# A Likelihood Ratio-Based Approach to Segmenting Unknown Objects

Nazir Nayal[1,2†], Youssef Shoeb[3,4†], Fatma Güney[1,2*]

[1*]Computer Engineering Department, Koç University, Turkey.
[2] KUIS AI Center, Turkey.
[3]Continental AG, Germany.
[4]Technische Universität Berlin, Germany.

*Corresponding author(s). E-mail(s): fguney@ku.edu.tr;
Contributing authors: nnayal21@ku.edu.tr; youssef.shoeb@continental.com;
[†]These authors contributed equally to this work.

## Abstract

Addressing the Out-of-Distribution (OoD) segmentation task is a prerequisite for perception systems operating in an open-world environment. Large foundational models are frequently used in downstream tasks, however, their potential for OoD remains mostly unexplored. We seek to leverage a large foundational model to achieve robust representation. Outlier supervision is a widely used strategy for improving OoD detection of the existing segmentation networks. However, current approaches for outlier supervision involve retraining parts of the original network, which is typically disruptive to the model's learned feature representation. Furthermore, retraining becomes infeasible in the case of large foundational models. Our goal is to retrain for outlier segmentation without compromising the strong representation space of the foundational model. To this end, we propose an adaptive, lightweight unknown estimation module (UEM) for outlier supervision that significantly enhances the OoD segmentation performance without affecting the learned feature representation of the original network. UEM learns a distribution for outliers and a generic distribution for known classes. Using the learned distributions, we propose a likelihood-ratio-based outlier scoring function that fuses the confidence of UEM with that of the pixel-wise segmentation inlier network to detect unknown objects. We also propose an objective to optimize this score directly. Our approach achieves a new state-of-the-art across multiple datasets, outperforming the previous best method by 5.74% average precision points while having a lower false-positive rate. Importantly, strong inlier performance remains unaffected.

**Keywords:** Anomaly Segmentation, Out-of-Distribution Detection, Likelihood Ratio, Unknown Segmentation, OoD Segmentation, Foundational Models for OoD.

## 1 Introduction

Semantic segmentation represents a significant advancement in deep learning. Learned features are densely mapped to a pre-defined set of classes by a pixel-level classifier. The remarkable performance of end-to-end models on this closed set has led researchers to consider the next challenge: extending semantic segmentation to the open-world setting where objects of unknown classes

also need to be segmented. One of the biggest challenges in segmenting unknown objects is the lack of outlier data.

In this work, we first attack the lack of data for unknown segmentation by utilizing a large foundation model, DINOv2 [36], for a robust representation space. The availability of internet-scale data has enabled the training of large visual foundation models, known for their generalization capabilities across various tasks [50, 5, 2, 35]. Despite these promising generalization capabilities, their potential for unknown object segmentation remains mostly unexplored. Only recently, PixOOD [45] has used DINOv2 without any training to avoid biases in industrial settings, however, their performance falls significantly behind the methods that use outlier supervision on commonly used SMIYC benchmark.

While collecting representative data for all possible classes in an open-world setting is impracticable, existing methods perform significantly better when trained using proxy outlier data [16, 33, 38], for example, obtained with the cut-and-paste method. Retraining with outlier supervision improves unknown segmentation but causes problems for known classes due to the reshaping of the representation space. Furthermore, retraining the entire model becomes infeasible in the case of large foundational models. We propose a novel way of utilizing proxy outlier data to improve the segmentation of unknown classes without compromising the performance of known classes.

Semantic segmentation models are typically trained to predict class probabilities with a softmax classifier. With a cross-entropy loss on the predicted class probabilities, the model learns to discriminate features of a certain class from the others. Such models excel in learning *discriminative* representations for the known classes but struggle to generalize to unknown classes due to partitioning the feature space between known classes. As an alternative, deep generative models directly learn a density model to predict the likelihood of a data sample. This likelihood is expected to be lower for outliers, such as samples from unknown classes. However, generative models often require more computational resources and can be challenging to train effectively.

Due to their potential to learn well-calibrated scores, deep generative models have been widely explored for out-of-distribution (OoD) tasks. However, in segmentation, their performance is often inferior to that of discriminative counterparts [24, 18, 47, 43]. To benefit from the best of both worlds, GMMSeg [26] presents a hybrid approach by augmenting the GMM-based generative model with discriminatively learned features. While discriminative features boost the inlier performance, GMM helps achieve an impressive OoD performance without explicitly training for it.

Nalisnick et al. [32] test the ability of deep generative models to detect OoD. They show that a generative model trained on one dataset assigns higher likelihoods to samples from another than those from the training dataset itself. Zhang et al. [48] first explain this phenomenon by showing that the expected log-likelihood is mathematically larger for out-of-distribution data and then propose to differentiate between outlier and OoD detection. While the learned density function can be used to detect outliers with respect to a single distribution, *OoD detection requires comparing two distributions*. As initially proposed by Bishop [3], OoD detection can be considered model selection between in-distribution and out-of-distribution data. Although an out-of-distribution is not often explicitly modeled, Zhang et al. [48] show that several existing works in OoD perform a likelihood ratio test with a proxy distribution for OoD, e.g. from auxiliary OoD datasets [20] or using background statistics [40].

In this paper, we propose applying the likelihood ratio as a principled way of detecting OoD in semantic segmentation. To calculate the likelihood ratio, we propose to train a lightweight unknown estimation module (UEM) on top of an already trained semantic segmentation model with a fixed number of semantic classes. UEM estimates an OoD distribution using proxy outlier data and a class-agnostic inlier distribution to calculate the likelihood ratio score. We also propose an objective to optimize the likelihood ratio score and train UEM with this objective. We show that our formulation is general enough to apply to both discriminative and generative segmentation models, with an example for each in the experiments. Our proposed method achieves state-of-the-art performance on multiple benchmarks while maintaining the same inlier performance.

# 2 Related Work

**OoD without Outlier Data:** Earlier approaches for OoD detection rely on uncertainty estimation methods to model predictive uncertainty. The uncertainty of a model can be estimated through maximum softmax probabilities [19], ensembles [23], MC-dropout [12], or by learning to estimate the confidence directly [22]. However, posterior probabilities in a closed-set setting may not always be well-calibrated for an open-world setting, potentially leading to overly confident predictions for unfamiliar categories [17, 21, 31].

**OoD with Outlier Data:** Hendrycks et al. [20] introduce outlier exposure to improve OoD detection. Outlier exposure leverages a proxy dataset of outliers to discover signals and learn heuristics for OoD samples. Chan et al. [7] use a proxy dataset and entropy maximization to fine-tune the model to give high entropy scores to unknown samples. Similarly, RbA [33] uses a proxy dataset to fine-tune the model to produce low logit scores on unknown objects. We follow a similar approach in our work and use a proxy dataset to learn a proxy distribution of OoD. However, our proxy dataset is only used to adjust the parameters of a small discriminator model, so it does not affect the performance of the inlier model.

**Deep Generative Models for OoD:** Generative models have been used to identify outliers based on the estimated probability density of the inlier training data distribution. Liang et al. [26] use a mixture of Gaussians to represent the data distribution within each class and model OoD instances as low-density regions. Other methods use normalizing flow [4, 15] or an energy-based model [14] to estimate inlier data density. However, estimating a data density of inliers only does not behave as expected for OoD detection, as Nalisnick et al. [32] show in their analysis of several deep generative models. Instead of a single density estimation, we treat OoD detection as model selection between two distributions as proposed in [48]. We directly train the model to optimize the likelihood ratio between an in-distribution and an out-of-distribution for a better separation of outliers. To our knowledge, this is the first work to consider the likelihood ratio for segmenting outliers.

**Mask-Based OoD:** A recent trend in OoD segmentation is to use mask-based models by predicting and classifying masks [9, 8, 25]. In masked-based models such as Mask2Former [8], each query specializes in detecting a certain known class [33, 1]. Based on this property of mask-based models, RbA [33] proposes an outlier scoring function based on the probability of not belonging to any known classes. Utilizing the same property, Maskomaly [1] selects outlier masks by thresholding the per-class mIoU on a validation set. Mask2Anomaly [38] augments Mask2Former with a global masked-attention mechanism and trains it using a contrastive loss on outlier data. EAM [16] performs OoD detection via an ensemble over mask-level scores. Almost all of these methods, except for Maskomaly [1], which is a simple inference-time post-processing technique, show the importance of utilizing OoD data during training. In this paper, we propose a better way of utilizing outlier data with the likelihood ratio, outperforming mask-based models in most metrics with pixel-based classification.

**Foundational Models for OoD:** Foundational models trained on large datasets have shown impressive zero-shot performance on downstream tasks like classification and segmentation [37, 36, 39]. For image-level OoD classification, Vojir et al. [44] leverage generic pre-trained representation from CLIP [37]. Wang et al. [46] train a negation text-encoder to equip CLIP with the ability to separate OoD samples from in-distribution samples. Recently, PixOoD [45] utilizes DINOv2 [36] for modeling the in-distribution data and achieves competitive results for OoD segmentation without using any outlier training. Initial work started exploring the potential of foundational models for OoD by building on their powerful representations. In this work, we take it further and improve outlier performance by retraining with outlier supervision without affecting the representation space of the foundational model.

# 3 Methodology

## 3.1 Overview

We propose a two-stage approach: In the first stage, a semantic segmentation model is trained solely on the known data with the standard segmentation losses. In the second stage, the semantic
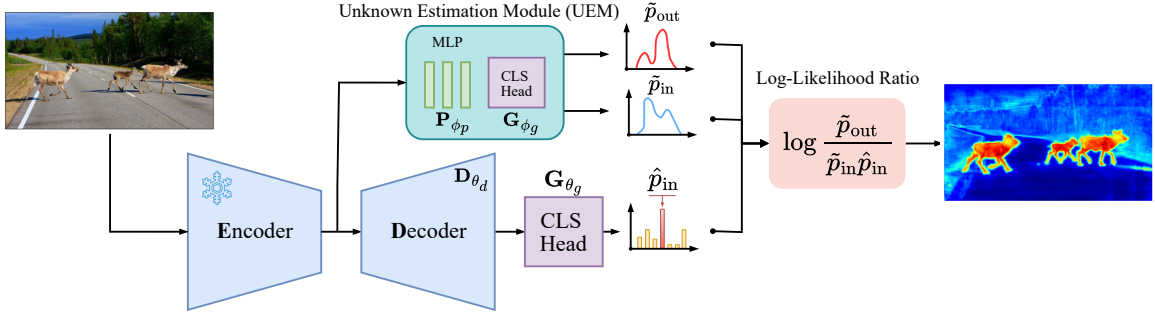
**Fig. 1 Overview.** Our proposed unknown estimation module (UEM) takes the input from the frozen encoder backbone and learns the outlier and inlier distributions $\tilde{p}_{\text{out}}$ and $\tilde{p}_{\text{in}}$. Then, we calculate the log-likelihood ratio by combining the outputs of UEM with the class probabilities of the inlier model $\hat{p}_{\text{in}}$.

segmentation model is fully frozen to maintain its exact inlier performance. We train an adaptive, lightweight unknown estimation module that estimates $\tilde{p}_{\text{out}}$ and a generic inlier distribution $\tilde{p}_{\text{in}}$ after injecting the training dataset with pseudo-unknown pixels in the second stage. With this setup, we propose an OoD scoring function based on the likelihood ratio by combining the output of this module and the inlier part and propose a loss function to optimize it. In Fig. 1, we provide an overview of our proposed approach.

## 3.2 Notation and Preliminaries

Given an input image $\mathbf{x} \in \mathbb{R}^{3 \times H \times W}$ and its corresponding label map $\mathbf{y} \in \mathcal{Y}^{H \times W}$, a closed-set semantic segmentation model learns a mapping from the input pixels to the class logits $\mathbf{F}_\theta(\mathbf{x})$ : $\mathbb{R}^{3 \times H \times W} \rightarrow \mathbb{R}^{K \times H \times W}$, where $\mathcal{Y} = \{1, \ldots, K\}$ is the set of known class labels during training. In OoD segmentation, we extend the label space to $\mathcal{Y}' = \mathcal{Y} \cup \{K + 1\}$, where $K + 1$ represents semantic categories unseen during training or the OoD class. To identify pixels belonging to the class $K + 1$, we define a scoring function $\mathcal{S}_{\text{out}}(\mathbf{x}) \in \mathbb{R}^{H \times W}$ that assigns high values to OoD pixels and low values to inlier pixels belonging to $\mathcal{Y}$.

**Likelihood Ratio:** Previous work in image-level OoD detection [32] has shown that when $\mathcal{S}_{\text{out}}(x)$ for an image is defined using the likelihood density of the training data, it assigns high likelihood values to some OoD samples. This limitation of likelihood-based methods has been mitigated by defining $\mathcal{S}_{\text{out}}(x)$ as the likelihood ratio (LR) between two distributions: $p_{\text{in}}$ representing the

likelihood of the sample belonging to the inlier distribution, and $p_{\text{out}}$ representing the likelihood of a pixel $x$ is an outlier. Formally:

$$\text{LR}(x) = \frac{p_{\text{out}}(x)}{p_{\text{in}}(x)} \qquad (1)$$

In this formulation, the likelihood of a sample being an inlier is reinforced by the likelihood of it not being an outlier, and vice versa. While defining $p_{\text{in}}$ is done using the inlier dataset, defining $p_{\text{out}}$ is challenging due to the unbounded diversity of $p_{\text{out}}$ compared to $p_{\text{in}}$. Therefore, using different assumptions, previous work explores approximating $p_{\text{out}}$ [40, 51]. In this work, we explore representing $p_{\text{out}}$ by means of both a generative and discriminative class-conditional distribution. We achieve these by utilizing synthetic pseudo-OoD objects based on driving scenes as done in previous works [38, 33, 16].

## 3.3 Learning an Inlier Segmentor

The existing pixel-level inlier segmentation models typically consist of three parts:

i. a feature extractor $\mathbf{E}$ : $\mathbb{R}^{3 \times H \times W} \mapsto \mathbb{R}^{C_e \times \dot{H} \times \dot{W}}$, reducing spatial dimension to $\dot{H} \times \dot{W}$,

ii. a decoder $\mathbf{D}_{\theta_d}$ : $\mathbb{R}^{C_e \times \dot{H} \times \dot{W}} \mapsto \mathbb{R}^{C_d \times H \times W}$, increasing it back to the original $H \times W$, and

iii. a classification head $\mathbf{G}_{\theta_g}$ : $\mathbb{R}^{C_d \times H \times W} \mapsto \mathbb{R}^{K \times H \times W}$ mapping features to class logit scores.

$C_e$ and $C_d$ denote the encoder and decoder's hidden dimension size, respectively. Hence, the mapping $\mathbf{F}_\theta(\mathbf{x}) : \mathbb{R}^{3 \times H \times W} \mapsto \mathbb{R}^{K \times H \times W}$ is defined

as $\mathbf{F}_\theta = \mathbf{E} \circ \mathbf{D}_{\theta_d} \circ \mathbf{G}_{\theta_g}$. In this notation, $\theta$ is the set of all learnable parameters and contains the union of $\theta_d$ and $\theta_g$. In some cases, features from multiple layers of the encoder $\mathbf{E}$ can be passed on to the decoder $\mathbf{D}$ to process features in a multi-scale fashion. We omit this in the notation for simplicity.

For the backbone, we use DINOv2 [36], which is a self-supervised ViT [11] that has been shown to produce robust and rich visual representations [39]. To maintain its rich representation, we freeze the backbone throughout all stages of training. For the decoder, we utilize a standard Feature Pyramid Network (FPN) [27] that takes features from multiple layers of the encoder and fuses them to produce an output feature map. For the classification head, we explore using two types of classifiers: generative and discriminative. Although the discriminative version seems less suitable for likelihood computations, we show that it performs exceptionally well under certain assumptions.

**Generative Classifier:** We adopt the generative classification formula proposed in [26], which replaces the linear softmax classification head by learning class densities of each pixel $p(x|k)$ with Gaussian Mixture Models (GMMs), where each class is represented with a separate GMM with a uniform prior on the component weights. Formally:

$$p(x|k, \theta_g) = \sum_{c=1}^{C} \pi_{kc} \, \mathcal{N}(x; \mu_{kc}, \Sigma_{kc}) \qquad (2)$$

where $C$ is the number of components per GMM, $\pi_{kc}$ is the component mixture weight for component $c$ of class $k$, $\mu_{kc}, \Sigma_{kc}$ are the mean and covariance matrix respectively, and $\mathcal{N}$ is the Gaussian distribution. The GMM parameters are learned with a variant of the Expectation-Maximization (EM) algorithm called Sinkhorn EM, which adds constraints that enforce an even assignment of features to mixture components, thereby improving the training stability. For more details, please refer to [26].

**Discriminative Classifier:** We train a single linear layer as a discriminative classifier. In this version, the parameters of the model $\theta$ are

supervised by the cross-entropy loss:

$$\theta^* = \mathrm{argmin}_\theta - \sum_{(x,k) \in \mathcal{D}} \log p(k|x, \theta) \qquad (3)$$

where $\mathcal{D}$ is the set of image-label pairs, and $p(k|x, \theta)$ is the softmax output of class $k$ after mapping it to class logits first, $\hat{y} = \mathbf{F}_\theta(x)$:

$$p(k|x, \theta) = \frac{\exp(\hat{y}_k)}{\sum_{k'} \exp(\hat{y}_{k'})} \qquad (4)$$

## 3.4 Unknown Estimation Module (UEM)

At this stage, we assume the existence of an inlier segmentation model trained as described in Section 3.3. The unknown estimation module (UEM) consists of a projection module $\mathbf{P}_{\phi_p} \in \mathbb{R}^{C_p \times H \times W}$, where $C_p$ is the hidden size output for the projection module, which is a 3-layer Multi-Layer Perceptron (MLP) that takes the output of the frozen backbone and produces a projected feature map as follows:

$$\mathbf{P}_{\phi_p}(\mathbf{x}) = \mathrm{MLP}\big(\mathbf{E}(\mathbf{x})\big) \qquad (5)$$

After that, the projected features are fed to a classification head $\mathbf{G}_{\phi_g}$ with two classes: one class maps to the OoD distribution and another to a generic inlier distribution learned directly from the backbone. The classifier head can be defined in a generative or discriminative fashion as explained in Section 3.3. Hence, the output of UEM: $\mathbf{U} \in \mathbb{R}^{2 \times H \times W}$ is defined as follows:

$$\mathbf{U}_\phi(\mathbf{x}) = \mathbf{G}_{\phi_g}\big(\mathbf{P}_{\phi_p}(\mathbf{x})\big) \qquad (6)$$

From the outputs of this module, we denote $\tilde{p}_{\mathrm{out}} = \mathbf{U}_1(\mathbf{x})$ and $\tilde{p}_{\mathrm{in}} = \mathbf{U}_0(\mathbf{x})$ as the likelihood of a sample $\mathbf{x}$ being an outlier or an inlier respectively. Hence, in the case of the generative classifier, the likelihoods would be those of the learned GMMs. In the case of a discriminative classifier, these would correspond to the class posterior probabilities $p(c|x)$ as the output of $\mathbf{G}_{\phi_g}$.

## 3.5 Log-Likelihood Ratio Score

First, we outline the formulation assuming a generative classifier for both the inlier model and the UEM module.

**Generative:** We propose the log-likelihood ratio as an OoD scoring function $\mathcal{S}_{\text{out}}$ where the likelihood ratio is defined in (1). For this, we need to define $p_{\text{out}}$ and $p_{\text{in}}$. We simply set the outlier distribution $p_{\text{out}} = \tilde{p}_{\text{out}}$, where $\tilde{p}_{\text{out}}(\mathbf{x}) \sim$ GMM with a uniform prior $\pi_c = \frac{1}{C}$:

$$\tilde{p}_{\text{out}}(\mathbf{x}) = \sum_{c=1}^{C} \pi_c \mathcal{N}(\mathbf{x}; \mu_c, \Sigma_c) \tag{7}$$

where $C$ is the number of components and $\mathcal{N}$ is the Gaussian distribution. As for the inlier distribution $p_{\text{in}}$, we define it by combining $\tilde{p}_{in}$ with the likelihood that a sample is inlier based on the inlier segmentation model. Due to the independence of the two sources of inlier confidence, $p_{\text{in}}$ can be defined as their product: $p_{\text{in}} = \tilde{p}_{\text{in}} \cdot \hat{p}_{\text{in}}$. In this case, $\tilde{p}_{\text{in}}(\mathbf{x})$ follows the same form of $\tilde{p}_{\text{out}}(\mathbf{x})$ in (7). As for $\hat{p}_{\text{in}}(\mathbf{x})$, we have:

$$\hat{p}_{\text{in}}(\mathbf{x}) = \max_{k} \log p(k|\mathbf{x}) \tag{8}$$

However, since the inlier segmentation model is a generative classifier, we have $p(k|\mathbf{x}) \sim$ GMM. And following [26], the log probability is used as class logit scores, hence $\mathbf{F}_k(\mathbf{x}) = \log p(k|\mathbf{x})$, which makes the full log-likelihood ratio score as follows:

$$\begin{aligned} \text{LLR}(\mathbf{x}) &= \log \frac{p_{\text{out}}(\mathbf{x})}{p_{\text{in}}(\mathbf{x})} \\ &= \log p_{\text{out}}(\mathbf{x}) - \log p_{\text{in}}(\mathbf{x}) \\ &= \log \tilde{p}_{\text{out}}(\mathbf{x}) - \log \tilde{p}_{\text{in}}(\mathbf{x}) - \max_{k} \log p(k|\mathbf{x}) \\ &= \log \tilde{p}_{\text{out}}(\mathbf{x}) - \log \tilde{p}_{\text{in}}(\mathbf{x}) - \max_{k} \mathbf{F}_k(\mathbf{x}) \end{aligned} \tag{9}$$

LLR unifies the confidence values of the inlier model and our proposed UEM in a single objective.

**Discriminative:** In this case, $p_{\text{in}}$ and $p_{\text{out}}$ are defined as $\tilde{p}_{\text{in}}$ and $\tilde{p}_{\text{out}}$ respectively. The difference compared to the generative case is that these terms are not defined as GMMs, but rather as logits computed through a linear classifier. As for $\hat{p}_{\text{in}}$, we define it to be the maximum of class logits defined in (4). Hence, the LLR score in this case can be written as:

$$\begin{aligned} \text{LLR}(\mathbf{x}) &= \log \frac{p_{\text{out}}(\mathbf{x})}{p_{\text{in}}(\mathbf{x})} \\ &= \log p_{\text{out}}(\mathbf{x}) - \log p_{\text{in}}(\mathbf{x}) \\ &= \log \tilde{p}_{\text{out}}(\mathbf{x}) - \log \tilde{p}_{\text{in}}(\mathbf{x}) - \max_{k} \log p(k|x) \\ &= \log \tilde{p}_{\text{out}}(\mathbf{x}) - \log \tilde{p}_{\text{in}}(\mathbf{x}) \\ &\quad - \max_{k} \left( \mathbf{F}_k(\mathbf{x}) + \log \sum_{k' \in \mathcal{Y}} \exp \left( \mathbf{F}_{k'}(\mathbf{x}) \right) \right) \end{aligned} \tag{10}$$

As the normalization term $\log \sum_{k \in \mathcal{Y}} \exp(\mathbf{F}(\mathbf{x}))$ does not affect the maximum, we obtain the final form of the scoring function as follows:

$$\text{LLR}(\mathbf{x}) = \log \tilde{p}_{\text{out}}(\mathbf{x}) - \log \tilde{p}_{\text{in}}(\mathbf{x}) - \max_{k} \mathbf{F}_k(\mathbf{x}) \tag{11}$$

This shows that the generative and discriminative case formulations converge to the same equation.

## 3.6 Log-Likelihood Ratio Loss

The proposed unknown estimation module $\mathbf{U}_\phi$ is supervised by the LLR loss defined as follows:

$$\mathcal{L}_{\text{LLR}}(\mathbf{x}, \tilde{\mathbf{y}}) = \text{BCE}(\text{LLR}(\mathbf{x}), \tilde{\mathbf{y}}) + \alpha \, \mathcal{L}_{\text{GMM}} \tag{12}$$

where $\tilde{\mathbf{y}}$ is a binary label map denoting known and pseudo-outlier pixels, BCE is the Binary Cross Entropy Loss, and $\mathcal{L}_{\text{GMM}}$ is the loss used to train the GMM component in case the generative classifier used as in [26]. The $\mathcal{L}_{\text{GMM}}$ consists of two terms as follows:

$$\mathcal{L}_{\text{GMM}} = \mathcal{L}_{\text{CE}} + \beta \mathcal{L}_{\text{contrast}} \tag{13}$$

where $\beta$ is a weighting coefficient. $\mathcal{L}_{\text{CE}}$ is the cross-entropy loss which is applied on the output logit scores of the generative classifier head $\mathbf{F}(\mathbf{x})$, which as shown in [26] corresponds to the class log probabilities of the class GMMs. As for $\mathcal{L}_{\text{contrast}}$, this loss is applied to contrast between every component within every class GMM with all other components, including those with the same class and of the other classes. In other words, using the Sinkhorn algorithm, each feature in the image is assigned to a unique component within the GMM

of its ground truth class. If there are $K$ classes and $C$ components with each $GMM$, then the cross entropy loss optimizes the likelihood feature belonging to its assigned component and minimizes it for the other $CK - 1$ components within the classification head.

# 4 Experiments

## 4.1 Experimental Setup

In our experiments, we use an inlier segmentation network composed of a feature extractor, an FPN [27] pixel decoder, and a generative classification head (GMMSeg [26]). The feature extractor is frozen, and we train the pixel decoder and segmentation head in the first stage on random patches of size $518 \times 1036$ taken from the Cityscapes dataset [10]. In the second stage, we train our unknown estimation module using a modified version of Anomaly Mix [42], where we randomly cut and past objects from the COCO dataset [28] on the training data for outlier supervision. During outlier supervision, all the trained parameters of the main segmentation network are frozen to maintain inlier performance. Finally, we maintain the training resolution during inference but with a sliding window approach to cover the whole image.

**Evaluation Datasets and Metrics:** We report the performance on SMIYC [6] Anomaly Track (SMIYC-AT), Obstacle Track (SMIYC-OT), RoadAnomaly [29], and the validation set of Fishyscapes LostandFound (FS LaF) [4]. SMIYC-AT and RoadAnomaly are real-world images featuring one or several OoD objects of varying sizes and categories. SMIYC-OT and FS LaF assess the model's capability to identify small-sized obstacles on the road. We evaluate the performance of our method using common pixel-wise anomaly segmentation metrics: Average Precision ($AP$) and False Positive Rate ($FPR$) at True Positive Rate of 95%.

## 4.2 Quantitative Results

**Backbone Feature Extractor:** The proposed Unknown Estimation Module (UEM) builds on a strong backbone model as the feature extractor. The backbone plays a critical role by encoding images into a rich representation space, which

helps first model the inliers and then differentiate the outliers with the UEM. We compare the performance of three different backbones for feature extraction, including a self-supervised one, DINOv2 [36]; a contrastive one, CLIP [37]; and a supervised one, Hierarchical Swin Transformer [30] for the baseline segmentation network.

Table 1 shows the mIoU performance of the inlier network using different backbones on Cityscapes and the anomaly segmentation performance of UEM trained on top of the inlier network on RoadAnomaly and Fishyscapes. While DINO and Swin show comparable performance on inlier data, DINO significantly outperforms Swin in handling outliers. CLIP shows lower inlier performance than both but surpasses Swin in outlier detection. This difference in outlier performance can be attributed to the pre-training of DINO and CLIP on larger and more diverse datasets, which results in more robust feature representations capable of effectively modeling both in-distribution and out-of-distribution.

**Improvements from Likelihood Ratio:**

We question whether the likelihood ratio is necessary for unknown segmentation. To investigate this, we push the performance of a generative model that requires no additional outlier training by using more powerful backbones. GMMSeg estimates class probability densities, allowing it to directly compute an anomaly score based on the likelihood of the maximum component without requiring outlier training. We use GMMSeg's density estimate (ID) as a baseline. We also consider the density estimate of the proxy OoD distribution alone as a scoring function (OoD). Lastly, we use the likelihood ratio scoring (LR), which integrates information from both distributions. Table 2 illustrates the performance improvements of the two scoring functions compared to the density estimates from the GMMSeg. The results consistently demonstrate that the likelihood ratio formulation provides better performance over inlier density estimates or the OoD scoring alone, highlighting the advantages of our approach.

We qualitatively compare the three scoring functions in Fig. 2. The in-distribution (ID) score demonstrates lower precision due to its tendency to favor known classes. In contrast, the out-of-distribution (OoD) scoring detects outliers very

| Backbone | mIoU↑ | Road Anomaly | | | FS LaF | | |
|---|---|---|---|---|---|---|---|
| | | $AUROC\uparrow$ | $AP\uparrow$ | $FPR\downarrow$ | $AUROC\uparrow$ | $AP\uparrow$ | $FPR\downarrow$ |
| Swin-b [30] | 81.6 | 92.80 | 65.42 | 27.96 | 96.65 | 43.53 | 18.88 |
| CLIP [37] | 77.8 | 97.84 | 87.84 | 10.81 | 99.95 | 55.84 | 4.12 |
| DINOv2-b [36] | 82.8 | 98.37 | 92.86 | 8.95 | 98.64 | 64.02 | 1.62 |

**Table 1  Ablation of Backbone Feature Extractor.** We compare the inlier and OoD detection performance using different backbones. We find the DINOv2 backbone to offer the best inlier and outlier performance.

| Backbone | Scoring | Road Anomaly | | | FS LaF | | |
|---|---|---|---|---|---|---|---|
| | | $AUROC\uparrow$ | $AP\uparrow$ | $FPR\downarrow$ | $AUROC\uparrow$ | $AP\uparrow$ | $FPR\downarrow$ |
| Swin-b [30] | ID | 80.14 | 32.34 | 56.53 | 82.54 | 5.51 | 69.16 |
| | OoD | 91.63 ↑11.49 | 61.48 ↑29.14 | 32.34 ↑24.19 | 96.46 ↑13.92 | 38.30 ↑32.78 | 18.10 ↑51.06 |
| | LR | 92.80 ↑12.66 | 65.42 ↑33.08 | 27.96 ↑28.57 | 96.65 ↑14.11 | 43.53 ↑38.02 | 18.88 ↑50.28 |
| CLIP [37] | ID | 90.35 | 50.71 | 34.72 | 91.97 | 13.28 | 38.70 |
| | OoD | 96.78 ↑6.43 | 83.54 ↑32.83 | 17.41 ↑17.31 | 98.56 ↑6.59 | 35.70 ↑22.42 | 6.05 ↑32.65 |
| | LR | 97.84 ↑7.49 | 87.84 ↑37.13 | 10.81 ↑23.91 | 99.95 ↑7.98 | 55.84 ↑42.56 | 4.12 ↑34.58 |
| DINOv2-b [36] | ID | 92.94 | 64.77 | 28.25 | 90.83 | 15.00 | 40.70 |
| | OoD | 97.45 ↑4.51 | 88.54 ↑23.77 | 13.55 ↑14.7 | 99.19 ↑8.36 | 63.53 ↑48.53 | 3.76 ↑36.94 |
| | LR | 98.37 ↑5.43 | 92.86 ↑28.09 | 8.95 ↑19.3 | 99.47 ↑8.64 | 79.02 ↑64.02 | 1.62 ↑39.08 |

**Table 2  Likelihood Ratio Gains Ablation.** We compare the OoD detection performance o different backbones and with different scoring functions on the Road Anomaly and FS LaF datasets. Gains/losses to the base ID scoring are highlighted in green/red, respectively.

confidently but at the cost of increasing false positives so as not to miss any outliers. The proposed likelihood ratio (LR) balances the two, leveraging the strengths of each to achieve the best results in terms of both inliers and outliers.

**Comparison to State-of-the-Art:** Table 3 shows our results compared to state-of-the-art methods on four datasets with the average performance over datasets in the last column. As each dataset has different characteristics, the existing methods behave differently across the datasets. The top-performing methods include recently proposed masked-based models RbA [33], EAM [14], and Mask2Anomaly [38]. While these methods achieve impressive performance in terms of accuracy, reasoning at the mask level hurts FPR, as considering a mask outlier introduces several false positives at once. Its negative effect on small objects can be seen by high FPR on SMIYC-OT and FS LaF. Our method achieves significantly lower FPR on these two datasets while being among the top-performing methods in terms of AP. Our method also achieves impressive accuracy levels on real-world images of SMIYC-OT and Road Anomaly, increasing AP by 1.85 and 1.18,

respectively, without causing high FPR. Averaged across the four datasets in the last column, our method sets a new state-of-the-art in both metrics, outperforming the previous state-of-the-art by 3.71% in AP and 0.27% in FPR.

Training data impacts the performance significantly. Both RbA and EAM are trained on Mapillary and Cityscapes datasets, whereas we train our inlier model only on Cityscapes. Additionally, EAM uses ADE20K [52] for outlier supervision, which contains a broader range of classes than COCO. We only use COCO to ensure a fair comparison to other methods. We also note that outlier supervision used in most other methods negatively impacts the performance of the inlier segmentation network as reported in [42, 33].

**Is DINOv2 All You Need?** To assess the backbone's impact, we compare our approach to PEBAL [42] and RbA [33]. We use their scoring functions to train our segmentation network with the DINOv2 backbone and adjust the outlier supervision process to their original implementations. As shown in Table 4, DINOv2 significantly improves PEBAL's performance on Road Anomaly across both metrics and results in a
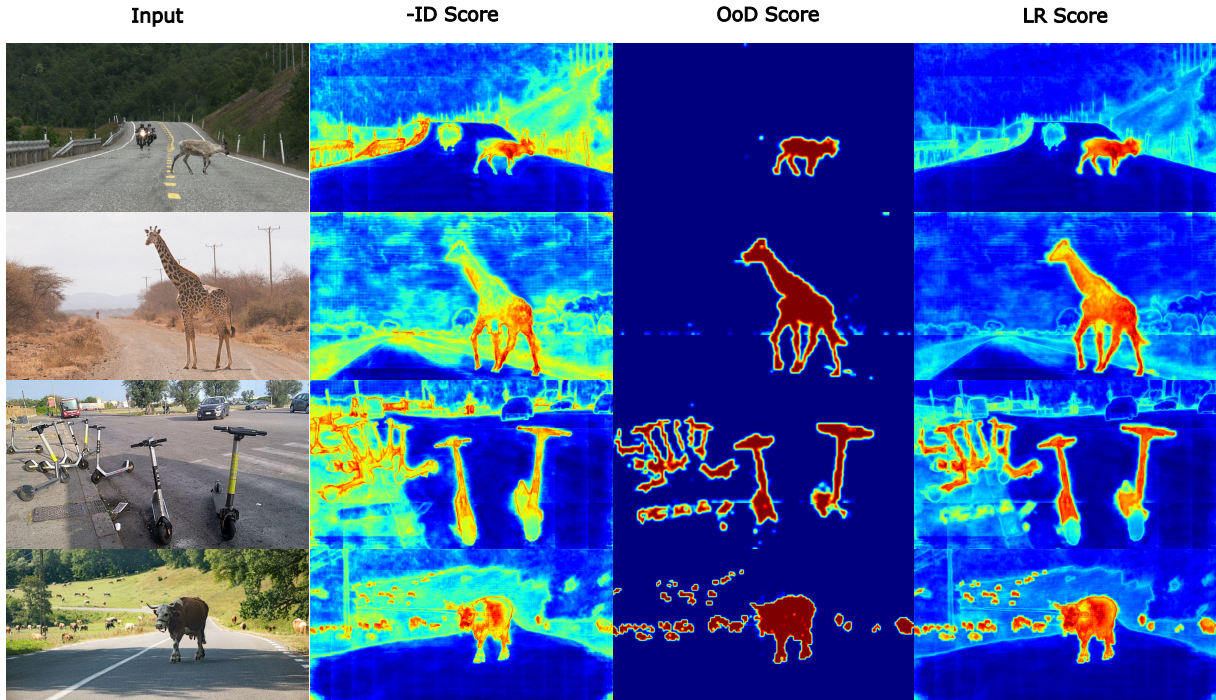
| Input | -ID Score | OoD Score | LR Score |

**Fig. 2 Qualitative Results on SMIYC-AT (G-G).** The second column (-ID score) shows the outlier score from using the GMM without outlier supervision. The third column (OoD Score) shows the anomaly score from the fine-tuned OoD detection head. The fourth column (LR Score) shows our proposed likelihood formulation. The likelihood formulation combines information from both and predicts more accurate OoD score maps.

lower FPR on FS LaF. We can attribute these improvements to the more robust backbone.

For RbA, AP on Road Anomaly improves, but other metrics are better using the original model with Mask2Former. This is likely due to the implicit one vs. all behavior in mask classification models, which the RbA scoring function is specifically designed for. Finally, our outlier scoring function performs best overall without modifying any original network parameters, a critical constraint for real-world applications. Both other methods are reported to lose at least 1% mIoU during fine-tuning. We found this effect exacerbated with DINO, which requires careful adjustments to mitigate inlier performance loss.

**Discriminative vs. Generative Modeling of Estimator module:** Our unknown estimation module models two distributions during fine-tuning. Each distribution can be modeled as a data density using generative GMMs or explicitly as a linear layer mapping function. In Table 3, we evaluated different possible combinations for each.

We find both discriminative and generative classifiers to outperform the previous state-of-the-art methods, with the fully discriminative classifier for the OoD modeling being slightly better. We omit the discriminative inlier and generative outlier (D-G) combination as we find the GMM takes too long to converge due to the unbounded range of values coming from the MLP.

**On the Number of Parameters in UEM:** The original segmentation network consists of 101M parameters. Our UEM module introduces an additional 788K parameters, representing a less than 1% increase in the overall model size. Despite this minimal parameter overhead, the UEM module significantly enhances OoD detection performance.

# 5 Conclusion and Future Work

In this work, we propose a novel strategy to utilize proxy outlier data for improved OoD detection without retraining the entire network. This allows us to build on the robust representation space of

**Table 3** comparison table:

| Method | SMIYC-AT | | SMIYC-OT | | RoadAnomaly | | FS LaF | | Average | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $AP \uparrow$ | $FPR \downarrow$ | $AP \uparrow$ | $FPR \downarrow$ | $AP \uparrow$ | $FPR \downarrow$ | $AP \uparrow$ | $FPR \downarrow$ | $\overline{AP} \uparrow$ | $\overline{FPR} \downarrow$ |
| cDNP [13] | - | - | - | - | 79.78 | 18.18 | 69.80 | 7.50 | - | - |
| UGainS [34] | - | - | - | - | 88.98 | 10.42 | 80.08 | 6.61 | - | - |
| Maximized Entropy [7] | 85.47 | 15.00 | 85.07 | 0.75 | - | - | 29.96 | 35.14 | - | - |
| CSL [49] | 80.08 | 7.16 | 87.10 | 0.67 | 61.38 | 43.80 | - | - | - | - |
| Maskomaly [1] | 93.35 | 6.87 | - | - | 70.90 | 11.90 | - | - | - | - |
| PEBAL [42] | 49.14 | 40.82 | 4.98 | 12.68 | 45.10 | 44.58 | 58.81 | 4.76 | 39.51 | 25.71 |
| Mask2Anomaly [38] | 88.72 | 14.63 | 93.22 | 0.20 | 79.70 | 13.45 | 46.04 | 4.36 | 76.92 | 8.16 |
| RbA [33] | 90.90 | 11.60 | 91.80 | 0.50 | 85.42 | 6.92 | 70.81 | 6.30 | 84.73 | 6.33 |
| EAM [16] | 93.75 | **4.09** | 92.87 | 0.52 | 69.40 | 7.70 | **81.50** | 4.20 | 84.38 | 4.13 |
| UEM (G-G) | 94.10 | 6.90 | 88.30 | 0.40 | 93.75 | **6.32** | 71.34 | 6.04 | 86.87 | 4.92 |
| UEM (G-D) | 92.50 | 11.30 | 92.0 | 0.20 | 92.86 | 8.95 | 79.02 | **1.62** | **89.10** | 5.52 |
| UEM (D-D) | **95.60** | 4.70 | **94.40** | **0.10** | 90.94 | 8.03 | 72.83 | 2.60 | 88.44 | **3.86** |

**Table 3 Quantitative Results on SMIYC-AT, SMIYC-OT, Road Anomaly, and FS LaF.** We compare our approach against existing OoD segmentation methods. Averaged across the four datasets, our approach sets a new state-of-the-art for both AP and FPR. The best result for each dataset is highlighted in **bold**, and the second best is underlined. Our method UEM (X-Y) has the flexibility to use a discriminative (D) or generative (G) modeling for the inlier classification head (X) and unknown estimation module (Y). We report the results of three possible combinations.

| Backbone | Road Anomaly | | | FS LaF | | |
|---|---|---|---|---|---|---|
| | $AUROC \uparrow$ | $AP \uparrow$ | $FPR \downarrow$ | $AUROC \uparrow$ | $AP \uparrow$ | $FPR \downarrow$ |
| PEBAL [42] | 98.32 ↑10.69 | 92.98 ↑47.88 | 7.13 ↑37.45 | 98.95 ↓0.01 | 51.11 ↓7.7 | 2.89 ↑1.87 |
| RbA [33] | 97.91 ↓0.08 | 89.37 ↑3.95 | 10.96 ↓4.04 | 98.39 ↓0.23 | 49.78 ↓21.03 | 4.64 ↑1.66 |
| UEM (Ours) | 98.08 | 92.86 | 8.95 | 98.75 | 79.02 | 1.62 |

**Table 4 Other Methods with DINOv2.** We compare the OoD performance of two other scoring functions using DINOv2 on Road Anomaly and FS LaF datasets. Gains/losses to the original method are highlighted in green/red, respectively. While DINOv2 improves the results of other methods, especially PEBAL on Road Anomaly, our method consistently achieves top results on both datasets, without affecting the inlier performance.

large foundational models, significantly enhancing the generalization capability of the proposed approach. We propose an unknown estimation module (UEM) that can be integrated into the existing segmentation networks to identify OoD objects effectively. We develop an OoD scoring function based on the likelihood ratio by combining UEM's outputs with inlier predictions. Our method sets a new state-of-the-art in outlier segmentation across multiple datasets, without causing any drops in the inlier performance.

For future work, we aim to investigate how the choice of proxy out-of-distribution (OoD) dataset influences the generalization performance of our method. In this study, we utilized the COCO dataset as the proxy OoD data for fair comparison with the other approaches. We plan to investigate the effect of mining more realistic outliers from real-world OoD datasets [41] as future work.

# References

[1] Ackermann, J., Sakaridis, C., Yu, F.: Maskomaly: Zero-shot mask anomaly segmentation. In: BMVC (2023)

[2] Aydemir, G., Xie, W., Guney, F.: Self-supervised object-centric learning for videos. In: NeurIPS (2023)

[3] Bishop, C.M.: Novelty detection and neural network validation. In: ICANN (1994)

[4] Blum, H., Sarlin, P.E., Nieto, J.I., Siegwart, R.Y., Cadena, C.: The fishyscapes benchmark: Measuring blind spots in semantic segmentation. IJCV **129**, 3119–3135 (2021)

[5] Blumenkamp, J., Morad, S., Gielis, J., Prorok, A.: Covis-net: A cooperative visual spatial foundation model for multi-robot applications. In: CoRL (2024)

[6] Chan, R., Lis, K., Uhlemeyer, S., Blum, H., Honari, S., Siegwart, R., Fua, P., Salzmann, M., Rottmann, M.: Segmentmeifyoucan: A benchmark for anomaly segmentation. In: NeurIPS (2021)

[7] Chan, R., Rottmann, M., Gottschalk, H.: Entropy maximization and meta classification for out-of-distribution detection in semantic segmentation. In: ICCV (2021)

[8] Cheng, B., Misra, I., Schwing, A.G., Kirillov, A., Girdhar, R.: Masked-attention mask transformer for universal image segmentation. In: CVPR (2022)

[9] Cheng, B., Schwing, A.G., Kirillov, A.: Per-pixel classification is not all you need for semantic segmentation. In: NeurIPS (2021)

[10] Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: CVPR (2016)

[11] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: ICLR (2021)

[12] Gal, Y., Ghahramani, Z.: Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In: ICML (2016)

[13] Galesso, S., Argus, M., Brox, T.: Far away in the deep space: Dense nearest-neighbor-based out-of-distribution detection. In: ICCV Workshops (2023)

[14] Grcić, M., Bevandić, P., Šegvić, S.: Densehybrid: Hybrid anomaly detection for dense open-set recognition. In: ECCV (2022)

[15] Grcić, M., Bevandić, P., Šegvić, S.: Dense anomaly detection by robust learning on synthetic negative data. Sensors (2024)

[16] Grcić, M., Šarić, J., Šegvić, S.: On advantages of mask-level recognition for outlier-aware segmentation. In: CVPR Workshops (2023)

[17] Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q.: On calibration of modern neural networks. In: ICML (2017)

[18] Haldimann, D., Blum, H., Siegwart, R.Y., Cadena, C.: This is not what I imagined: Error detection for semantic segmentation through visual dissimilarity. arXiv.org (2019)

[19] Hendrycks, D., Gimpel, K.: A baseline for detecting misclassified and out-of-distribution examples in neural networks. In: ICLR (2017)

[20] Hendrycks, D., Mazeika, M., Dietterich, T.G.: Deep anomaly detection with outlier exposure. In: ICLR (2019)

[21] Jiang, H., Kim, B., Guan, M., Gupta, M.: To trust or not to trust a classifier. In: NeurIPS (2018)

[22] Kendall, A., Gal, Y.: What uncertainties do we need in bayesian deep learning for computer vision? NeurIPS (2017)

[23] Lakshminarayanan, B., Pritzel, A., Blundell, C.: Simple and scalable predictive uncertainty estimation using deep ensembles. In: NeurIPS (2017)

[24] Lee, K., Lee, K., Lee, H., Shin, J.: A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In: NeurIPS (2018)

[25] Li, F., Zhang, H., xu, H., Liu, S., Zhang, L., Ni, L.M., Shum, H.Y.: Mask DINO: towards a unified transformer-based framework for object detection and segmentation. In: CVPR (2023)

[26] Liang, C., Wang, W., Miao, J., Yang, Y.: GMMSeg: Gaussian mixture based generative semantic segmentation models. In: NeurIPS (2022)

[27] Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: CVPR

(2017)

[28] Lin, T.Y., Maire, M., Belongie, S.J., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common objects in context. In: ECCV (2014)

[29] Lis, K., Nakka, K.K., Fua, P.V., Salzmann, M.: Detecting the unexpected via image resynthesis. In: ICCV (2019)

[30] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: ICCV (2021)

[31] Minderer, M., Djolonga, J., Romijnders, R., Hubis, F.A., Zhai, X., Houlsby, N., Tran, D., Lucic, M.: Revisiting the calibration of modern neural networks. In: NeurIPS (2021)

[32] Nalisnick, E.T., Matsukawa, A., Teh, Y.W., Görür, D., Lakshminarayanan, B.: Do deep generative models know what they don't know? In: ICLR (2019)

[33] Nayal, N., Yavuz, M., Henriques, J.F., Güney, F.: RbA: Segmenting unknown regions rejected by all. In: ICCV (2023)

[34] Nekrasov, A., Hermans, A., Kuhnert, L., Leibe, B.: UGainS: Uncertainty Guided Anomaly Instance Segmentation. In: GCPR (2023)

[35] Nguyen, V.N., Groueix, T., Ponimatkin, G., Lepetit, V., Hodan, T.: Cnos: A strong baseline for cad-based novel object segmentation. In: CVPR (2023)

[36] Oquab, M., Darcet, T., Moutakanni, T., Vo, H.V., Szafraniec, M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-Nouby, A., Assran, M., Ballas, N., Galuba, W., Howes, R., Huang, P.Y., Li, S.W., Misra, I., Rabbat, M., Sharma, V., Synnaeve, G., Xu, H., Jegou, H., Mairal, J., Labatut, P., Joulin, A., Bojanowski, P.: DINOv2: Learning robust visual features without supervision. Transactions on Machine Learning Research (2024)

[37] Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: ICML (2021)

[38] Rai, S., Cermelli, F., Fontanel, D., Masone, C., Caputo, B.: Unmasking anomalies in road-scene segmentation. In: ICCV (2023)

[39] Ranzinger, M., Heinrich, G., Kautz, J., Molchanov, P.: Am-radio: Agglomerative visual foundation model – reduce all domains into one. In: CVPR (2024)

[40] Ren, J.J., Liu, P.J., Fertig, E., Snoek, J., Poplin, R., DePristo, M.A., Dillon, J.V., Lakshminarayanan, B.: Likelihood ratios for out-of-distribution detection. In: NeurIPS (2019)

[41] Shoeb, Y., Chan, R., Schwalbe, G., Nowzad, A., Güney, F., Gottschalk, H.: Have we ever encountered this before? retrieving out-of-distribution road obstacles from driving scenes. In: WACV (2024)

[42] Tian, Y., Liu, Y., Pang, G., Liu, F., Chen, Y., Carneiro, G.: Pixel-wise energy-biased abstention learning for anomaly segmentation on complex urban driving scenes. In: ECCV (2022)

[43] Vojir, T., Šipka, T., Aljundi, R., Chumerin, N., Reino, D.O., Matas, J.: Road anomaly detection by partial image reconstruction with segmentation coupling. In: ICCV (2021)

[44] Vojíř, T., Šochman, J., Aljundi, R., Matas, J.: Calibrated out-of-distribution detection with a generic representation. In: ICCV Workshops (2023)

[45] Vojíř, T., Šochman, J., Matas, J.: PixOOD: Pixel-level out-of-distribution detection. In: ECCV (2024)

[46] Wang, H., Li, Y., Yao, H., Li, X.: CLIPN for zero-shot ood detection: Teaching CLIP to say no. In: ICCV (2023)

[47] Xia, Y., Zhang, Y., Liu, F., Shen, W., Yuille, A.: Synthesize then compare: Detecting failures and anomalies for semantic segmentation. In: ECCV (2020)

[48] Zhang, A., Wischik, D.: Falsehoods that ml researchers believe about ood detection. In: NeurIPS Workshops (2022)

[49] Zhang, H., Li, F., Qi, L., Yang, M.H., Ahuja, N.: Csl: Class-agnostic structure-constrained learning for segmentation including the unseen. AAAI (2024)

[50] Zhang, J., Herrmann, C., Hur, J., Cabrera, L.P., Jampani, V., Sun, D., Yang, M.H.: A tale of two features: Stable diffusion complements dino for zero-shot semantic correspondence. In: NeurIPS (2023)

[51] Zhang, M., Zhang, A., McDonagh, S.G.: On the out-of-distribution generalization of probabilistic image modelling. In: NeurIPS (2021)

[52] Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralba, A.: Scene parsing through

ade20k dataset. In: CVPR (2017)