

YoOOD: Utilizing Object Detection Concepts for Multi-Label Out-of-Distribution Detection

Alon Zolfi¹, Guy Amit¹, Amit Baras¹, Satoru Koda², Ikuya Morikawa², Yuval Elovici¹, Asaf Shabtai¹

¹Ben-Gurion University of the Negev, Israel

{zolfi, guy5, barasa}@post.bgu.ac.il, {elovici, shabtaia}@bgu.ac.il

²Fujitsu Limited, Japan

{koda.satoru, morikawa.ikuya}@fujitsu.com

Abstract

Out-of-distribution (OOD) detection has attracted a large amount of attention from the machine learning research community in recent years due to its importance in deployed systems. Most of the previous studies focused on the detection of OOD samples in the multi-class classification task. However, OOD detection in the multi-label classification task, a more common real-world use case, remains an underexplored domain. In this research, we propose YoOOD – a method that utilizes concepts from the object detection domain to perform OOD detection in the multi-label classification task. Object detection models have an inherent ability to distinguish between objects of interest (in-distribution data) and irrelevant objects (OOD data) in images that contain multiple objects belonging to different class categories. These abilities allow us to convert a regular object detection model into an image classifier with inherent OOD detection capabilities with just minor changes. We compare our approach to state-of-the-art OOD detection methods and demonstrate YoOOD’s ability to outperform these methods on a comprehensive suite of in-distribution and OOD benchmark datasets.

1. Introduction

Machine learning and particularly deep learning-based networks have become a state-of-the-art solution for computer vision tasks, such as image classification [8, 16], object detection [27, 29], and image segmentation [3, 4]. However, it has been shown that these models can produce overconfident predictions on samples that are not within the distribution they were trained on, *i.e.*, out-of-distribution (OOD) samples [26].

In the last few years, many solutions have been proposed to address this problem, most of which focus on the separation of in-distribution and OOD data [1, 13, 17, 18]. How-

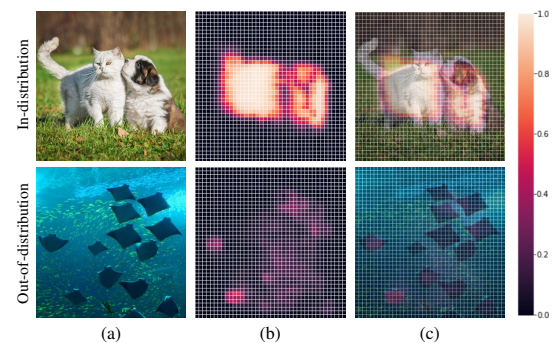


Figure 1. Examples of (a) in-distribution (top) and OOD (bottom) images with the (b) corresponding YoOOD confidence heatmap, and (c) heatmap placed on top of the image.

ever, these studies only proposed solutions for the multi-class classification task, in which an input is associated with a single class category. The problem of OOD detection in the *multi-label* classification domain has been overlooked and remains underexplored. Despite its importance, only two studies specifically addressed this problem [11, 34].

Multi-label image classification and object detection are two closely related tasks. The former involves assigning multiple labels to an image, while the latter goes a step further by not only recognizing the objects present in an image but also localizing them with bounding boxes. Therefore, object detectors have the inherent ability to distinguish between objects of interest and irrelevant objects [22, 27, 29]. This ability, along with object detection’s similarity to multi-label classification, can be leveraged to create an OOD detection mechanism for the multi-label setting.

In this paper, we propose *YoOOD*, a multi-label classifier that utilizes the main concepts of state-of-the-art object detectors, and specifically, the YOLO object detector [2, 14, 27]. To convert YOLO’s network into an image classifier, we simply replace the last layer of each detection

head with a simplified one. *YoLOOD* is based on the objectness score concept, which is commonly used in object detectors [28, 29] to understand the relevance of different parts of an image. This concept is realized by training the network to predict low scores for background areas or areas that contain irrelevant objects. In terms of OOD detection, this can be interpreted as assigning low scores to OOD data and high scores to in-distribution data, as shown in Figure 1.

As opposed to existing OOD methods, in which the networks are actively trained on negative data from an external OOD data source [12, 25], YoLOOD offers a major advantage, as it exploits all parts of the image to model both the original data distribution and objects that might be OOD, without the need for external data sources. To tackle our approach’s requirement of bounding box annotations, we leverage the tremendous progress large language models (LLMs) have made lately to perform fully-automated dataset labeling. Using Grounding DINO [21], a multi-modal open-set object detection model that can detect arbitrary objects with textual inputs, we automatically generate bounding boxes for each image with just standard image classification annotations (*i.e.*, category names).

We perform extensive evaluations (more than 500 trained models) which demonstrate YoLOOD’s state-of-the-art performance and superiority over commonly used OOD detection methods on large-scale benchmark image datasets (*e.g.*, MS-COCO [19]). In addition, we propose new in-distribution and OOD dataset benchmarks for OOD detection in the multi-label domain. To create the in-distribution dataset, we extract a subset of 20 classes from the Objects365 dataset. With regard to the OOD datasets, we aim to create datasets that better capture the complexity of the multi-label setting in which images may contain multiple objects belonging to different class categories. The new datasets consist of: (a) a subset of the Objects365 dataset [31] in which the classes are different than the classes in the in-distribution subset mentioned above, and (b) a subset of the NUS-WIDE dataset [5]. Our results show that YoLOOD substantially improves the FPR95 compared to other OOD detection methods (*e.g.*, when using the Objects365 subset as the in-distribution dataset and the NUS-WIDE subset as the OOD dataset, the FPR95 value decreases by 12.27%).

We summarize our contributions as follows:

- We propose YoLOOD, a novel OOD detection technique for the multi-label image classification domain, which is powered by the YOLO object detection system and outperforms state-of-the-art techniques.
- We are the first to adopt the use of bounding box annotations for OOD detection - a unique technique that exploits all parts of an input image to model both the original data distribution and OOD data, without depending on external data sources.

- We introduce new benchmark datasets: (a) a large-scale in-distribution dataset, and (b) two new OOD datasets that better reflect the complexity of OOD detection in the multi-label domain (*i.e.*, images may contain multiple objects of different class categories) and make them available to the scientific community.

2. Background

2.1. Multi-Label Classification

The multi-label classification problem is defined as follows: let \mathcal{X} be the input space and \mathcal{Y} be the output space of classifier $f : \mathcal{X} \rightarrow \mathbb{R}^{|\mathcal{Y}|}$ trained on samples drawn from distribution $\mathcal{D}(\mathcal{X}, \mathcal{Y})$. Every input $x \in \mathcal{X}$ can be associated with a subset of labels $\mathcal{Y} = \{1, 2, \dots, N_c\}$, which is represented by a binary indicator vector $y = \{0, 1\}^{N_c}$, where $y_n = 1$ if the input is associated with class n .

Object detection is a variant of the multi-label classification task, in which the model also determines the location of existing objects, *i.e.*, bounding box coordinates.

2.2. Out-of-Distribution Detection

Let \mathcal{D}_{in} denote the marginal distribution of \mathcal{D} over \mathcal{X} , which represents the distribution of in-distribution data. At inference time, the system may encounter an input drawn from a different distribution \mathcal{D}_{out} over \mathcal{X} . For OOD detection in the multi-label setting, a decision function G is defined such that:

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

where x is considered OOD if none of the objects present in it are in-distribution objects.

2.3. YOLO Object Detector

In this paper, we focus on the state-of-the-art one-stage YOLO object detector and leverage its capabilities. Our proposed approach utilizes one of its latest versions, YOLOv5 [14].

YOLO’s architecture. YOLO’s architecture is comprised of two components: (a) a backbone network used to extract features from the input image, and (b) three detection heads which process the image’s features at three different scales. The size of a detection head is determined by the size of the input image and the network’s *stride* (downsampling factor) – 32, 16, and 8. This allows the network to detect objects of different sizes: the first detection head (with the largest stride) has a broader context, specializing in the detection of large objects, while the smallest one has finer resolution and specializes in the detection of small objects.

The last layer of each detection head predicts a 3D tensor of size $W \times H \times (4 + 1 + N_c)$, where $W \times H$ is the grid size (W and H are the width and height, respectively) and

$4 + 1 + N_c$ (which will be referred to as a *candidate* in the remainder of the paper) encodes three parts: the bounding box offsets, objectness score, and class scores.

Training YOLO. Every ground-truth object is associated with a single cell in each detection head. The input image is divided into an $W \times H$ grid, and the responsible cell is determined by the grid cell that the object’s center falls in. YOLO does not assume that the class categories are mutually exclusive; therefore the class score vector is trained with the multi-label configuration (the sigmoid function is applied on each output neuron).

Postprocessing. YOLO outputs a fixed amount of candidates (the amount depends on the size of the input image) which are then filtered in three sequential steps: objectness score filtering, class score filtering, and non-maximum suppression (NMS).

Additional details about YOLO’s architecture, training, and inference can be found in the supplementary material.

3. Method

Object detectors provide a natural solution for OOD detection, because they have an inherent ability to distinguish between objects of interest (in-distribution data) and irrelevant objects (OOD data). Additionally, object detectors use *passive negative learning* during training, utilizing the unlabeled data present in the training images (*i.e.*, areas that are not included within labeled bounding boxes), which allows them to better generalize to the sub-task of discarding irrelevant objects. In YOLO, this is accomplished using the objectness score. Overall, the combination of these factors makes object detectors perfect candidates for OOD detection tasks.

In this section, we introduce *YoLOOD*, our novel OOD detection method for the multi-label domain, inspired by the YOLO object detector [14]. With just minor changes, we convert a YOLO network into a multi-label image classifier.

3.1. YoLOOD Detection Layer

Since we propose a method for the image classification domain, our model does not have to predict bounding box coordinates. Therefore, we replace the last layer of each detection head with a 3D tensor of size $W_k \times H_k \times (1 + N_c)$, where k denotes the k^{th} detection head ($k \in \{1, 2, 3\}$). The size $W_k \times H_k$ grid remains identical to the grid in the original YOLO architecture, where each cell only predicts the objectness score and N_c class scores.

Formally, let $f_{\text{YoLOOD}} : \mathcal{X} \rightarrow \mathcal{C}$ be a *YoLOOD* image classifier that receives an input image $x \in \mathcal{X}$ and outputs a set of candidates $f_{\text{YoLOOD}}(x) = \mathcal{C}$ such that $\mathcal{C} = \bigcup_k \mathcal{C}_k$, where \mathcal{C}_k denotes the candidates of the k^{th} detection head. Each set of candidates \mathcal{C}_k is comprised of $W_k \times H_k$ candidates (based on the downsampling factor mentioned in Section 2), and thus $|\mathcal{C}| = \sum_{k=1}^3 |\mathcal{C}_k| = \sum_{k=1}^3 W_k \cdot H_k$.

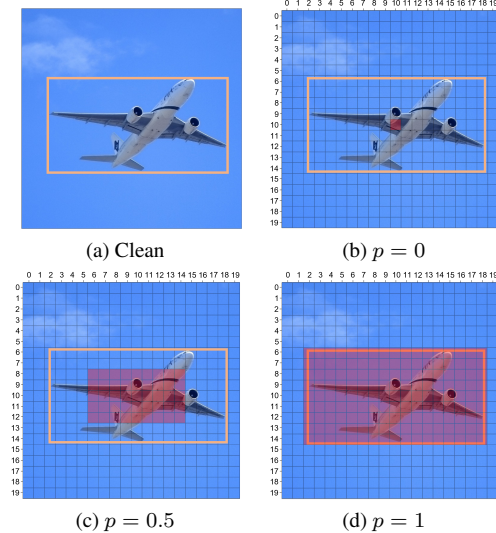


Figure 2. An example of an image divided into a size 20×20 grid with various p values (k is omitted for simplicity). The ground-truth bounding box is framed in orange. Grid cells highlighted in red represent the responsible cells (explained in Section 3.1).

Objectness score. As explained in Section 2, in the original YOLO model, only a single cell grid is responsible for the prediction of an object. However, in image classification, we are not limited to the case of a specific cell predicting the object’s bounding box. Therefore, we propose enlarging the responsible areas, such that grid cells that contain *any* part of an object are responsible for its detection, allowing the model to capture broader representations of objects.

However, an important aspect to consider is that bounding boxes do not accurately segment the object (*i.e.*, objects’ shapes are not necessarily rectangular), which may result in background areas that are not part of the object included in the bounding box. This might result in incorrect associations of irrelevant areas within the object’s bounding box with the object. Therefore, we propose using only a portion of the grid cells that cover the object’s area.

More formally, let $(W_o, H_o) \in [0, 1]^2 \subset \mathbb{R}^2$ be the normalized width and height of an object. The relative width and height of an object in a specific grid is defined as: $(W_r, H_r) = (W_o \cdot W_k, H_o \cdot H_k)$. In addition, let $(x_{k,\text{center}}, y_{k,\text{center}})$ denote the indices of the cell that the object’s center falls in in the k^{th} grid. We define p_k as the portion of grid cells relative to the grid’s size with regard to the object’s center, *i.e.*, the responsible grid cells expand from the object’s center to the object’s boundaries as a function of p_k . As shown in Figure 2, when $p_k = 0$, the responsible cells are identical to those of the original YOLO model (single cell), while when $p_k = 1$, the responsible cells cover the area of the entire object.

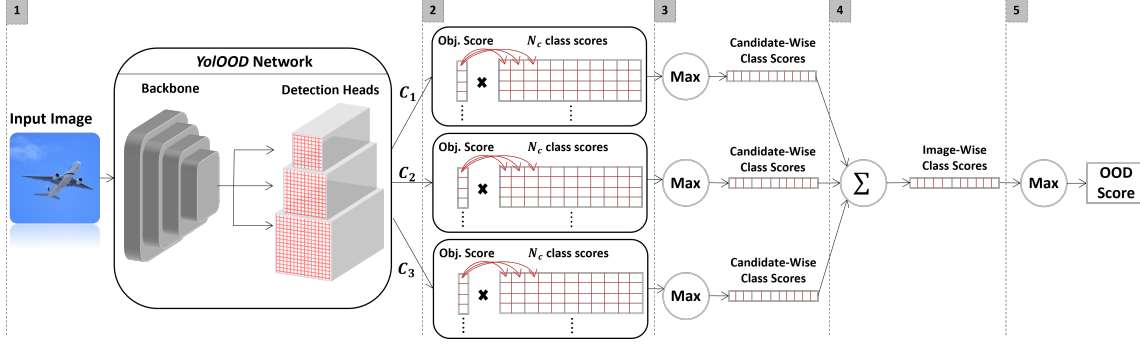


Figure 3. An overview of *YoLOOD*'s OOD detection pipeline: (1) the input image is fed to the model, which returns three groups of candidates, each of which corresponds to a different detection head; (2) each candidate's objectness score is multiplied by its corresponding class scores; (3) the highest score for each class is extracted from each group of candidates (candidate-wise class scores); (4) the candidate-wise class scores are aggregated into a single vector of image-wise class scores; and (5) the highest score, which will be referred to as the OOD score, is extracted.

Let \hat{c} denote the corresponding ground-truth value of a predicted candidate $c \in \mathcal{C}_k$. The objectness score value of \hat{c} at location (i, j) is defined as follows:

$$\begin{aligned} \varphi_u &= i \geq x_{k,\text{center}} - p_k \cdot \frac{W_r}{2}, \varphi_l = j \geq y_{k,\text{center}} - p_k \cdot \frac{H_r}{2} \\ \varphi_d &= i \leq x_{k,\text{center}} + p_k \cdot \frac{W_r}{2}, \varphi_r = j \leq y_{k,\text{center}} + p_k \cdot \frac{H_r}{2} \end{aligned} \quad (2)$$

$$\hat{c}_{\text{obj}}(i, j) = \begin{cases} 1 & \varphi_u \wedge \varphi_d \wedge \varphi_l \wedge \varphi_r \\ 0 & \text{else} \end{cases} \quad (3)$$

The effect of the different p_k values is discussed in Section 4.

Class scores. Each predicted candidate $c \in \mathcal{C}_k$ contains N_c class scores which represent the model's confidence in the presence of each class. Since multiple cells may be responsible for the prediction of an object, a cell can be assigned with multiple classes. As in the original YOLO training procedure, we train each candidate in a multi-label fashion, *i.e.*, classes are not mutually exclusive. Formally, the class score for a class $n \in \{1, \dots, N_c\}$ in \hat{c} at location (i, j) is:

$$\hat{c}_{\text{cls } n}(i, j) = \begin{cases} 1 & \text{class } n \text{ is in cell } (i, j) \\ 0 & \text{else} \end{cases} \quad (4)$$

3.2. YoLOOD Loss Function

To train *YoLOOD*, we devised a custom loss function that is composed of two components:

- Objectness score loss -

$$\mathcal{L}_{\text{obj}} = \sum_{c \in \mathcal{C}} \mathcal{L}_{\text{BCE}}(c_{\text{obj}}, \hat{c}_{\text{obj}}) \quad (5)$$

where \mathcal{L}_{BCE} denotes the binary cross-entropy loss.

- Class score loss -

$$\mathcal{L}_{\text{cls}} = \sum_{c \in \mathcal{C}} \sum_{n \in \{1, \dots, N_c\}} \hat{c}_{\text{obj}} \cdot \mathcal{L}_{\text{BCE}}(c_{\text{cls } n}, \hat{c}_{\text{cls } n}) \quad (6)$$

Finally, the total loss function is:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{obj}} + \mathcal{L}_{\text{cls}} \quad (7)$$

3.3. YoLOOD as a Multi-Label Classifier

Since the core definition of a multi-label classifier's output is a single vector containing class probabilities, we aggregate *YoLOOD*'s output such that the score for each class n is the highest score across all the candidates and is formally defined as:

$$y_n = \max_{c \in \mathcal{C}} \{\sigma(c_{\text{obj}}) \cdot \sigma(c_{\text{cls } n})\} \quad (8)$$

where σ denotes the sigmoid function.

3.4. YoLOOD for Multi-Label OOD Detection

Similar to the YOLO candidates' postprocessing (see Section 2.3), we utilize both the objectness and class scores and propose using *YoLOOD* for OOD detection in the following way (visualized in Figure 3):

$$\text{YoLOOD}(x) = \max_{n \in \{1, \dots, N_c\}} \sum_{c_k \in f_{\text{YoLOOD}}(x)} \max_{c \in \mathcal{C}_k} \{\sigma(c_{\text{obj}}) \cdot \sigma(c_{\text{cls } n})\} \quad (9)$$

$$G(x, \tau) = \begin{cases} 1 & \text{YoLOOD}(x) \geq \tau \\ 0 & \text{YoLOOD}(x) < \tau \end{cases} \quad (10)$$

where τ denotes the threshold, which can be selected according to the value that yields a high percentage (*e.g.*, 95%) of in-distribution data correctly classified by $G(x, \tau)$.

4. Evaluation

4.1. Experimental Setup

In-distribution datasets. We consider the following in-distribution datasets, originally proposed in [11]:

- PASCAL VOC [7] - consists of 5,717 training, 5,823 validation, and 10,991 test images across 20 class categories.
- MS-COCO [19] (2017 version) - consists of 117,266 training, 4,952 validation, and 40,670 test images across 80 class categories.

In addition, we propose a new benchmark, which is a subset of the Objects365 dataset [31]:

- Objects365_{in} - a subset of the original dataset which is comprised of the 20 most frequent classes that do not overlap with the classes in the OOD datasets (presented below). It consists of 68,723 training, 5,000 validation, and 10,000 test images. Further details can be found in the supplementary material.

The training and validation images are used in the training process, while the test set is used as an in-distribution set in the OOD detection evaluation. Furthermore, to demonstrate that our approach is not limited to datasets that contain bounding box annotations, we employ Grounding DINO [21], a multi-modal open-set object detection model. Grounding DINO takes a pair of $\langle \text{IMAGE}, \text{CAPTION} \rangle$ and returns the location of objects in the image based on the provided caption. In our research, we use Grounding DINO to automatically annotate the training set images, where the caption is simply a concatenation of the category names present in the image (*i.e.*, standard image classification annotations).

Out-of-distribution datasets. In previous studies [11, 34] proposing solutions for OOD detection in the multi-label setting, the effectiveness of the proposed method was evaluated on a subset of images from the ImageNet-22K [30] and Textures [6] datasets. However, these datasets only contain images with a single class category, resulting in an oversimplified setup. Therefore, we propose two new benchmarks constructed from datasets which contain images associated with multiple class categories and instances, thus reflecting the complexity of the multi-label setting:

- Objects365_{out} - a subset of the Objects365 dataset which contains ~ 200 classes that do not overlap with any of the classes present in the in-distribution datasets (*e.g.*, lamp, tomato), and specifically with the Objects365_{in} subset. This subset contains 11,669 images.
- NUS-WIDE_{out} - a subset of the original NUS-WIDE dataset [5]. We remove overlapping class categories, which leaves us with a subset of 54 categories (*e.g.*, toy, tree, whale). This subset contains 13,149 images.

Further details about the datasets are presented in the supplementary material.

Metrics. In our evaluation, performance is measured with metrics commonly used in the OOD detection domain: (a) *FPR95* - the false positive rate of OOD samples when the true positive rate is at 95%; (b) *AUROC* - the area under the receiver operating characteristic curve; and (c) *AUPR* - the area under the precision-recall curve.

Networks' architecture. For YolOOD, we use the latest version of the YOLO object detector, YOLOv5 [14], pre-trained on the MS-COCO dataset. As explained in Section 2.3, the network is comprised of a backbone and three detection heads. YOLOv5 provides several model sizes: nano, small, medium, *etc.*, each of which contains a different number of learnable parameters for the backbone and detection heads. We use the YOLOv5 small version (YOLOv5s), which contains $\sim 4.1M$ and $\sim 3M$ learnable parameters in the backbone and detection heads, respectively, which amounts to a total of $\sim 7.1M$ parameters. To apply our OOD detection approach, we replace the last layer of each detection head with the detection layer described in Section 3.1.

To perform a fair comparison between the proposed method and other state-of-the-art OOD detection methods, we also train a multi-label image classifier based on the backbone of YOLOv5s. To this end, we replace the detection heads with three fully connected layers, resulting in a similar-sized network (total of $\sim 7.1M$ parameters). This network is referred to as *YOLO-cl*s in the evaluation.

Training details. For each in-distribution dataset, we fine-tune a pair of *YolOOD* and *YOLO-cl*s models using the backbone's pretrained weights for a total of 30 epochs. More precisely, we fine-tune five pairs of models, each initialized with a different seed. The results presented in the paper are averaged across the models (the complete results, including the standard deviation values, are included in the supplementary material). We use the Adam optimizer [15] with an initial learning rate of 10^{-5} and 10^{-4} respectively for the backbone and the remaining layers. The learning rate is reduced by a factor of 10 if the mAP on the validation set does not improve for two consecutive epochs. For the *YOLO-cl*s model, we apply the logistic sigmoid function on the outputs of the classification layer (*i.e.*, logits) for multi-label training. For both models, the images are resized to 640×640 pixels and applied with color-based augmentations and geometric transformations. The mAPs for YolOOD are on par with those of YOLO-cl on all of the in-distribution datasets evaluated (~ 1 -2% difference). Detailed results are presented in the supplementary material. The source code can be found online.¹

4.2. Results

Effect of responsible grid cell percentage p_k . We characterize the effect of the percentage $\{p_k | k \in \{1, 2, 3\}\}$ of

¹<https://github.com/AlonZolfi/YolOOD>

Table 1. Comparison of the OOD detection performance of YoLODD vs. state-of-the-art methods. \downarrow indicates that lower values are better, and \uparrow indicates that higher values are better. Bold indicates superior results.

¹Trained using the auto-generated annotations. ²Trained using the original annotations.

Method	\mathcal{D}_{out} \mathcal{D}_{in}	Objects365 _{out}		Objects365 _{in}	NUS-WIDE _{out}		Objects365 _{in}
		PASCAL-VOC	MS-COCO	FPR95 \downarrow / AUROC \uparrow / AUPR \uparrow	PASCAL-VOC	MS-COCO	
MaxLogit [11]		28.91 / 94.96 / 95.32	16.39 / 96.90 / 99.17	29.95 / 94.33 / 94.38	23.60 / 95.99 / 96.05	12.16 / 97.53 / 99.24	38.07 / 92.62 / 91.48
MSP [10]		50.78 / 88.36 / 88.61	46.25 / 86.78 / 95.63	65.20 / 83.99 / 84.13	47.34 / 89.34 / 88.71	40.89 / 88.33 / 95.53	78.08 / 78.42 / 76.91
Mahalanobis [17]		73.34 / 73.90 / 70.94	87.69 / 52.13 / 77.53	83.30 / 63.23 / 56.69	77.23 / 73.76 / 67.76	88.74 / 60.35 / 80.85	88.54 / 62.34 / 54.12
ODIN [18]		28.91 / 94.96 / 95.32	16.39 / 96.90 / 99.17	29.95 / 94.33 / 94.38	23.60 / 95.99 / 96.05	12.16 / 97.53 / 99.24	38.07 / 92.62 / 91.48
JointEnergy [34]		27.90 / 95.37 / 96.04	14.80 / 97.16 / 99.28	23.13 / 95.84 / 96.20	20.19 / 96.53 / 96.76	8.29 / 97.90 / 99.39	24.46 / 95.34 / 94.96
YoLODD-a ¹		18.37 / 96.10 / 95.85	11.70 / 97.21 / 99.19	18.40 / 95.76 / 95.15	21.24 / 96.29 / 96.08	7.62 / 98.13 / 99.43	12.19 / 97.64 / 97.29
YoLODD-o ²		16.38 / 96.60 / 96.49	11.53 / 97.29 / 99.23	17.24 / 95.97 / 95.42	18.47 / 96.85 / 96.77	4.40 / 98.56 / 99.57	9.54 / 98.01 / 97.61

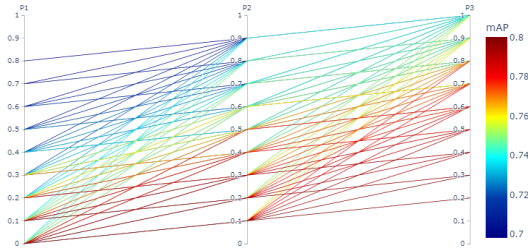


Figure 4. Models’ mAP for different p_k combinations on the COCO dataset. Each line represents a model trained with a single combination.

responsible cells (described in Section 3.1), where p_1 (resp. p_3) represents the smallest (resp. largest) detection head. We perform an extensive evaluation to determine the effect of different p_k combinations, where p_k is selected from 11 evenly spaced numbers in the range $[0, 1]$. To limit the number of possible combinations, we set a constraint such that $p_3 > p_2 > p_1$, based on the fact that the grid’s resolution increases in each subsequent detection head, resulting in 165 different combinations. An illustration of the results obtained on the COCO dataset is provided in Figure 4. After training, we sort all of the models according to their in-distribution mAP and select the 20 best-performing models. Then, we count the number of occurrences of each p_k value (each p_k is counted independently, not as a triplet) and select the most frequent values. After aggregating the results over all in-distribution datasets, we found that the best configuration is $(p_1, p_2, p_3) = (0.0, 0.1, 0.5)$. As expected, the detection head with the highest resolution ($k = 3$) benefits greatly from a large portion of responsible cells, while the detection head with the smallest resolution works best with just a single responsible cell. It should be emphasized that we recommend using this configuration for all datasets, *i.e.*, p_k is not a hyperparameter that should be tuned. The complete results can be found in the supplementary material.

YoLODD vs. state-of-the-art OOD detection methods.

We compare our approach to state-of-the-art OOD detection methods using the *YOLO-cls* network: (a) MaxLogit [11], (b) Maximum Softmax Probability (MSP) [10], (c) ODIN [18], (d) Mahalanobis [17], and (e) JointEnergy [34].

Further details about each method can be found in the supplementary material. In Table 1 we can see that YoLODD outperforms all baselines and state-of-the-art approaches, including the multi-label OOD detection method, JointEnergy. Specifically, when the networks are trained on the PASCAL-VOC, MS-COCO, and Objects365_{in} datasets with the *auto-generated* annotations (*i.e.*, constructed using Grounding DINO with solely standard classification annotations), and OOD detection is evaluated on the Objects365_{out} OOD dataset, YoLODD reduces the FPR95 by 9.53%, 3.10%, and 4.73%, compared to the best-performing method.

It should also be noted that the Mahalanobis method shows poor performance compared to the other methods examined. In [34], the authors hypothesized that the Mahalanobis method may not be well suited for the multi-label task, since it is based on the assumption that feature representation forms class-conditional Gaussian distributions. This assumption also holds in our case, however from a different perspective - the *YOLO-cls* backbone used for the evaluation was originally trained for the object detection task and therefore learned different feature representations. Furthermore, the performance of ODIN and MaxLogit is similar, since the best hyperparameter configuration found for ODIN is equivalent to MaxLogit (a special case where the temperature is set at one and the magnitude of noise is set at zero), which correlates with the results in [34].

Objectness score vs. class score vs. joint score. We also perform an analysis to examine the effect of different aggregation methods on the candidates’ output scores, *i.e.*, objectness and class scores. We consider two different approaches in addition to the regular YoLODD approach described in Equation 9 (which uses a joint probability between the objectness and class scores):

$$\text{YoLODD}_{\text{Cls}}(x) = \max_{n \in \{1, \dots, N_c\}} \sum_{c_k \in f_{\text{YoLODD}}(x)} \max_{c \in C_k} \{\sigma(c_{\text{cls}} n)\} \quad (11)$$

$$\text{YoLODD}_{\text{Obj}}(x) = \sum_{c_k \in f_{\text{YoLODD}}(x)} \max_{c \in C_k} \{\sigma(c_{\text{obj}})\} \quad (12)$$

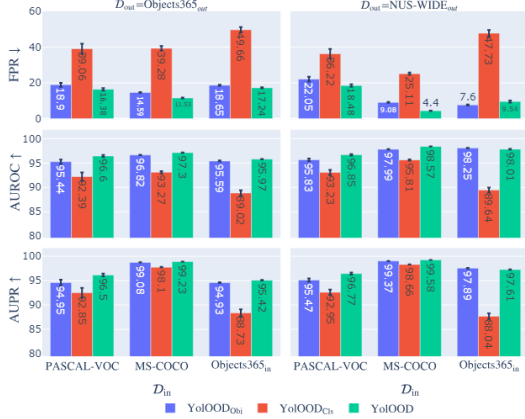


Figure 5. A comparison of the impact of different aggregation methods on YoLOOD candidates’ output scores (Section 4.2).

In Figure 5, we can see that when using the objectness and class scores separately (YoLOOD_{Obj} and YoLOOD_{Cls}, respectively), the objectness score contributes more to the ability to distinguish between in-distribution and OOD samples. However, when using a joint probability (*i.e.*, a combination of the scores), the OOD detection performance exceeds the performance of both individually. Interestingly, when using Objects365_{in} as the in-distribution dataset, YoLOOD_{Obj} achieves similar performance to YoLOOD. This can be attributed to the greater complexity of the image distribution in Objects365_{in} compared to PASCAL-VOC and COCO. This complexity results in the model learning less robust class scores, ultimately impacting the joint score.

Effectiveness of different aggregation functions. Since YoLOOD’s output is comprised of candidates from three detection heads, we consider the use of different aggregation functions between them, *i.e.*, instead of combining all of the candidates from all three detection heads into one large set $\mathcal{C} = \bigcup_k \mathcal{C}_k$, we extract the best candidate from each set \mathcal{C}_k and then apply the aggregation function. In addition, we also examine the effect of different aggregation functions on the class score output vector (similar to [34] which proposed summing energies over all of the labels). More formally, we replace Equation 9 with the following:

$$\text{agg}_1 \quad \text{agg}_2 \quad \max_{c \in \mathcal{C}_k} \{ \sigma(c_{\text{obj}}) \cdot \sigma(c_{\text{cls } n}) \} \quad (13)$$

$n \in \{1, \dots, N_c\}, \mathcal{C}_k \in f_{\text{YoLOOD}}(x)$

where agg_1 can either be the summation or max function, and agg_2 can be one of the following functions: summation, multiplication, or max.

In general, the results presented in Table 2 show that on most in-distribution datasets, the approach that yields the best results is summing the scores over the different candidate sets and then extracting the maximum value of all of the class scores. It is interesting to see that when summing over the class scores (Table 2, bottom), multiplication between the candidate sets works better than the other functions.

Table 2. OOD detection performance when using different combinations of aggregation functions for YoLOOD’s detection heads and class scores output vector. The results are averaged across the OOD datasets.

Class Agg.	Head Agg.	D_{in}		PASCAL-VOC		MS-COCO		Objects365 _{in}	
		FPR95 ↓ / AUROC ↑ / AUPR ↑		FPR95 ↓ / AUROC ↑ / AUPR ↑		FPR95 ↓ / AUROC ↑ / AUPR ↑			
Max	Max	24.10 / 95.63 / 95.61	8.55 / 98.01 / 99.43	21.01 / 95.70 / 95.08					
	Multiply	19.76 / 96.02 / 95.74	10.55 / 97.59 / 99.30	18.22 / 96.23 / 95.68					
	Sum	17.43 / 96.73 / 96.64	7.97 / 97.93 / 99.41	13.39 / 96.99 / 96.52					
Sum	Max	44.56 / 78.31 / 70.89	54.70 / 80.59 / 92.21	35.26 / 85.67 / 79.28					
	Multiply	19.58 / 96.04 / 95.80	9.61 / 97.81 / 99.38	17.55 / 96.39 / 95.94					
	Sum	40.11 / 84.91 / 81.83	39.50 / 90.20 / 96.71	26.18 / 93.68 / 92.39					

Combining YoLOOD with JointEnergy. Wang *et al.* [34] presented the JointEnergy technique in the following way:

$$E_{y_n}(x) = -\log(1 + e^{f_{y_n}(x)}) \quad (14)$$

$$E_{\text{joint}}(x) = \sum_{n=1}^{N_c} -E_{y_n}(x)$$

where $f_{y_n}(x)$ denotes the logit of the n^{th} class.

Since YoLOOD’s output can be transformed into a single vector of class probabilities (Section 3.3), JointEnergy can be applied and used with the YoLOOD architecture. Therefore, instead of applying the sigmoid function to the objectness and class scores (Equation 9), we apply the energy function E_{y_n} and combine label-wise energies over all labels for a single OOD score. Formally, the combination of YoLOOD and JointEnergy can be written as follows:

$$\text{YoLOOD}_{\text{JointEnergy}}(x) = \sum_{n \in \{1, \dots, N_c\}} \sum_{\mathcal{C}_k \in f_{\text{YoLOOD}}(x)} -\max_{c \in \mathcal{C}_k} \{ E_{y_n}(c_{\text{obj}}) \cdot E_{y_n}(c_{\text{cls } n}) \}$$

We compare the results obtained with YoLOOD_{JointEnergy} to the performance of JointEnergy when using the YOLO-clc network. We observe that YoLOOD_{JointEnergy} outperforms JointEnergy when applied to the YOLO-clc network on most in-distribution and OOD examined datasets across all metrics. For example, when using Objects365_{out} as the OOD dataset, the FPR95 metric improves by 10.49%, 7.7%, and 7.22% when the networks are trained on the PASCAL-VOC, MS-COCO, and Objects365_{in} datasets, respectively. Moreover, in some cases, the performance of YoLOOD_{JointEnergy} even exceeds the performance obtained using the regular YoLOOD (*e.g.*, when MS-COCO and Objects365_{out} are used as the in-distribution and OOD datasets, respectively, the FPR95 decreases from 11.53% to 7.1%).

YoLOOD vs. a vanilla YOLO We also examine the advantage of using YoLOOD over a vanilla YOLO (*i.e.*, a standard object detector). We train a YOLO detector using the same configuration used to train YoLOOD for a fair comparison: pretrained weights are only used in the backbone layers and no custom augmentations are used (*e.g.*, mosaic [2]).

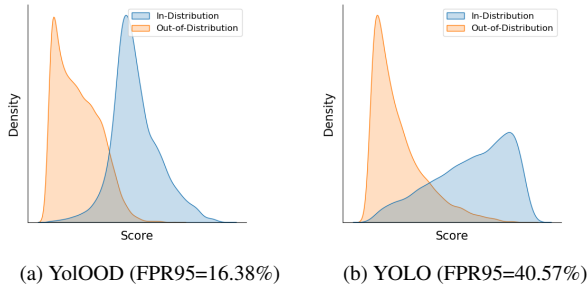


Figure 6. Score distribution when using PASCAL-VOC as the in-distribution dataset and Objects365_{out} as the OOD dataset.

The OOD score functions for the vanilla YOLO model are similar to those of YoLOOD, differing only at the anchor box level, which results in $3\times$ more candidates.

Our evaluation demonstrates that YoLOOD outperforms the vanilla YOLO on all of the in-distribution and OOD datasets examined. An example of this is provided in Figure 6, which presents the distribution of YoLOOD and vanilla YOLO scores when using PASCAL-VOC and Objects365_{out} as the in-distribution and OOD datasets, respectively. As can be seen, there are notable differences in the models’ performance. The vanilla YOLO (Figure 6b) excels at assigning low scores to OOD data but faces challenges in definitively identifying in-distribution data, as evidenced by the uniformly distributed scores obtained. On the other hand, YoLOOD (Figure 6a) performs well on OOD data, while substantially enhancing in-distribution detection. We hypothesize that the difference in the models’ performance is mainly the result of the fact that (a) the loss function in the vanilla YOLO largely focuses on improving the bounding box coordinates’ regression, and (b) YoLOOD uses more responsible grid cells for the detection of each object, capturing broader representations.

5. Related Work

5.1. Object Detection

Object detectors have been studied extensively over the last few years, demonstrating state-of-the-art performance, and various solutions aimed at identifying the objects present in an image and their precise location (*i.e.*, bounding boxes) have been proposed. Modern deep learning-based object detectors are usually composed of two components, a backbone network and a head, which is used to predict the bounding boxes and classes of existing objects in an image. Broadly, there are two types of models: one-stage detectors (*e.g.*, SSD [22] and YOLO [2, 14, 27]) and two-stage detectors (*e.g.*, Mask R-CNN [9] and Faster R-CNN [29]). Object detectors developed in recent years often insert additional layers between the backbone and the head, which are used to collect feature maps from different stages [20, 33].

5.2. Out-of-Distribution Detection

Deep neural networks’ overconfidence for OOD data was first noted by [26]. Several baseline approaches have been proposed to tackle this problem; for example, Hendrycks *et al.* proposed two baselines: MSP [10] and MaxLogit, which uses the highest score from the classifier’s last layer as an OOD score [11]. In recent years, OOD detection has attracted the attention of the machine learning research community, and researchers have proposed methods aimed at improving OOD uncertainty estimation, such as: (1) ODIN [18], which combines input preprocessing and temperature scaling; (2) the Mahalanobis [17] distance-based approach, which utilizes the network’s internal feature representations; (3) the gradient-based GradNorm [13] score; (4) ReAct [32], which rectifies activation values; (5) IsoMax [24], which proposes isotropy maximization loss; and (6) the energy score [23]. These studies only addressed the multi-class classification task, and the topic of OOD detection in the multi-label domain remains largely underexplored. OOD detection in the multi-label domain has only been addressed by [11], who proposed the MaxLogit baseline, and [34], who proposed JointEnergy which combines label-wise energies over all labels.

It should be noted that the OOD detection methods proposed in many prior studies require *external* OOD data (in addition to the OOD test dataset) to improve the detector’s robustness. In some cases, the model is provided with OOD samples for hyperparameter tuning [18], while in other cases, these samples are used for negative learning in which the model is explicitly trained on images that do not contain in-distribution data [12, 25]. In contrast, object detection models in general, and YOLO in particular, do not require any type of external OOD data to model irrelevant objects and background areas.

6. Conclusion

In this paper, we presented YoLOOD – a novel approach for OOD detection in the underexplored multi-label classification domain, leveraging key concepts from object detectors, and demonstrated how all parts of an input image can be utilized to model both in-distribution and OOD data. In our evaluation, we presented new benchmark datasets that better capture the complexity of the multi-label domain. Our comprehensive evaluation demonstrated that YoLOOD outperforms existing OOD detection methods, as well as a standard YOLO object detector, achieving state-of-the-art performance. Moreover, our findings highlight the significance of the objectness score, a characteristic shared by YOLO and other object detectors, in enhancing OOD detection performance. These findings imply broader potential for the utilization of object detectors in the OOD detection task. In future work, we plan to explore modeling images at a finer granularity level (pixel-level annotations).

References

- [1] Guy Amit, Moshe Levy, Ishai Rosenberg, Asaf Shabtai, and Yuval Elovici. Food: Fast out-of-distribution detector. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2021. **1**
- [2] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, 2020. **1, 7, 8**
- [3] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiao-xiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4974–4983, 2019. **1**
- [4] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062*, 2014. **1**
- [5] Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. Nus-wide: a real-world web image database from national university of singapore. In *Proceedings of the ACM international conference on image and video retrieval*, pages 1–9, 2009. **2, 5**
- [6] Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3606–3613, 2014. **5**
- [7] Mark Everingham, SM Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015. **5**
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015. **1**
- [9] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. **8**
- [10] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations*, 2017. **6, 8**
- [11] Dan Hendrycks, Steven Basart, Mantas Mazeika, Mohammadreza Mostajabi, Jacob Steinhardt, and Dawn Song. Scaling out-of-distribution detection for real-world settings. *arXiv preprint arXiv:1911.11132*, 2019. **1, 5, 6, 8**
- [12] Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*, 2019. **2, 8**
- [13] Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional shifts in the wild. *Advances in Neural Information Processing Systems*, 34:677–689, 2021. **1, 8**
- [14] Glenn Jocher. ultralytics/yolov5: v6.0 - yolov5n 'nano' models, roboflow integration, tensorflow export, opencv dnn support, 2021. **1, 2, 3, 5, 8**
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. **5**
- [16] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. **1**
- [17] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems*, 31, 2018. **1, 6, 8**
- [18] Shiyu Liang, Yixuan Li, and R Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *6th International Conference on Learning Representations, ICLR 2018*, 2018. **1, 6, 8**
- [19] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. **2, 5**
- [20] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. **8**
- [21] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. **2, 5**
- [22] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. **1, 8**
- [23] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems*, 33:21464–21475, 2020. **8**
- [24] David Macêdo, Tsang Ing Ren, Cleber Zanchettin, Adriano LI Oliveira, and Teresa Ludermit. Entropic out-of-distribution detection. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2021. **8**
- [25] Sina Mohseni, Mandar Pitale, JBS Yadawa, and Zhangyang Wang. Self-supervised learning for generalizable out-of-distribution detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5216–5223, 2020. **2, 8**
- [26] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 427–436, 2015. **1, 8**
- [27] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. **1, 8**

- [28] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. [2](#)
- [29] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. [1](#), [2](#), [8](#)
- [30] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. [5](#)
- [31] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8430–8439, 2019. [2](#), [5](#)
- [32] Yiyou Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations. *Advances in Neural Information Processing Systems*, 34:144–157, 2021. [8](#)
- [33] Mingxing Tan, Ruoming Pang, and Quoc V Le. Efficientdet: Scalable and efficient object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10781–10790, 2020. [8](#)
- [34] Haoran Wang, Weitang Liu, Alex Bocchieri, and Yixuan Li. Can multi-label classification networks know what they don’t know? *Advances in Neural Information Processing Systems*, 34:29074–29087, 2021. [1](#), [5](#), [6](#), [7](#), [8](#)