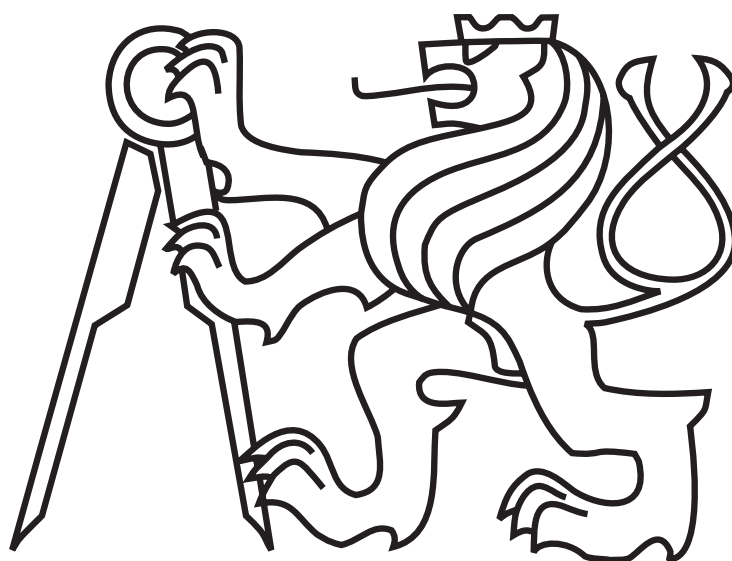


CZECH TECHNICAL UNIVERSITY IN PRAGUE

Faculty of Electrical Engineering

SEMESTRAL PROJECT



Pavel Linder

Detecting Out-Of-Distribution Samples with Object Detectors

Department of Computer Science

Supervisor: Ing. Bc. Radim Špetlík

May, 2023

Abstract

Contents

1	Introduction	1
2	Background	2
3	OOD detection in multi-class setting	3
3.1	MSP detector	3
3.2	Out-of-Distribution detector for Neural networks (ODIN)	4
3.3	Mahalanobis distance detector	4
3.4	Energy-based detector	4
3.5	Fast Out-Of-Distribution Detector (FOOD)	5
3.6	GradNorm	5
4	OOD detection in multi-label setting	6
4.1	JointEnergy detector	6
4.2	YolOOD	6
5	Datasets	7
6	Experiments	9
6.1	Experimental setup	9
6.2	OOD detection	10
6.3	Evaluation settings	10
6.4	Results	11
7	Conclusion	12
8	Future Work	12

1 Introduction

Artificial Intelligence (AI) has become a topic of great interest among the general public especially in recent times. Machine Learning (ML) models are being widely adopted across various domains to handle a wide range of tasks, and novel applications are being discovered every day. When deploying machine learning models in real-world scenarios, our primary concern is typically centered around the ultimate precision of the predictive outcomes. However, it is equally crucial to take into account the reliability and validity of these predictions. One must assess whether the model's high response is a result of its exposure to comparable data during training or if it is yielding unreliable predictions for unexplored data that was not previously encountered during the training process.

Modern deep learning models can easily produce these overconfident predictions. This issue not only decreases a model's robustness but also raises significant concerns in areas such as medical care, where incorrect diagnoses can result in severe outcomes. Further, it can also question the safety in AI [1]. A new area of research called Out-Of-Distribution (OOD) detection aims to get rid of this vulnerability by determining whether an input is in-distribution (ID) or OOD. [2] [3] [4] [5] [6] [7] By identifying OOD samples, models can decrease the risk of inaccurate predictions, speed up human intervention when necessary, and establish a dependable and secure incorporation of machine learning technologies across a range of domains.

Over the years, extensive research has been carried out on multi-class classification; however, the multi-label task remains an area that has been largely underexplored. The goal of this work is to provide a comprehensive analysis of the different techniques and trends utilized in multi-label OOD detectors, while also categorizing and discussing the various methods employed. Additionally, we offer a concise overview of the multi-class context to show the fundamental concepts.

2 Background

Multi-label classification Multi-label classification is an instance of supervised learning problem where each input sample can be associated with multiple labels simultaneously. This approach differs from the traditional multi-class classification, which assigns only a single label to each sample. Multi-label classification allows the presence of multiple or even zero labels for a given input. This makes it a more adaptable and comprehensive approach, making it suitable for several real-world scenarios.

Formally, let \mathcal{X} (respectively, \mathcal{Y}) denote the input (respectively, output) space of classifier $f : \mathcal{X} \rightarrow \mathbb{R}^{|\mathcal{Y}|}$ trained on samples drawn from distribution \mathcal{P} , where \mathcal{P} is a distribution over $\mathcal{X} \times \mathcal{Y}$. Each input sample $x \in \mathcal{X}$ is associated with a set of labels $Y = 1, 2, \dots, K$ represented as a binary vector $y = [y_1, y_2, \dots, y_K]$, where $y_i = 1$ if input sample x is associated with class i , and 0 otherwise.

To illustrate, consider the example of image classification. In a multi-label setting, an image can have multiple objects or attributes of interest. For instance, an image may contain a cat, a dog, and a tree simultaneously. Instead of assigning a single label to the image, a multi-label classifier would predict a binary vector where the elements corresponding to "cat," "dog," and "tree" are all set to 1.

For OOD detection we will utilize neural networks with a shared feature space to obtain a multi-label output prediction. As opposite to the approach of training disjoint classifiers as proposed in literature [8], implementing the end-to-end training method with a shared feature space is more computationally efficient than training K completely independent models. Multi-label classification models are now mostly trained using this technique, which has been shown to work effectively in various domains [9] [10] [11].

Multi-class classification The difference between multi-label and multi-class classification [12] is determined by the flexibility and granularity of the label assignments. Multi-label classification acknowledges the possibility of multiple labels being equivalently relevant simultaneously, accommodating more complex and diverse scenarios.

Object detection The goal of an object detector extends the task of classification and except classifying the objects present in an image it also provides their precise spatial localization by drawing bounding boxes around them. This allows for accurate identification and tracking of objects of interest. There are several well-known object detectors such as SSD [13], YOLO [14], R-CNN [15] and Mask R-CNN [16].

Out-Of-Distribution detection Out-Of-Distribution (OOD) detection can be formulated as an instance of binary classification. Let \mathcal{D}_{in} denote the marginal distribution of \mathcal{P} over \mathcal{X} which represents the distribution of in-distribution (ID) data. In practise, OOD is

frequently characterized by a distribution that emulates the uncertainties that arise during deployment, such as samples from an irrelevant distribution whose label set does not intersect with \mathcal{Y} and thus should not be predicted by the model. We will denote this distribution over \mathcal{X} as \mathbb{D}_{out} . For OOD detector G in multi-label setting with classifier f we define a decision function for the binary classification as follows:

$$G(\mathbf{x}, f) = \begin{cases} 0 & \text{if } \mathbf{x} \sim \mathbb{D}_{out} \\ 1 & \text{if } \mathbf{x} \sim \mathbb{D}_{in} \end{cases} \quad (1)$$

Energy function An energy function is a mathematical function that maps a scalar value to a given input. In the context of the paper, the energy function is used to represent the cost or likelihood of a particular label. Formally, the energy function $E(x) : \mathcal{X} \rightarrow \mathbb{R}$ maps each point x of an input space to a scalar value called the energy. It assigns a high energy score to OOD samples and a low energy score to in-distribution samples.

Kullback-Leibler (KL) divergence Kullback-Leibler (KL) divergence [17] is a measure of the difference between two probability distributions. KL divergence is defined as the expectation of the logarithmic difference between the model-predicted distribution $q = q_i$ and the reference distribution $p = p_i$. Formally:

$$D_{KL} = (p \parallel q) = \sum_i p_i \log \frac{p_i}{q_i} = \sum_i p_i \log q_i + \sum_i p_i \log p_i = H(p, q) - H(p) \quad (2)$$

3 OOD detection in multi-class setting

In 2015, Nguyen et al. [18] exploit the vulnerability of deep networks. They found that one cause of overconfidence for OOD data is due to the nature of the fast-growing exponential function used in computing softmax probabilities. Small changes to the logits can result in significant changes to the output distribution. The direct correlation between prediction probability and confidence in a softmax distribution is poor.

3.1 MSP detector

Maximum Softmax Probability (MSP) score In 2016, Hendrycks and colleagues [2] took a significant step towards addressing the vulnerabilities of neural networks. The main findings can be concluded as follows. The probability of correctly classified examples is higher than that of misclassified and out-of-distribution examples. Thus, by capturing the statistics of prediction probabilities of correct examples, we can use them to detect whether an example is in error or abnormal. It is important to note that relying solely on softmax probabilities may lead to unreliable outcomes. However, analyzing the statistics of

these probabilities is yet effective across various domains, including computer vision. This work has become a baseline for OOD detection in multi-class settings and is still used for benchmarking state-of-the-art methods.

3.2 Out-of-Distribution detector for Neural networks (ODIN)

Temperature scaling and input pre-processing In 2017, Liang et al. presented another output-based approach known as ODIN [19] (Out-of-Distribution detector for Neural networks). ODIN is a technique that uses temperature scaling and input pre-processing (adding small perturbations). It aims to effectively separate the softmax probability distributions of in-and-out-of-distribution images. This separation enables more accurate and efficient OOD detection. Temperature scaling involves scaling the logits (inputs to the softmax function) by a temperature parameter T before computing the softmax function. This has the effect of sharpening the softmax output, making it easier to distinguish between in-distribution and out-of-distribution images.

3.3 Mahalanobis distance detector

Mahalanobis distance In 2018, Lee et al. [20] introduced a new feature-based approach for detecting out-of-distribution (OOD) samples, which can be applied to any pre-trained softmax neural classifier without requiring to re-train. The method uses a ‘generative’ (distance-based) classifier to measure the probability density of test samples on feature spaces of Deep Neural Networks (DNNs). The authors assume that pre-trained low and upper level features can be fitted well by a class-conditional Gaussian distribution. Confidence score is defined using the Mahalanobis distance with respect to the closest class-conditional distribution.

3.4 Energy-based detector

Energy score In 2021, Liu, Wang, et al. [21] presented a new method that uses outputs of the DNNs. Instead of employing softmax scores, as with the MSP approach, they employ energy scores, which better distinguish between in- and out-of-distribution samples. Unlike softmax confidence scores, energy scores are theoretically aligned with the probability density of the inputs and are less susceptible to the overconfidence issue. Energy can be flexibly used as a scoring function for any pre-trained neural classifier, as well as a trainable cost function to explicitly shape the energy surface for OOD detection. Energy score function was introduced in section 2 and for a given input (\mathbf{x}, y) is defined as $E(\mathbf{x}, y) = -f_y(\mathbf{x})$. Without altering the parameterization of the neural network $f(\mathbf{x})$, we can define the **free energy function** $E(\mathbf{x}; f)$ over $\mathbf{x} \in \mathbb{R}^D$ in terms of the denominator of the softmax activation:

$$E(\mathbf{x}, f) = -T * \log \sum_i^K e^{f_i(\mathbf{x})/T} \quad (3)$$

where K is the number of classes and T is the temperature parameter.

3.5 Fast Out-Of-Distribution Detector (FOOD)

Statistical testing and additional output neuron In 2021, Amit and Levy [22] proposed a Fast Out-Of-Distribution Detector (FOOD) that efficiently detects out-of-distribution (OOD) samples with minimal inference time overhead. The paper presents a DNN model with a final Gaussian layer that models a density function for each class and a rapid OOD detector that does not require OOD samples for training or hyperparameter tuning. The paper uses a log likelihood ratio statistical test and an additional output neuron for OOD detection.

3.6 GradNorm

Categorizing OOD detectors We can categorize the OOD detection methods that have been implemented into two main groups: **output-based** and **feature-based**. Output-based detectors employ the output of neural networks which is then computed using an aggregation function (such as max or sum) and a scoring function to detect OOD instances. Examples of these methods include ODIN[19], Energy-based[21], and MSP[2]. On the other hand, feature-based detectors utilize the feature space to distinguish between OOD and ID samples. Mahalanobis distance based detection[23] is an example of this method.

Using KL divergence In 2021, Huang, Rui and Geng [24] introduced a novel approach for OOD detection called GradNorm. Their method utilizes information extracted from the **gradient space**. The main idea proposed in this paper is to use the vector norm of gradients, which are backpropagated from the **KL divergence** between the softmax output and a uniform probability distribution, as a means of detecting OOD inputs.

Definition for KL divergence can be find in Equation 2. For this method we set the reference distribution to be uniform $\mathbf{u} = [1/C, 1/C, \dots, 1/C] \in \mathbb{R}^C$. The predictive probability distribution is represented by the softmax output. KL divergence then becomes:

$$D_{KL} = (\mathbf{u} \parallel \text{softmax}(f(\mathbf{x}))) = \frac{1}{C} \sum_{c=1}^K \log \frac{e^{f_c(\mathbf{x})/T}}{\sum_{j=1}^C e^{f_j(\mathbf{x})/T}} - H(\mathbf{u}) \quad (4)$$

where the first term is the cross-entropy loss with a uniform vector \mathbf{u} , and the second term is a constant $H(\mathbf{u})$. KL divergence measures how far the predictive distribution is from the uniform distribution. ID data is expected to have a larger KL divergence because the

prediction tends to concentrate on one of the ground-truth classes and is therefore less uniformly distributed.

GradNorm score This technique does not directly employ KL divergence, but instead utilizes the gradient vector norm, which is propagated backwards from the KL divergence. The OOD score is then defined as:

$$S(\mathbf{x}) = \left\| \frac{\partial D_{KL}(\mathbf{u} \parallel \text{softmax}(f(\mathbf{x})))}{\partial \mathbf{w}} \right\| \quad (5)$$

where \mathbf{w} is the set of parameters in vector form and $\|\cdot\|$ denotes L_1 norm. . .

At the end exploring gradient space proved to be useful and this become a new state-of-the-art method.

4 OOD detection in multi-label setting

As mentioned previously, research on detecting Out-Of-Distribution in multi-label classifiers has quite recently just begun, resulting in a limited number of approaches. A few of the used methods adopt concepts and notions from multi-class OOD detection. In this section, we will provide a more detailed overview description of state-of-the-art OOD detectors in multi-label setting.

4.1 JointEnergy detector

using energy function with summation

4.2 YolOOD

Using object detector

5 Datasets

MS-COCO [25] The MS-COCO dataset comprises images of everyday scenes with common objects to advance object recognition within the context of scene understanding. There are 82,783 training, 40,504 validation, and 40,775 testing images, all of which contain 80 common object categories. This dataset is commonly used for state-of-the-art benchmarking in object recognition and classification including multi-label classification.



Figure 1: Examples of images from dataset MS-COCO and Pascal-VOC. MS-COCO: Picture containing objects from categories ‘person’, ‘tie’, ‘bowl’, ‘chair’, ‘dining table’ and ‘clock’.

VOC-Pascal [26] It is good for multi-label classification because it includes a variety of object categories, allowing for the evaluation of algorithms that can recognize and detect multiple objects in an image. PASCAL-VOC has 22,531 images, which are divided into 20 classes.

ImageNet [27] ImageNet is a large-scale hierarchical image database built upon the WordNet structure. Compared to small image datasets like Caltech101/256, MSRC, and PASCAL, ImageNet is much larger in scale and diversity, offering 20× the number of categories and 100× the number of total images. ImageNet is often use not only for benchmarking for the variety of image but also for pre-trained models.



Figure 2: Examples of images with label annotations from dataset TACO and Textures.

TACO [28] The TACO dataset is designed for litter detection and segmentation with 1500 high-resolution images and 4784 annotations. Litter detection is a challenging problem for multiple reasons so this dataset can be use to show the robustness of classifiers.

6 Experiments

In this section, we evaluate some of the existing state-of-the-art methods for multi-label OOD detection on public datasets in order to assess how they compare on identical settings. To compare state-of-the-art methods, we will be using pre-trained networks with identical weights and experimentation methodology.

To provide a broad overview of the methods we select state-of-the-art method from each respective family. **JointEnergy** as an instance of output-based, **Mahalanobis** for feature-based and **GradNorm** for gradient-based detectors. All three methods were implemented to have the same evaluation pipeline.

The evaluation was performed on datasets that have been utilized across multiple papers about multi-label OOD detection. Consequently, two key outcomes will be obtained:

- new comparison of the methods because we used different datasets
- evaluation of GradNorm method in multi-label settings which has not been done yet

For the purpose of evaluation, datasets utilized in both of these works were utilized. The labels of the ID datasets are different from the labels of the OOD datasets in this evaluation to reduce bias.

6.1 Experimental setup

ID datasets We will use two multi-label datasets: MS-COCO[25] and PASCAL-VOC[26]. These datasets are used to train multi-label classifier and for producing ID reliability scores.

OOD datasets To assess the models trained on the ID datasets, we adopt the same approach as described in [3]. We select a subset of ImageNet[27] dataset, the whole Textures[29] dataset and the whole TACO[28] dataset. For ImageNet, we use the same set of 20 classes as in [3] to have different labels than ID datasets. Textures dataset was used both in [3] and [24]. Finally, TACO dataset was used in [5] and other papers. By this selection we should cover different class domains and provide new insights.

Classification model training We deploy two classifiers sourced from the study by Wang et al.[3] featuring DenseNet-121 architecture. These classifiers were initially trained on ImageNet-1K and subsequently fine-tuned through the utilization of sigmoid function. Data was augmented with random crops and flips. Achieved mAP is 87.51% for PASCAL-VOC, 73.83% for MS-COCO.

6.2 OOD detection

JointEnergy, Mahalanobis We use the public implementation of original authors to produce OOD scores.

GradNorm Since GradNorm has yet to be adopted for multi-label classifiers, it is necessary to adjust the original implementation. For our implementation, we have set the parameters in the same fashion as the authors of GradNorm did, with temperature T being equal to 1 and using the gradients weights of the last layer only. Although binary cross-entropy loss with logits is used for the multi-label classifier, it is still noteworthy that we compute cross-entropy loss here. The multi-label classifier does not utilize softmax on the network’s output, but rather employs the sigmoid function to generate a K -long vector of class probabilities. This vector is exactly what we need to calculate GradNorm, and hence, we can utilize it in a way displayed below.

To generate GradNorm scores, one must execute the following steps:

- use the classifier to produce the predicted probability distribution over classes $[K \times 1]$ vector
- divide the probability distribution by temperature T
- compute gradients of KL divergence as average the derivative of the cross-entropy loss for all labels
- perform backward pass of the network
- get the last layer’s output and compute the first norm

6.3 Evaluation settings

Evaluation pipeline For every classifier, scores are computed from In-Distribution (ID) datasets as **inScores** and scores for Out-Of-Distribution (OOD) datasets as **outScores**. OOD detection is a binary classification so it requires the generation of ground truth labels. In the case of *inScores*, the ground truth is set to 0 since we know that they contain labels which they were trained for. Alternatively, for *outScores*, the ground truth is set to 1. Metrics are then computed from both *inScores* and *outScores* along with their corresponding ground truth labels.

Metrics The evaluation of performance is executed by means of metrics that are widely utilized in the field of Out-of-Distribution (OOD) detection. These metrics include:

- (a) FPR95 - the false positive rate of OOD samples when the true positive rate is at 95%
- (b) AuROC - the area beneath the receiver operating characteristic curve
- (c) AuPR - the area beneath the precision-recall curve

OOD score	detector family	MS-COCO			PASCAL		
		FPR95		AUROC	AUPR		
		↓		↑	↑		
JointEnergy	output-based	1	2	3	4	56	6
Mahalanobis	feature-based	1	1	1	1	1	1
GradNorm	gradient-based	1	1	1	1	1	1

Table 1: Quantitative results for 3 Out-Of-Distribution (OOD) detector: JointEnergy, Mahalanobis and GradNorm. We want the minimal FPR95 and maximal AUROC and AUPR metrics.

6.4 Results

In table bla you can see the quantitative results in Table 1.

7 Conclusion

8 Future Work

References

- [1] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in AI safety. June 2016.
- [2] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. October 2016.
- [3] Haoran Wang, Weitang Liu, Alex Bocchieri, and Yixuan Li. Can multi-label classification networks know what they don’t know? September 2021.
- [4] Sina Mohseni, Arash Vahdat, and Jay Yadawa. Shifting transformation learning for out-of-distribution detection. June 2021.
- [5] Alon Zolfi, Guy Amit, Amit Baras, Satoru Koda, Ikuya Morikawa, Yuval Elovici, and Asaf Shabtai. YolOOD: Utilizing object detection concepts for out-of-distribution detection. December 2022.
- [6] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. October 2021.
- [7] Rui Huang and Yixuan Li. MOS: Towards scaling out-of-distribution detection for large semantic space. May 2021.
- [8] Grigorios Tsoumakas and Ioannis Katakis. Multi-label classification. *International Journal of Data Warehousing and Mining*, 3(3):1–13, July 2007.
- [9] Wenjie Zhang, Junchi Yan, Xiangfeng Wang, and Hongyuan Zha. Deep extreme multi-label learning. April 2017.
- [10] Konstantinos Trohidis, Grigorios Tsoumakas, George Kalliris, and Ioannis Vlahavas. Multi-label classification of music by emotion. *EURASIP Journal on Audio, Speech, and Music Processing*, 2011(1):4, December 2011.
- [11] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M. Summers. ChestX-ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. May 2017.
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, June 2021.
- [13] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. SSD: Single Shot MultiBox Detector. volume 9905, pages 21–37. 2016.

REFERENCES

- [14] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. YOLOv4: Optimal Speed and Accuracy of Object Detection, April 2020.
- [15] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation, October 2014.
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN, January 2018.
- [17] Richard Dykstra. Kullback–Leibler Information. July 2005.
- [18] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images, April 2015.
- [19] Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *6th International Conference on Learning Representations, ICLR 2018 - Conference Track Proceedings*, June 2017.
- [20] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks, October 2018.
- [21] Weitang Liu, Xiaoyun Wang, John D. Owens, and Yixuan Li. Energy-based Out-of-distribution Detection, April 2021.
- [22] Guy Amit, Moshe Levy, Ishai Rosenberg, Asaf Shabtai, and Yuval Elovici. FOOD: Fast out-of-distribution detector. August 2020.
- [23] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. July 2018.
- [24] Rui Huang, Andrew Geng, and Yixuan Li. On the Importance of Gradients for Detecting Distributional Shifts in the Wild, October 2021.
- [25] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft COCO: Common Objects in Context, February 2015.
- [26] Mark Everingham, S. M. Ali Eslami, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The Pascal Visual Object Classes Challenge: A Retrospective. *International Journal of Computer Vision*, 111(1):98–136, January 2015.
- [27] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, June 2009.
- [28] Pedro F. Proença and Pedro Simões. TACO: Trash Annotations in Context for Litter Detection, March 2020.

REFERENCES

- [29] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing Textures in the Wild, November 2013.