural Information Processing Systems*, 34:14980–14992, 2021.
- Will Cukierski. Dogs vs. cats, 2013. URL <https://kaggle.com/competitions/dogs-vs-cats>.
- Ailin Deng, Miao Xiong, and Bryan Hooi. Great models think alike: Improving model reliability via inter-model latent agreement. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 7675–7693. PMLR, 23–29 Jul 2023.
- Sepideh Esmailpour, Bing Liu, Eric Robertson, and Lei Shu. Zero-shot out-of-distribution detection based on the pre-trained model clip. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 6568–6576, 2022.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.

- Saurabh Garg, Sivaraman Balakrishnan, Zachary Chase Lipton, Behnam Neyshabur, and Hanie Sedghi. Leveraging unlabeled data to predict out-of-distribution performance. In *International Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?id=o\\_HsiMPYh\\_x](https://openreview.net/forum?id=o_HsiMPYh_x).
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- Sachin Goyal, Ananya Kumar, Sankalp Garg, Zico Kolter, and Aditi Raghunathan. Finetune like you pretrain: Improved finetuning of zero-shot vision models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19338–19347, 2023.
- Devin Guillory, Vaishaal Shankar, Sayna Ebrahimi, Trevor Darrell, and Ludwig Schmidt. Predicting with confidence on unseen distributions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1134–1144, 2021.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.
- Jiaxian Guo, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Boyang Li, Dacheng Tao, and Steven Hoi. From images to textual prompts: Zero-shot visual question answering with frozen large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10867–10877, 2023.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=HJz6tiCqYm>.
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *Proceedings of International Conference on Learning Representations*, 2017.
- Dan Hendrycks, Norman Mu, Ekin D. Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. AugMix: A simple data processing method to improve robustness and uncertainty. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- Saachi Jain, Hannah Lawrence, Ankur Moitra, and Aleksander Madry. Distilling model failures as directions in latent space. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=99RpBVpLiX>.



- Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019.
- Yiding Jiang, Vaishnavh Nagarajan, Christina Baek, and J Zico Kolter. Assessing generalization of SGD via disagreement. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=Wv0GCEAQhx1>.
- Nitish Joshi, Xiang Pan, and He He. Are all spurious features in natural language alike? an analysis through a causal lens. *arXiv preprint arXiv:2210.14011*, 2022.
- Andreas Kirsch, Jishnu Mukhoti, Joost van Amersfoort, Philip H. S. Torr, and Yarin Gal. On pitfalls in ood detection: Entropy considered harmful, 2021. Uncertainty & Robustness in Deep Learning Workshop, ICML.
- Yoonho Lee, Annie S Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, and Chelsea Finn. Surgical fine-tuning improves adaptation to distribution shifts. *arXiv preprint arXiv:2210.11466*, 2022.
- Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pages 5542–5550, 2017.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022.
- Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. *Advances in neural information processing systems*, 33:21464–21475, 2020.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- Jack Merullo, Louis Castricato, Carsten Eickhoff, and Ellie Pavlick. Linearly mapping from image to text space. *arXiv preprint arXiv:2209.15162*, 2022.
- Felix Michels, Nikolas Adaloglou, Tim Kaiser, and Markus Kollmann. Contrastive language-image pretrained (clip) models are powerful out-of-distribution detectors. *arXiv preprint arXiv:2303.05828*, 2023.
- Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil Houlsby, Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks. *Advances in Neural Information Processing Systems*, 34:15682–15694, 2021.
- Yifei Ming, Ziyang Cai, Jiuxiang Gu, Yiyou Sun, Wei Li, and Yixuan Li. Delving into out-of-distribution detection with vision-language representations. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=KnCS9390Va>.



- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Vivek Narayanaswamy, Rushil Anirudh, Irene Kim, Yamen Mubarka, Andreas Spanias, and Jayaraman J. Thiagarajan. Predicting the generalization gap in deep models using anchoring. In *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4393–4397, 2022.
- Nathan Ng, Kyunghyun Cho, Neha Hulkund, and Marzyeh Ghassemi. Predicting out-of-domain generalization with local manifold smoothness. *arXiv preprint arXiv:2207.02093*, 2022.
- Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1406–1415, 2019.
- Sarah Pratt, Ian Covert, Rosanne Liu, and Ali Farhadi. What does a platypus look like? generating customized prompts for zero-shot image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15691–15701, 2023.
- Haoxuan Qu, Yanchao Li, Lin Geng Foo, Jason Kuen, Jiuxiang Gu, and Jun Liu. Improving the reliability for confidence estimation. In *European Conference on Computer Vision*, pages 391–408. Springer, 2022.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021a.
- 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021b.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge. In *European Conference on Computer Vision*, pages 146–162. Springer, 2022.

- Haoyu Song, Li Dong, Wei-Nan Zhang, Ting Liu, and Furu Wei. Clip models are few-shot learners: Empirical studies on vqa and visual entailment. *arXiv preprint arXiv:2203.07190*, 2022.
- Rakshith Subramanyam, TS Jayram, Rushil Anirudh, and Jayaraman J Thiagarajan. Crepe: Learnable prompting with clip improves visual relationship prediction. *arXiv preprint arXiv:2307.04838*, 2023.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Puja Trivedi, Danai Koutra, and Jayaraman J Thiagarajan. A closer look at scoring functions and generalization prediction. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. In *Advances in Neural Information Processing Systems*, pages 10506–10518, 2019.
- Hualiang Wang, Yi Li, Huifeng Yao, and Xiaomeng Li. Clipn for zero-shot ood detection: Teaching clip to say no. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1802–1812, 2023.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=gEZrGCozdqR>.
- Yixuan Wei, Han Hu, Zhenda Xie, Ze Liu, Zheng Zhang, Yue Cao, Jianmin Bao, Dong Chen, and Baining Guo. Improving clip fine-tuning performance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5439–5449, 2023.
- Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7959–7971, 2022.
- Yuzhe Yang, Haoran Zhang, Dina Katabi, and Marzyeh Ghassemi. Change is hard: A closer look at subpopulation shift. In *International Conference on Machine Learning*, 2023.
- Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. Self-chained image-language model for video localization and question answering. *Advances in Neural Information Processing Systems*, 36, 2024.
- Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *International Conference on Computer Vision (ICCV)*, 2019.

Fei Zhu, Zhen Cheng, Xu-Yao Zhang, and Cheng-Lin Liu. Rethinking confidence calibration for failure prediction. In *European Conference on Computer Vision*, pages 518–536. Springer, 2022.

## Appendix A. Additional Analysis of DECIDER

### A.1 Impact of PIM’s predictive performance

We note that, the quality of attributes has a direct impact on PIM. That said, even with

Table 2: Regardless of the predictive performance w.r.t. the task model, PIM is very useful for failure detection.

Source	Target	Accuracy (%)		Failure Recall	
		Baseline	PIM	Baseline	PIM
Waterbirds	Waterbirds	75.06	<b>82.11</b>	0.49	<b>0.61</b>
CIFAR-100	CIFAR-100C	<b>31.84</b>	31.25	0.81	<b>0.85</b>
DomainNet-R	DomainNet-C	<b>33.14</b>	30.12	0.68	<b>0.84</b>
DomainNet-R	DomainNet-S	<b>20.02</b>	<b>20.11</b>	0.70	<b>0.85</b>
DomainNet-R	DomainNet-P	31.25	<b>31.72</b>	0.71	<b>0.77</b>
DomainNet-R	DomainNet-I	11.16	<b>12.63</b>	0.70	<b>0.84</b>
ImageNet	ImageNet-Sketch	<b>23.57</b>	22.49	0.90	<b>0.95</b>

slightly lower performance, the core-nuisance attribute disambiguation, which is the most critical for failure detection, is not compromised. To demonstrate this, we show in Table 2 the failure recall performance on Waterbirds, CIFAR-100, ImageNet, and DomainNet (Real to Clipart, Sketch, Painting, Infograph domains) datasets. Regardless of its predictive performance, the failure recall of DECIDER is consistently higher than the best performing baseline.

### A.2 Replacing PIM with CLIP classifiers

Through PIM, we create a variant of the task model that disambiguates core attributes (identified using LLM) from nuisance attributes. This means that if the task model and PIM do not agree on a prediction, it is likely a failure. However, with a zero-shot CLIP classifier, the visual features and their correlations to the core attributes can be drastically different from the task model. This renders the model disagreement based failure detection highly ineffective. To show this, we tried two versions of the CLIP classifier: one using text prompts for each class label (CLIP-cl) and another using core attributes (CLIP-att) like DECIDER. From Table 3, we see that PIM performs much better than both versions.

Table 3: Replacing PIM with CLIP classifiers

	CLIP-cls	CLIP-att	PIM
DECIDER MCC	-0.01	0.02	<b>0.46</b>

### A.3 Ablating synthetic augmentations for PIM training

We conducted an ablation study on CIFAR100 to demonstrate the value of augmentations used during PIM training. As Table 4 shows, there is a significant drop in failure detection

Table 4: **Impact of Augmentations:** We compare the failure detection performance in the presence and absence of augmentations for PIM training.

	No Aug.	Augmix (0.2) + Cutmix (0.2)
DECIDER MCC	0.27	<b>0.53</b>

performance without augmentations.

## Appendix B. Algorithm Listing for PIM

---

### Algorithm 1 Training Procedure for Prior Induced Model $\phi$

---

**Input:** Training data  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^M$ , attribute set  $a_k$  for each class  $c \in \mathcal{Y}$  extracted from LLM, VLM (CLIP) text encoder  $T(\cdot)$ , classifier  $\mathcal{F}$ , cross-entropy loss  $\mathcal{L}(\cdot)$  parameters  $\phi$  initialized with ImageNetv1 weights.

**Output:** Optimized parameters  $\phi$ .

- 1: **for** each epoch **do**
  - 2:     **for** each batch  $\{(x_i, y_i)\}$  in  $\mathcal{D}$  **do**
  - 3:         Apply augmentation (Cutmix or Augmix) across batch with probability  $p$
  - 4:         Compute features  $\mathbf{h}_l$  for layer  $l$  from  $\mathcal{F}$ .
  - 5:         Use PIM to map  $\mathbf{h}_l$  on to the VLM latent space to obtain  $\phi(\mathbf{h}_l)$
  - 6:         Initialize sample-level loss weights to uniform
  - 7:         **for** each class  $c$  **do**
  - 8:             Compute cosine similarity between  $\phi(\mathbf{h}_l)$  and the CLIP text embeddings  $T(\cdot)$  of attributes from class  $c$ .
  - 9:             Aggregate similarities to derive class-level logits  $\tilde{y}$ .
  - 10:         **end for**
  - 11:         Compute sample-level loss weights based on the discrepancy in predictions between  $\mathcal{F}$  and  $\phi$ .
  - 12:         Update  $\phi$  using the objective -  $\min_{\phi} \mathcal{L}(y, \tilde{y})$
  - 13:         **end for**
  - 14:     **end for**
  - 15: **return** Optimized  $\phi$ .
-

## Appendix C. Training Details

### C.1 Classifier Training

Table 5 provides the hyper-parameter and optimization settings for every dataset employed for training the classifier  $\mathcal{F}$ . We use a multi-step LR decay scheduler, which reduces the learning rate by a factor of 0.2.

Table 5: Hyper-parameter and optimization settings for training classifier  $\mathcal{F}$  for different datasets.

Dataset	Epochs	Initial Learning Rate	LR Decay epochs	Momentum	Optimizer
CIFAR100	200	0.1	60, 120, 160	0.2	SGD
Waterbirds	100	0.001	30, 60	0.9	SGD
CelebA	20	0.1	-	0.9	SGD
Cats & Dogs	100	0.01	30, 60	0.9	SGD
PACS	200	0.01	60, 120, 160	0.9	SGD
Domainnet	30	0.001	-	0.9	Adam

### C.2 PIM (Prior Induced Model) Training Details

We adopt the following protocol to train PIM for all datasets. We train PIM for 200 epochs, starting with an initial learning rate of 0.1, and implement multi-step decay at epochs 60, 120, and 160, where we reduce the learning rate by factor of 0.1. We utilize the AdamW optimizer for optimization. Additionally, we apply both CutMix and AugMix transformations to the entire batch with probabilities of 0.2 each. Moreover, we carefully weight our loss function during training. The loss weights are increased by a factor of 2.0 for samples where the classifier  $\mathcal{F}$  succeeds but PIM fails, and by 1.5 for cases where both the classifier and PIM fail.

## Appendix D. Prompts Used to Query LLM (GPT3) for Attribute Generation

- **Waterbirds:** *"List 100 distinct two-word phrases that uniquely describe the visual characteristics (like type of feet, beak, wings, plumage, feathers, feather texture, body shape, body type etc) of {class\_name}. Make sure the phrases are not long descriptions."*
- **CIFAR100:** *"List 50 distinct two-word phrases that uniquely describe the visual characteristics (like shape, color, texture) of {class\_name}. Make sure the phrases are not long descriptions."*

- **PACS:** *"List 30 distinct two-word phrases that uniquely describe the visual characteristics of {class\_name}. Do not describe their colors. Make sure the phrases are not long descriptions."*
- **CelebA:** *"List 25 distinct two-word phrases that uniquely describe the visual characteristics of {class\_name} hair person. Make sure the phrases are not long descriptions."*
- **Cats and Dogs:** *"List 50 distinct two-word phrases that uniquely describe the visual characteristics of {class\_name}. Make sure the phrases are not long descriptions."*
- **DomainNet:** *"List 100 distinct two-word phrases that uniquely describe the visual characteristics of {class\_name}. Make sure the phrases are not long descriptions."*
- **ImageNet:** *"List 100 distinct two-word phrases that uniquely describe the visual characteristics of {class\_name}. Make sure the phrases are not long descriptions."*

## Appendix E. Additional Results

**Experiment with ViT-B-16:** We extend our study to incorporate the ViT architecture, specifically using the ViT-B-16 model, for the Waterbirds datasets. We provide these results in Table 6. For ViT-B-16, we explore two PIM variations: one using features from the first layer and another from the ninth layer of the ViT-B-16 classifier model. From Table 6 it is evident that DECIDER continues to outperform as a more reliable failure estimator, indicating its adaptability with various classifier architectures. Moreover, obtaining features from the initial layers of the classifier for constructing the PIM (Prior Induced Model) proves to be more effective than sourcing them from the deeper layers, aligning with our previous findings.

Table 6: Performance Comparison for ViT-B-16 architecture on the Waterbirds dataset

Dataset	Method	FR	SR	MCC	
Waterbirds	MSP	0.1587	0.9048	0.0954	
	Energy	0.4924	0.6732	0.1656	
	Ent	0.3076	0.8301	0.1613	
	DECIDER layer1				
	+ mean	0.5592	0.7743	0.3416	
	+ max	0.6056	0.8091	0.4235	
	DECIDER layer9				
	+ mean	0.4818	0.6789	0.1613	
	+ max	0.5380	0.7542	0.2970	

**Detailed Results with PACS:** Expanding on the results provided in the main paper, we provide the failure detection performance metrics under the settings where classifier  $\mathcal{F}$  and PIM  $\phi$  are trained on different domains. For all experiments, we used early layer features of

the classifier. For each of these experiments, the failure estimation threshold is established based on the validation set from the respective training domain. The additional results are tabulated in Table 7 to Table 10.

Table 7: Performance Comparison on PACS dataset, where the classifier and the PIM are trained and calibrated on the *Art Painting* domain

Eval. Domain	Method	FR	SR	MCC	
Art Painting	MSP	0.7345	0.7799	0.4564	
	Energy	0.6381	0.7698	0.3659	
	Ent	0.6959	0.7837	0.4294	
	DECIDER				
	+ mean	0.7516	0.9822	0.7928	
	+ max	0.8458	0.9784	<b>0.8498</b>	
Cartoon	MSP	0.5636	0.6675	0.2204	
	Energy	0.5033	0.7188	0.2141	
	Ent	0.5799	0.6687	0.2371	
	DECIDER				
	+ mean	0.8394	0.6430	0.4895	
	+ max	0.8945	0.6027	<b>0.5284</b>	
Photo	MSP	0.5660	0.7857	0.3617	
	Energy	0.5406	0.8265	0.3851	
	Ent	0.5305	0.8254	0.3743	
	DECIDER				
	+ mean	0.6942	0.8594	0.5637	
	+ max	0.7754	0.8424	<b>0.6200</b>	
Sketch	MSP	0.6412	0.6088	0.2445	
	Energy	0.3252	0.8238	0.1639	
	Ent	0.6001	0.6252	0.2196	
	DECIDER				
	+ mean	0.8642	0.4977	0.3944	
	+ max	0.9066	0.4576	<b>0.4187</b>	



Table 8: Performance Comparison on PACS dataset, where the classifier and the PIM are trained and calibrated on *Cartoon* domain

<b>Eval. Domain</b>	<b>Method</b>	<b>FR</b>	<b>SR</b>	<b>MCC</b>	
Art Painting	MSP	0.4988	0.6999	0.1938	
	Energy	0.5467	0.6335	0.1739	
	Ent	0.4772	0.7118	0.1855	
	<hr/>				
	DECIDER				
	+ mean	0.6602	0.7556	<b>0.4011</b>	
+ max	0.6270	0.7849	0.3977		
Cartoon	MSP	0.6280	0.9206	0.4343	
	Energy	0.5427	0.9211	0.3761	
	Ent	0.5061	0.9349	0.3818	
	<hr/>				
	DECIDER				
	+ mean	0.6341	0.9950	<b>0.7430</b>	
+ max	0.5854	0.9950	0.7092		
Photo	MSP	0.4561	0.7660	0.2281	
	Energy	0.4819	0.7418	0.2266	
	Ent	0.4355	0.7974	0.2431	
	<hr/>				
	DECIDER				
	+ mean	0.6656	0.8916	<b>0.5552</b>	
+ max	0.6316	0.9016	0.5354		
Sketch	MSP	0.6033	0.6570	0.2497	
	Energy	0.5668	0.7372	0.2926	
	Ent	0.5470	0.7202	0.2575	
	<hr/>				
	DECIDER				
	+ mean	0.7604	0.8871	<b>0.6220</b>	
+ max	0.7132	0.9006	0.5887		

Table 9: Performance Comparison on PACS dataset, where the classifier and the PIM are trained and calibrated on *Photo* domain

Eval. Domain	Method	FR	SR	MCC
Art Painting	MSP	0.5364	0.5983	0.1220
	Energy	0.6269	0.5254	0.1399
	Ent	0.5658	0.5847	0.1365
	DECIDER			
	+ mean	0.6272	0.6955	<b>0.2913</b>
	+ max	0.5653	0.7266	0.2630
Cartoon	MSP	0.43532	0.56802	0.00258
	Energy	0.59117	0.51313	0.08078
	Ent	0.44831	0.55131	-0.00029
	DECIDER			
	+ mean	0.47292	0.66274	0.10499
	+ max	0.42448	0.75236	<b>0.13940</b>
Photo	MSP	0.5278	0.9835	0.4537
	Energy	0.5000	0.9633	0.3189
	Ent	0.5556	0.9859	0.4965
	DECIDER			
	+ mean	0.7143	0.9939	<b>0.7082</b>
	+ max	0.6571	0.9927	0.6498
Sketch	MSP	0.2226	0.8689	0.0886
	Energy	0.3440	0.9324	0.2377
	Ent	0.2176	0.8919	0.1077
	DECIDER			
	+ mean	0.4229	0.9424	<b>0.2996</b>
	+ max	0.4103	0.9263	0.2774

Table 10: Performance Comparison on PACS dataset, where the classifier and the PIM are trained and calibrated on *Sketch* domain

Eval. Domain	Method	FR	SR	MCC	
Art Painting	MSP	0.3836	0.6026	-0.0112	
	Energy	0.3317	0.6462	-0.0184	
	Ent	0.4331	0.5513	-0.0124	
	DECIDER				
	+ mean	0.9156	0.1769	0.1200	
	+ max	0.9710	0.1179	<b>0.1670</b>	
Cartoon	MSP	0.4536	0.6892	0.1326	
	Energy	0.5270	0.6129	0.1279	
	Ent	0.5215	0.6633	0.1692	
	DECIDER				
	+ mean	0.8830	0.5640	<b>0.4717</b>	
	+ max	0.8563	0.5338	0.4065	
Photo	MSP	0.3107	0.6667	-0.0179	
	Energy	0.2750	0.7074	-0.0145	
	Ent	0.3479	0.6481	-0.0031	
	DECIDER				
	+ mean	0.9679	0.1333	0.1734	
	+ max	0.9850	0.1148	<b>0.2116</b>	
Sketch	MSP	0.6822	0.9532	0.4221	
	Energy	0.3458	0.9314	0.1702	
	Ent	0.6449	0.9464	0.3778	
	DECIDER				
	+ mean	0.4673	0.9950	<b>0.5729</b>	
	+ max	0.4299	0.9639	0.3034	