

A Benchmark for Anomaly Segmentation

Dan Hendrycks*
UC Berkeley

Steven Basart*
University of Chicago

Mantas Mazeika
UIUC

Mohammadreza Mostajabi
TTIC

Jacob Steinhardt
UC Berkeley

Dawn Song
UC Berkeley

Abstract

Detecting out-of-distribution examples is important for safety-critical machine learning applications such as self-driving vehicles. However, existing research mainly focuses on small-scale images where the whole image is considered anomalous. We propose to segment only the anomalous regions within an image, and hence we introduce the Combined Anomalous Object Segmentation benchmark for the more realistic task of large-scale anomaly segmentation. Our benchmark combines two novel datasets for anomaly segmentation that incorporate both realism and anomaly diversity. Using both real images and those from a simulated driving environment, we ensure the background context and a wide variety of anomalous objects are naturally integrated, unlike before. Additionally, we improve out-of-distribution detectors on large-scale multi-class datasets and introduce detectors for the previously unexplored setting of multi-label out-of-distribution detection. These novel baselines along with our anomaly segmentation benchmark open the door to further research in large-scale out-of-distribution detection and segmentation.

1. Introduction

Detecting out-of-distribution inputs is important in real-world applications of deep learning. When faced with anomalous inputs flagged as such, systems may initiate a conservative fallback policy or defer to human judgment. This is especially important in safety-critical applications of deep learning, such as self-driving cars or medical screening. Accordingly, research on out-of-distribution detection has a rich history, and a host of methods have been developed [12, 30]. Recent work leveraging deep neural representations has proven successful at distinguishing anomalous inputs in complex domains, such as image data [18, 19, 24]. However, scaling up these approaches to work with higher-resolution scenes has proven challenging, in part due to the lack of

high-quality large-scale benchmarks.

Most previous formulations of anomaly detection on images treat entire images as anomalies. However, natural images are not monolithic entities, but rather are composed of numerous attributes, objects, and components, and are described by various orientations, lighting, etc. In practice, an image could be anomalous in certain specific regions while being in-distribution elsewhere. Knowing which regions of an image are anomalous could allow for safer handling of unfamiliar objects in the case of self-driving cars. Creating a benchmark for this task is difficult, though, as simply cutting and pasting anomalous objects into images introduces various unnatural giveaway cues that oversimplify and trivialize the task of anomaly segmentation, such as edge effects, mismatched orientation, and lighting [5].

To overcome these issues, we leverage a simulated driving environment to create a dataset for anomaly segmentation, which we call StreetHazards. Using the Unreal Engine and the open-source CARLA simulation environment [11], we insert a diverse array of foreign objects that the model has not encountered before into driving scenes and re-render the scenes with these novel objects. This enables integration of the foreign objects into their surrounding context with correct lighting and orientation.

To complement the StreetHazards dataset, we also introduce the BDD-Anomaly dataset of real images with anomalous objects. This is a dataset derived from the BDD100K semantic segmentation dataset [34]. Leveraging the large scale of BDD100K, we reserve infrequent object classes to be anomalies. We combine this dataset with StreetHazards to form the Combined Anomalous Object Segmentation (CAOS) benchmark. The CAOS benchmark improves over previous evaluations for anomaly segmentation in driving scenes by evaluating detectors on realistic and diverse anomalies. We evaluate several baselines on the CAOS benchmark and discuss problems with porting existing approaches from earlier formulations of out-of-distribution detection.

In addition to introducing the CAOS benchmark, we also explore large-scale settings for more traditional whole-image anomaly detection. Large-scale datasets such as ImageNet-

*Equal Contribution.

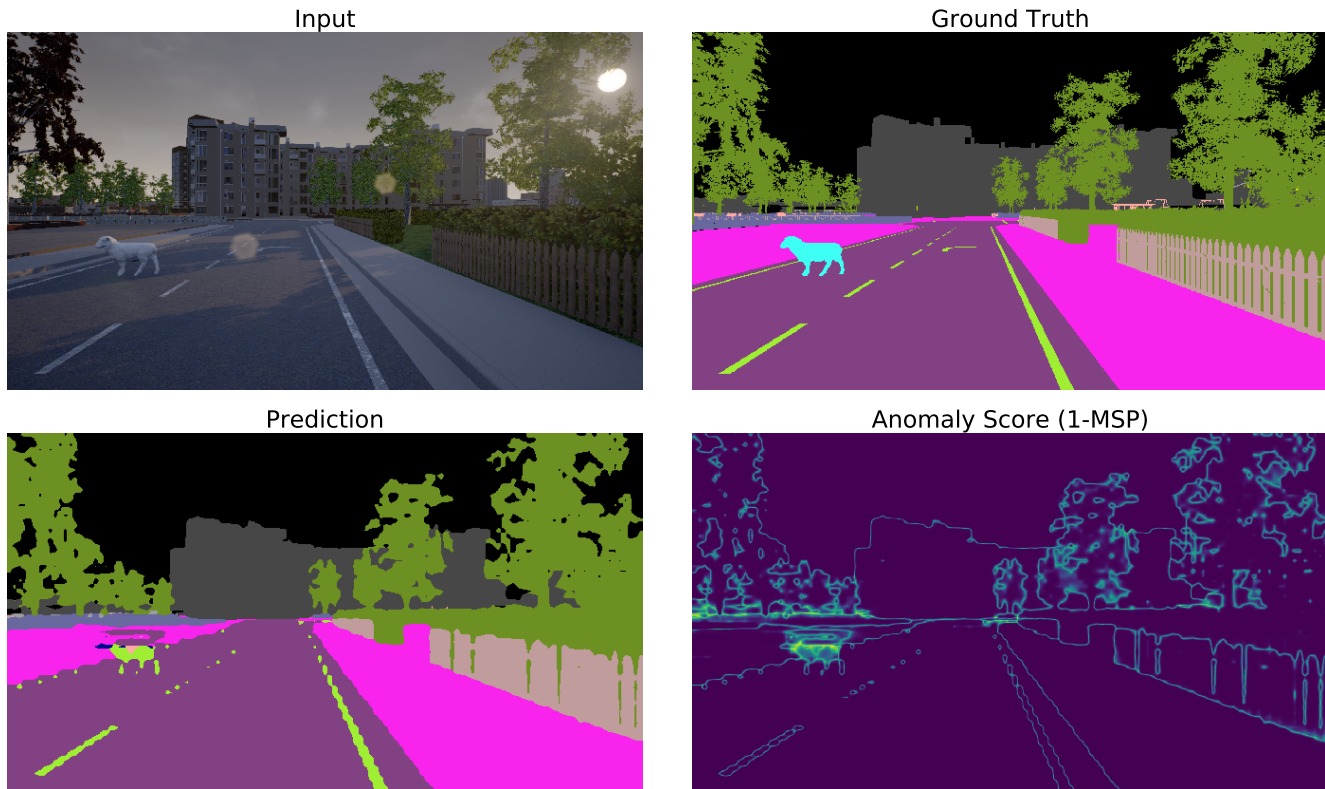


Figure 1: An example image from the StreetHazards dataset with a sheep anomaly in view. The maximum softmax probability (MSP) scores only partially detect this anomaly, and there are many false positives. Thus, much future work remains to develop effective techniques for anomaly segmentation.

1K [9] and Places365 [37] have many similar classes. We identify that the previous state-of-the-art out-of-distribution detector, when scaled up to datasets with similar classes, provides less meaningful anomaly scores. Hence we introduce a KL-divergence-based matching technique to improve standard out-of-distribution detection. Finally, we explore the use of multi-label classifiers for out-of-distribution detection and evaluate several novel detection mechanisms in this setting. As multi-label is a more natural setting than multi-class, this is an unexplored avenue that could grow in importance as research on out-of-distribution detection moves to large-scale, realistic evaluations. The CAOS benchmark datasets and code for experiments is available at <https://github.com/hendrycks/anomaly-seg>

2. Related Work

Anomaly Segmentation. Several prior works explore segmenting anomalous image regions. One line of work uses the WildDash dataset [35], which contains numerous annotated driving scenes in conditions such as snow, fog, and rain. The WildDash test set contains fifteen “negative images” from different domains, on which expected behavior is to mark the entire image as out-of-distribution. Thus, while the task

is segmentation, the anomalies do not exist as objects within an otherwise in-distribution scene. This setting is similar to that explored by [18], in which whole images from other datasets serve as out-of-distribution examples.

To approach anomaly segmentation on WildDash, Krešo et al. [22] train on multiple semantic segmentation domains and treat regions of images from the WildDash driving dataset as out-of-distribution if they are segmented as regions from different domains, i.e. indoor classes. Bevandić et al. [4] use ILSVRC 2012 images and train their network to segment the entirety of these images as out-of-distribution.

In medical anomaly segmentation and product fault detection, anomalies are regions of otherwise in-distribution images. Baur et al. [2] segment anomalous regions in brain MRIs using pixel-wise reconstruction loss. Similarly, Haselmann et al. [16] perform product fault detection using pixel-wise reconstruction loss. More recently, Bergmann et al. [3] introduce an expansive dataset for segmentation of product faults. The reconstruction-based approaches in these works require good modeling of the clean data to work. In contrast to medical anomaly segmentation and fault detection, we consider complex images from street scenes. These images have high variability in scene layout and lighting, and hence

	Fishyscapes	Lost and Found	BDD-Anomaly (Ours)	StreetHazards (Ours)
Train Images	0*	1036	6688	5125
Test Images	1000	1068	361	1500
Anomaly Types	12	9	2	250

Table 1: Quantitative comparison of the CAOS benchmark with related datasets. Although the BDD-Anomaly dataset has two object types, it has many unique object instances, while Lost and Found has only nine unseen objects at test time. *Fishyscapes models are trained on Cityscapes.

are less amenable to reconstruction-based techniques.

The two works closest to our own are the Lost and Found [27] and Fishyscapes [5] datasets. In Table 1, we quantitatively compare the CAOS benchmark to these datasets. The Lost and Found dataset consists of real images in a driving environment with small road hazards. The images were collected to mirror the Cityscapes dataset [8] but are only collected from one city and so have less diversity. It contains 35 unique anomalous objects, and methods are allowed to train on many of these. For Lost and Found, only nine unique objects are truly unseen at test time. Crucially, this is a different evaluation setting from our own, where anomalous objects are not revealed at training time, so their dataset is not directly comparable. Nevertheless, the BDD-Anomaly dataset that we introduce fills several gaps in Lost and Found. First, the images are more diverse because they are sourced from BDD100K, a more recent and comprehensive semantic segmentation dataset than Cityscapes. Second, the anomalies are not restricted to small, sparse road hazards. Concretely, anomalous regions in Lost and Found take up 0.11% of the image on average, whereas anomalous regions in the BDD-Anomaly dataset are larger and fill 0.83% of the image on average. Finally, although the BDD-Anomaly dataset has two anomaly types, compared to Lost and Found it has far more unique anomalous objects.

The Fishyscapes benchmark for anomaly segmentation consists of cut-and-paste anomalies from out-of-distribution domains. This is problematic, because the anomalies stand out as clearly unnatural in context. For instance, the orientation of anomalous objects is unnatural, and the lighting of the cut-and-paste patch differs from the lighting in the original image, providing an unnatural cue to anomaly detectors that would not exist for real anomalies. Techniques for detecting image manipulation [38] are competent at detecting artificial image elements of this kind. Our StreetHazards dataset overcomes these issues by leveraging a simulated driving environment to naturally insert anomalous *3D models* into a scene rather than by overlaying 2D images. These anomalies are integrated into the scene with proper lighting and orientation, mimicking real-world anomalies and making them significantly more difficult to detect.

Multi-Class Out-of-Distribution Detection. A recent line of work leverages deep neural representations from multi-class classifiers to perform out-of-distribution (OOD)

detection on high-dimensional data, including images, text, and speech data. Hendrycks and Gimpel [18] formulate the task and propose the simple baseline of using the maximum softmax probability of the classifier on an input to gauge whether the input is out-of-distribution. In particular, the task is formulated as distinguishing between examples from an in-distribution dataset and various out-of-distribution datasets. Importantly, entire images are treated as out-of-distribution.

Continuing this line of work, Lee et al. [24] propose to improve the neural representation of the classifier to better separate in- and out-of-distribution examples. They use generative adversarial networks to produce near-distribution examples. In training, they encourage the classifier to output a uniform posterior on these synthetic out-of-distribution examples. Hendrycks et al. [19] observe that outliers are often easy to obtain in large quantity from diverse, realistic datasets and demonstrate that training out-of-distribution detectors to detect these outliers generalizes to completely new, unseen classes of anomalies. Other work investigates improving the anomaly detectors themselves given a fixed classifier [10, 25]. However, as observed in [19], most of these works tune hyperparameters on a particular type of anomaly that is also seen at test time, so their evaluation setting is more lenient. In this paper, we ensure that all anomalies seen at test time come from entirely unseen categories and are not tuned on in any way, hence we do not compare to techniques such as [25]. Additionally, in a point of departure from prior work, we focus primarily on large-scale images and also datasets with many classes.

Multi-Label Classification. Natural images often contain many objects of interest with complex relationships of co-occurrence. Multi-label image classification acknowledges this more realistic setting by allowing each image to have multiple overlapping labels. This problem has long been of interest [13], and recent web-scale multi-label datasets demonstrate its growing importance, including Tencent ML-Images [32] and Open Images [23]. Prior work addresses multi-label classification in a variety of ways, such as by leveraging label dependencies [31], but none have explored using multi-label classifiers for out-of-distribution detection.

3. The CAOS Benchmark

The Combined Anomalous Object Segmentation (CAOS) benchmark is comprised of two complementary datasets for

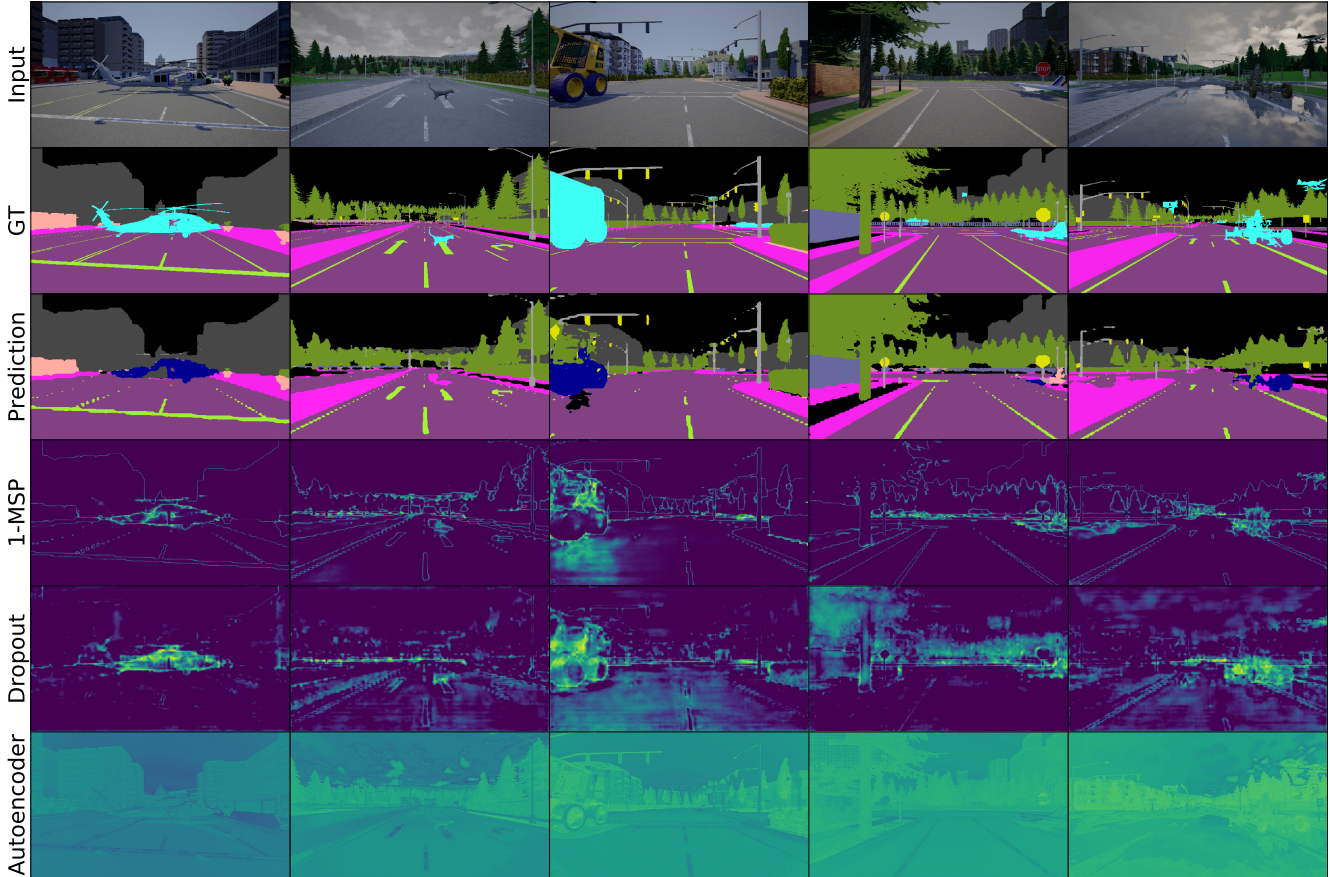


Figure 2: An assortment of conspicuously anomalous scenes, model predictions, and model uncertainties. GT is ground truth, MSP is the maximum softmax probability [18], and Dropout is the model prediction variance under different evaluations with dropout activated during inference [15]. The autoencoder model is based on the spatial autoencoder used in [2].

evaluating anomaly segmentation systems on diverse, realistic anomalies. First is the StreetHazards dataset, which leverages simulation to provide a large variety of anomalous objects realistically inserted into driving scenes. Second is the BDD-Anomaly dataset, which consists of real images. StreetHazards contains a highly diverse array of anomalies; BDD-Anomaly contains anomalies in real-world images. Together, these datasets allow researchers to judge techniques on their ability to segment diverse anomalies as well as anomalies in real images. All images have 720×1280 resolution, and we recommend evaluating with the AUROC, AUPR, and FPRK metrics, which we describe in Section 4.

The StreetHazards Dataset. StreetHazards is an anomaly segmentation dataset that leverages simulation to provide diverse, realistically-inserted anomalous objects. To create the StreetHazards dataset, we use the Unreal Engine along with the CARLA simulation environment [11]. From several months of development and testing including customization of the Unreal Engine and CARLA, we can insert foreign entities into a scene while having them be properly

integrated. Unlike previous work, this avoids the issues of inconsistent chromatic aberration, edge effects, differences in environmental lighting, and other simple cues that an object is anomalous. Additionally, using a simulated environment allows us to dynamically insert diverse anomalous objects in any location and have them render properly with changes to lighting and weather including time of day, cloudy skies, and rain.

We use 3 towns from CARLA for training, from which we collect RGB images and their respective semantic segmentation maps to serve as our training data for our semantic segmentation model. We generate a validation set from the fourth town. Finally, we reserve the fifth and sixth town as our test set. We insert anomalies taken from the Digation Model Bank Library and semantic ShapeNet (ShapeNet-Sem) [29] into the test set in order to evaluate methods for out-of-distribution detection. In total, we use 250 unique anomaly models of diverse types. There are 12 classes used for training: background, road, street lines, traffic signs, sidewalk, pedestrian, vehicle, building, wall, pole, fence, and

			MSP	MSP+CRF	Background	Dropout	AE	Branch
StreetHazards	FPR95	↓	33.7	29.9	69.0	79.4	91.7	68.4
	AUROC	↑	87.7	88.1	58.6	69.9	66.1	65.7
	AUPR	↑	6.6	6.5	4.5	7.5	2.2	1.5
BDD-Anomaly	FPR95	↓	31.9	26.0	78.4	33.4	82.2	78.4
	AUROC	↑	84.2	86.3	31.2	83.7	56.8	61.2
	AUPR	↑	6.3	8.2	2.3	6.5	1.4	2.8

Table 2: Results on the Combined Anomalous Object Segmentation benchmark. AUPR is low for all methods due to the large class imbalance, but most methods perform substantially better than chance. We find that segmentation-specific methods like CRF postprocessing before computing the MSP can improve segmentation of anomalies. All results are percentages.

vegetation. The thirteenth class is the anomaly class that is only used at test time. We collect 5,125 image and semantic segmentation ground truth pairs for training, 1,031 pairs without anomalies for validation, and 1,500 test pairs with anomalies.

The BDD-Anomaly Dataset. BDD-Anomaly is an anomaly segmentation with real images in diverse conditions. We source the BDD-Anomaly dataset from the BDD100K semantic segmentation dataset [34], a large-scale semantic segmentation dataset with diverse driving conditions. The original data consists in 7000 images for training and 1000 for validation. There are 18 original classes. We choose *motorcycle* and *train* as the anomalous object classes and remove all images with these objects from the training and validation sets. This yields 6,688 training pairs, 951 validation pairs without anomalies, and 361 testing pairs with anomalous objects.

4. Experiments

Metrics. To evaluate baseline methods for anomaly segmentation, we use three standard metrics of detection performance: area under the ROC curve (AUROC), false positive rate at 95% recall (FPR95), and area under the precision-recall curve (AUPR). The AUROC and AUPR are important metrics, because they give a holistic measure of performance when the cutoff for detecting anomalies is not a priori obvious or when we want to represent the performance of a detection method across several different cutoffs.

The AUROC can be thought of as the probability that an anomalous example is given a higher score than a ordinary example. Thus, a higher score is better, and an uninformative detector has a AUROC of 50%. AUPR provides a metric more attuned to class imbalances, which is relevant in anomaly and failure detection, when the number of anomalies or failures may be relatively small. Last, the FPR95 metric consists of measuring the false positive rate at 95%. This is important because it tells us how many false positives (i.e. false alarms) are necessary for a given method to achieve a desired recall. This desired recall may be thought of as a safety threshold. Moreover, anomalies and system

failures may require human intervention, so a detector requiring little human intervention while still detecting most anomalies is of pecuniary importance.

Since each pixel is treated as a prediction, there are many predictions to evaluate. To fit the evaluation in memory, we compute the metrics on each image and average over the images to obtain final values.

4.1. Segmentation

Methods. Our first baseline is pixel-wise Maximum Softmax Probability (MSP). Introduced in [18] for multi-class out-of-distribution detection, we directly port this baseline to anomaly segmentation. Alternatively, the background class might serve as an anomaly detector, because it contains everything not in the other classes. To test this hypothesis, “Background” uses the posterior probability of the background class to obtain the anomaly score. The Dropout method leverages MC Dropout [15] to obtain an epistemic uncertainty estimate. Following Kendall et al. [20], we compute the variance of the pixel-wise posteriors over multiple test-time dropout masks and average across all classes, which serves as the anomaly score. We also experiment with an autoencoder baseline similar to [2, 16] where pixel-wise reconstruction loss is used as the anomaly score. This method is called AE. The “Branch” method is a direct port of the confidence branch detector from [10] to pixel-wise prediction. Finally, we experiment with passing the posterior through a fully-connected conditional random field with Gaussian edge potentials followed by the MSP. To our knowledge, this method has not appeared in prior work on anomaly segmentation, and we call it MSP+CRF. We use preset hyperparameters for the CRF.

For all of the baselines except the autoencoder, we train a PSPNet [36] decoder with a ResNet-101 encoder [17] for 20 epochs. We train the encoder and decoder both with SGD with momentum using a learning rate of 2×10^{-2} , momentum of 0.9, and learning rate decay of 10^{-4} . The parameters for batch normalization are frozen from the initial ImageNet pre-training. For AE, we use a 4-layer U-Net [28] with a spatial latent code as in [2]. The U-Net also uses batch

D_{in}	FPR95 ↓				AUROC ↑			
	MaxLogit	LogitAvg	LOF	IForest	MaxLogit	LogitAvg	LOF	IForest
VOC	35.6	98.2	84.0	98.6	90.9	47.9	68.4	46.3
COCO	40.4	94.5	78.4	95.6	90.2	55.5	70.2	41.4

Table 3: Multi-class out-of-distribution detection comparison of the maximum logit, logit average, Local Outlier Factor, and Isolation Forest anomaly detectors on PASCAL VOC and MS-COCO. The same network architecture is used for all three detectors. All results shown are percentages.

norm and is trained for and is trained for 10 epochs.

Results and Analysis. The maximum softmax probability baseline performs the best. The intuitive baseline of using the posterior for the background class to detect anomalies performs worst and surprisingly has a lower AUROC than chance. This suggests that the background class may not align with rare visual features. Even though reconstruction-based scores succeed in product fault segmentation, we find that the AE method performs poorly on the CAOS benchmark, which may be due to the more complex domain. AUPR for all methods is low, indicating that the large class imbalance presents a serious challenge. However, we find that passing the posteriors through a CRF results in some improvements over the MSP baseline. Thus, progress on this task is possible and there is much room for improvement.

In Figure 2, we can qualitatively see that Dropout highlights some anomalous objects more strongly than MSP. However, both methods have a visibly high number of false positives, as they assign a high anomaly score to object edges, a problem also observed in the recent work of Blum et al. [5]. This is because they rely on the posterior probabilities output by the segmentation model, which transition smoothly across semantic boundaries, thereby assigning low confidence to predictions in these regions. Angus [1] proposes a fix to remove these spurious detections but finds that it barely improves quantitative performance. Larger performance gains may come from improved modeling of the support of in-distribution regions.

Autoencoder-based methods are qualitatively different from approaches using the posterior, because they model the input itself and can avoid boundary effects seen in the MSP and Dropout rows of Figure 2. While autoencoder methods are successful in medical anomaly segmentation and product fault detection, we find the AE baseline to be ineffective in the more complex domain of street scenes. The last row of Figure 2 shows pixel-wise reconstruction loss on example images from StreetHazards. Anomalies are not distinguished well from in-distribution elements of the scene. New methods must be developed to mitigate the boundary effects faced by posterior-based methods while also attaining good detection performance. Thus, there is much potential for improvement on the CAOS benchmark.

4.2. Multi-Label Prediction

Current work on out-of-distribution detection is in the multi-class or unsupervised setting. Yet as classifiers learn to classify more objects and process larger images, the multi-label formulation becomes increasingly natural. To our knowledge, this problem setting has yet to be explored. We provide a baseline and evaluation setup.

Methods. For multi-label classification we use the datasets PASCAL VOC [14] and MS-COCO [26]. To evaluate the models trained on these datasets, we use 20 out-of-distribution classes from ImageNet-22K. These classes have no overlap with ImageNet-1K, PASCAL VOC, nor MS-COCO. The 20 classes are chosen not to overlap with ImageNet-1K since the multi-label classifiers models are pre-trained on ImageNet-1K. Specifically, we choose the following classes from ImageNet-22K to serve as out-of-distribution data: dolphin (n02069412), deer (n02431122), bat (n02139199), rhino (n02392434), raccoon (n02508213), octopus (n01970164), giant clam (n01959492), leech (n01937909), Venus flytrap (n12782915), cherry tree (n12641413), Japanese cherry blossoms (n12649317), red wood (n12285512), sunflower (n11978713), croissant (n07691650), stick cinnamon (n07814390), cotton (n12176953), rice (n12126084), sugar cane (n12132956), bamboo (n12147226), and tumeric (n12356395).

For our experiments we use a ResNet-101 backbone architecture that is pre-trained on ImageNet-1K. We replace the final layer of ResNet-101 with 2 fully connected layers and apply the logistic sigmoid function for multi-label prediction. During training we freeze the batch normalization parameters due to an insufficient number of images for proper mean and variance estimation. We train each model for 50 epochs using the Adam optimizer [21] with hyperparameter values 10^{-4} and 10^{-5} for β_1 and β_2 respectively. For data augmentation we use standard resizing, random crops, and random flips to obtain images of size $256 \times 256 \times 3$. As a result of this training procedure, the mAP of the ResNet-101 on PASCAL VOC is 89.11% and 72.0% for MS-COCO.

We evaluate the trained MS-COCO and PASCAL VOC models using four different detectors. We measure the effectiveness of each baseline by computing the false positive rate at five percent recall (FPR95) and the area under the receiver operating characteristic curve (AUROC). Results

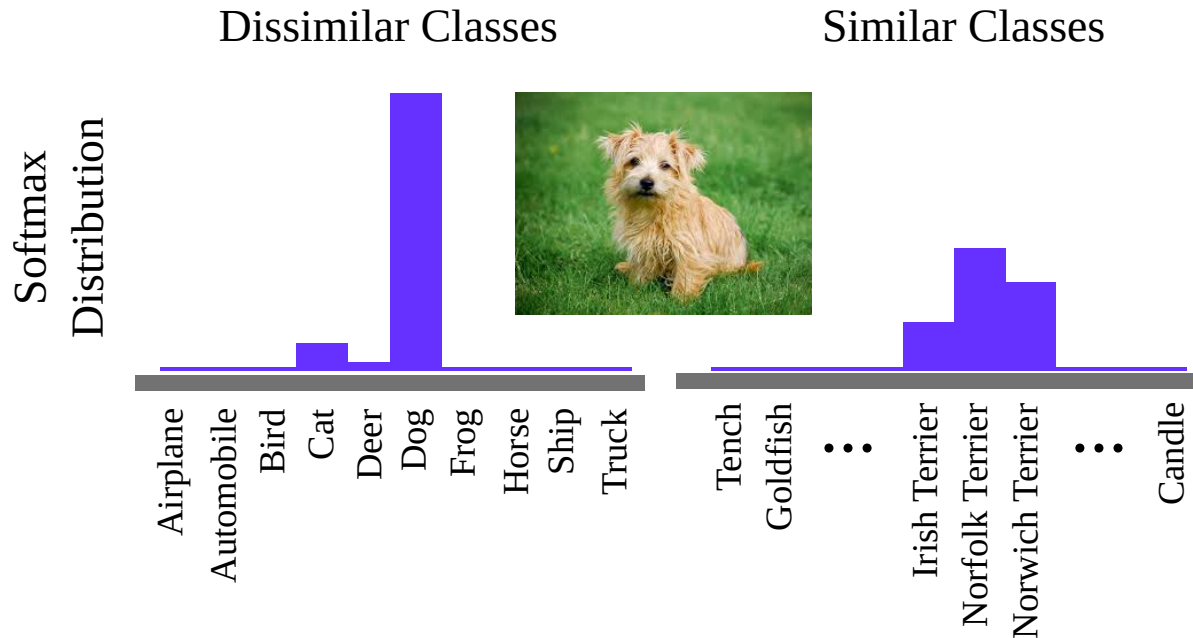


Figure 3: Small-scale datasets such as CIFAR-10 have relatively disjoint classes, but larger-scale datasets including ImageNet have several classes which are highly visually similar to several other classes. The implication is that large-scale classifiers disperse probability mass among several classes. If the prediction confidence is used for out-of-distribution detection, then images which have similarities to other classes will often wrongly be deemed out-of-distribution due to dispersed confidence.

are in Table 3. MaxLogit denotes taking the negative of the maximum value of a logit vector as the anomaly score. Alternatively, one can average the logit values, denoted by LogitAvg. Isolation Forest, denoted by IForest, works by randomly partitioning the space into half spaces to form a decision tree. The score is determined by how close a point is to the root of the tree. Finally we experiment with the local outlier factor (LOF) [6], which computes a local density ratio between every element and its neighbors. Here we set the number of neighbors considered to be 20. These serve as our baseline detectors for multi-label OOD detection.

Results. Results are shown in Table 3. We observe that the MaxLogit method outperforms the average logit and LOF by a significant margin. The MaxLogit method is based on the maximum softmax probability baseline [18], again validating the utility of this simple yet effective baseline. These results establish a baseline and evaluation setup to build on for future work in out-of-distribution detection with multi-label datasets.

4.3. Multi-Class Prediction

Problem. In large-scale image classification, a network is often tasked with predicting an object’s identity from one of hundreds or thousands of classes. As the number of classes increases, class distinctions tend to become more fine and subtle. An increase in similarity and overlap between classes spells a problem for the multi-class out-of-

distribution baseline [18]. The multi-class out-of-distribution detection baseline anomaly score is the negative maximum softmax probability or prediction confidence of model p on input x , $-\max_k p(y = k | x)$. Classifiers tend to have higher confidence on in-distribution examples than out-of-distribution examples, enabling OOD detection. Assuming single-model evaluation and no access to other anomalies or test-time adaptation, the maximum softmax probability (MSP) is the state-of-the-art multi-class out-of-distribution detection method. However, we show that the MSP is problematic for large-scale datasets including ImageNet-1K and Places365 [37]. These datasets contain some classes which are visually similar, so that the MSP is less suitable for OOD detection. For a sketch of why, consider Figure 3. Visually similar classes have their probability mass dispersed among several classes. Consequently, a classifier may produce a low confident prediction for an in-distribution image not because the image is unfamiliar or out-of-distribution but because the object’s exact class is difficult to determine. In turn, we need a different out-of-distribution detection baseline for large-scale multi-class systems.

Method. Some classes tend to be predicted with low confidence and others with high confidence. The shape of posterior is often class dependent. We depart from the maximum softmax probability and no longer implicitly assume the posterior entropies are similar. We capture the typical shape of each class’s posterior distribution and form posterior distribu-

\mathcal{D}_{in}	\mathcal{D}_{out}^{test}	FPR95 ↓		AUROC ↑		AUPR ↑	
		MSP	KL (Ours)	MSP	KL (Ours)	MSP	KL (Ours)
ImageNet	Gaussian	1.7	3.6	99.5	98.0	96.1	78.8
	Rademacher	21.7	14.6	88.8	92.7	43.5	54.2
	Blobs	26.8	7.5	87.4	98.5	41.2	93.4
	Textures	68.6	58.9	80.1	85.3	36.5	48.2
	LSUN	66.8	59.9	74.6	78.5	33.2	37.9
	Places365	68.9	72.0	77.3	79.4	39.1	45.8
	Mean	42.42	36.06	84.60	88.73	48.26	59.70
Places365	Gaussian	10.2	11.9	93.12	92.7	16.34	15.6
	Rademacher	19.8	0.5	89.02	99.8	11.10	88.3
	Blobs	58.9	27.0	71.50	92.9	5.11	30.7
	Textures	85.7	74.1	64.51	79.2	3.58	12.4
	Places69	88.8	91.0	60.18	64.6	4.50	5.7
	Mean	52.68	40.88	75.67	85.82	8.13	30.52

Table 4: Multi-class out-of-distribution detection results using the maximum softmax probability and our KL matching anomaly score. Results are on ImageNet and Places365. Values are rounded so that 99.95% rounds to 100%.

tion templates for each class. During test time, the network’s softmax posterior distribution is compared to these templates and an anomaly score is generated.

More concretely, we compute k different distributions d_k , one for each class. We write $d_k = \mathbb{E}_{x' \sim \mathcal{X}_{val}^k} [p(y|x')]$ where \mathcal{X}_{val}^k is the set of validation examples with label k . Then for a new test input x , we calculate the anomaly score $\min_k \text{KL}[p(y|x) \| d_k]$ rather than the MSP baseline $-\max_k p(y = k | x)$. Note that we utilize the validation dataset, but our KL matching method does not even require the validation dataset’s labels. Consequently, if a new Norfolk Terrier input has a low-confidence softmax posterior distribution $p(y|x)$ resembling the typical response $p_{\text{Norfolk Terrier}}$, it will then rightly have a low anomaly score.

Datasets. To evaluate the MSP baseline out-of-distribution detection and our KL matching detector, we use ImageNet-1K, an object recognition dataset, and Places365, a scene recognition dataset, as in-distribution datasets \mathcal{D}_{in} and several out-of-distribution test datasets \mathcal{D}_{out} . The first out-of-distribution dataset is *Gaussian* noise, where each pixel of these out-of-distribution examples are i.i.d. sampled from $\mathcal{N}(0, 0.5)$ and clipped to be contained within $[-1, 1]$. Another type of test-time noise is *Rademacher* noise, in which each pixel is i.i.d. sampled from $2 \cdot \text{Bernoulli}(0.5) - 1$, or each pixel is 1 or -1 with equal probability. *Blob* examples are more structured than noise; they are algorithmically generated blob images. Meanwhile *Textures* is a dataset consisting in images of describable textures [7]. When evaluating the ImageNet-1K detector, we use *LSUN* images which is a dataset for scene recognition [33]. Our final dataset is another scene recognition dataset. *Places69* are scene classes which are separate from the scene classes of Places365. In all, we evaluate against out-of-distribution examples span-

ning synthetic and realistic images.

Results. Results are shown in Table 4. Observe that our KL matching technique outperforms the maximum softmax probability baseline for all three metrics on both ImageNet and Places365. In the case of Places365, the AUROC improvement is over 10%. These results were computed using a ResNet-50 trained on either ImageNet-1K or Places365. In consequence, we suggest a new baseline for large-scale multi-class out-of-distribution detection.

5. Conclusion

In this paper, we investigated expanding out-of-distribution detection to more realistic, large-scale settings, which led us to develop a novel benchmark for anomaly segmentation. We introduced the CAOS benchmark for anomaly segmentation consisting of diverse, naturally-integrated anomalous objects in driving scenes. Baseline methods on the CAOS benchmark substantially improve on random guessing but are still lacking. We observe that incorporating segmentation-specific techniques like Gaussian CRFs further improves performance, indicating that progress is possible on this task, and that there is wide room for improvement. We also investigated using multi-label classifiers for out-of-distribution detection and established a baseline and experimental setup for this previously unexplored setting. Finally, on large-scale multi-class image datasets, we proposed an improved out-of-distribution detector that provides consistent and significant gains. In all, we hope that our new baselines and our new anomaly segmentation benchmark will enable further research on out-of-distribution detection in real-world safety-critical environments, a necessity for safe autonomous vehicles.

References

- [1] M. Angus. Towards pixel-level ood detection for semantic segmentation, 2019.
- [2] C. Baur, B. Wiestler, S. Albarqouni, and N. Navab. Deep autoencoding models for unsupervised anomaly segmentation in brain mr images. *Lecture Notes in Computer Science*, page 161–169, 2019. ISSN 1611-3349. doi: 10.1007/978-3-030-11723-8_16.
- [3] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger. Mvtec ad - a comprehensive real-world dataset for unsupervised anomaly detection. In *CVPR*, 2019.
- [4] P. Bevandić, I. Krešo, M. Oršić, and S. Šegvić. Discriminative out-of-distribution detection for semantic segmentation, 2018.
- [5] H. Blum, P.-E. Sarlin, J. Nieto, R. Siegwart, and C. Cadena. The fishyscapes benchmark: Measuring blind spots in semantic segmentation, 2019.
- [6] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. Lof: identifying density-based local outliers. In *ACM sigmod record*, volume 29, pages 93–104. ACM, 2000.
- [7] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. In *Computer Vision and Pattern Recognition*, 2014.
- [8] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [10] T. DeVries and G. W. Taylor. Learning confidence for out-of-distribution detection in neural networks. *arXiv preprint arXiv:1802.04865*, 2018.
- [11] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun. CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 1–16, 2017.
- [12] A. Emmott, S. Das, T. Dietterich, A. Fern, and W.-K. Wong. A meta-analysis of the anomaly detection problem, 2015.
- [13] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>.
- [14] M. Everingham, L. V. Gool, C. K. I. Williams, J. M. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88: 303–338, 2009.
- [15] Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. *International Conference on Machine Learning*, 2016.
- [16] M. Haselmann, D. P. Gruber, and P. Tabatabai. Anomaly detection using deep learning based image completion. In *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 1237–1242. IEEE, 2018.
- [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CVPR*, 2015.
- [18] D. Hendrycks and K. Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *ICLR*, 2017.
- [19] D. Hendrycks, M. Mazeika, and T. Dietterich. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*, 2019.
- [20] A. Kendall, V. Badrinarayanan, and R. Cipolla. Bayesian segnet: Model uncertainty in deep convolutional encoder-decoder architectures for scene understanding. *ArXiv*, abs/1511.02680, 2015.
- [21] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *ICLR*, 2014.
- [22] I. Krešo, M. Oršić, P. Bevandić, and S. Šegvić. Robust semantic segmentation with ladder-densenet models, 2018.
- [23] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Mallocci, T. Duerig, and V. Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *arXiv:1811.00982*, 2018.
- [24] K. Lee, H. Lee, K. Lee, and J. Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. *International Conference on Learning Representations*, 2018.
- [25] S. Liang, Y. Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *International Conference on Learning Representations*, 2018.
- [26] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollar. Microsoft COCO: Common objects in context. *ECCV*, 2014.
- [27] P. Pinggera, S. Ramos, S. K. Gehrig, U. Franke, C. Rother, and R. Mester. Lost and found: Detecting small road hazards for self-driving vehicles. *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1099–1106, 2016.
- [28] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, page 234–241, 2015. ISSN 1611-3349. doi: 10.1007/978-3-319-24574-4_28.
- [29] M. Savva, A. X. Chang, and P. Hanrahan. Semantically-Enriched 3D Models for Common-sense Knowledge. *CVPR 2015 Workshop on Functionality, Physics, Intentionality and Causality*, 2015.
- [30] B. Schölkopf, R. Williamson, A. Smola, J. Shawe-Taylor, and J. Platt. Support vector method for novelty detection. In *Proceedings of the 12th International Conference on Neural Information Processing Systems, NIPS’99*, pages 582–588, Cambridge, MA, USA, 1999. MIT Press.
- [31] J. Wang, Y. Yang, J. Mao, Z. Huang, C. Huang, and W. Xu. Cnn-rnn: A unified framework for multi-label image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2285–2294, 2016.
- [32] B. Wu, W. Chen, Y. Fan, Y. Zhang, J. Hou, J. Huang, W. Liu, and T. Zhang. Tencent ml-images: A large-scale multi-label image database for visual representation learning. *arXiv preprint arXiv:1901.01703*, 2019.
- [33] F. Yu, Y. Zhang, S. Song, A. Seff, and J. Xiao. LSUN: construction of a large-scale image dataset using deep learning with humans in the loop. *CoRR*, 2015.
- [34] F. Yu, W. Xian, Y. Chen, F. Liu, M. Liao, V. Madhavan, and

- T. Darrell. BDD100K: A diverse driving video database with scalable annotation tooling. *CoRR*, abs/1805.04687, 2018.
- [35] O. Zendel, K. Honauer, M. Murschitz, D. Steininger, and G. Fernandez Dominguez. Wilddash-creating hazard-aware benchmarks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 402–416, 2018.
- [36] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6230–6239, 2017.
- [37] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba. Places: A 10 million image database for scene recognition. *PAMI*, 2017.
- [38] P. Zhou, X. Han, V. I. Morariu, and L. S. Davis. Learning rich features for image manipulation detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1053–1061, 2018.