# Fast Post-Hoc Confidence Fusion for 3-Class Open-Set Aerial Object Detection

Spyridon Loukovitis, Vasileios Karampinis, Athanasios Voulodimos

*Abstract*— Developing reliable UAV navigation systems requires robust air-to-air object detectors capable of distinguishing between objects seen during training and previously unseen objects. While many methods address closed-set detection and achieve high-confidence recognition of in-domain (ID) targets, they generally do not tackle open-set detection, which requires simultaneous handling of both ID and out-of-distribution (OOD) objects. Existing open-set approaches typically rely on a single uncertainty score with thresholding, limiting flexibility and often conflating OOD objects with background clutter. In contrast, we propose a lightweight, model-agnostic post-processing framework that explicitly separates background from unknown objects while preserving the base detector's performance. Our approach extends open-set detection beyond binary ID/OOD classification to real-time three-way classification among ID targets, OOD objects, and background. To this end, we employ a fusion scheme that aggregates multiple confidence estimates and per-detection features using a compact multilayer perceptron (MLP). Incorporating different logit variants into the MLP consistently enhances performance across both binary and three-class classification without compromising throughput. Extensive ablation and comparative experiments confirm that our method surpasses threshold-based baselines in two-class classification by an average of 2.7% AUROC, while retaining or improving open-set mAP. Furthermore, our study uniquely enables robust three-class classification, a critical capability for safe UAV navigation, where OOD objects must be actively avoided and background regions safely ignored. Comparative analysis highlights that our method surpasses competitive techniques in AUROC across datasets, while improving closed-set mAP by up to 9 points, an 18% relative gain.

*Index Terms*— Computer Vision for Transportation, Autonomous Vehicle Navigation, Systems: Perception and Autonomy

## I. INTRODUCTION

Conventional object detectors are trained and evaluated under the closed-set assumption, recognizing only categories observed during training. Recent methods tackled the task of domain generalization [1], [2] seeking to preserve in-domain recognition under distribution shifts (e.g., synthetic-to-real transfer), while sustaining high-confidence predictions on known classes across varied conditions. However, these approaches do not address the presence of previously unseen objects in realistic aerial settings. In practice, such models often suppress or misclassify unfamiliar targets, a failure that can jeopardize safe navigation for autonomous systems.

To address the presence of previously unseen classes in real-world data, open-set [3] detection has emerged as a promising paradigm, enabling perception systems to reject unknown instances while reliably recognizing known categories. Prior work on open-set detection [4], [5] mainly focused on distinguishing the in-distribution from out-of-distribution (ID/OOD), with little emphasis on examining whether OOD predictions correspond to genuine objects from unseen classes or merely background clutter. This distinction is especially important in navigation settings [6], where OOD objects may pose safety risks and must be actively avoided, whereas background regions can be safely ignored. In practice, setting the confidence threshold for objectness too low results in numerous background detections being perceived as obstacles, reducing path planning efficiency, while setting it too high causes missed obstacles that threaten safety. The issue is further amplified in aerial detection, where a large portion of false positives originates from the background, and the threat of drones intruders remains a critical concern for safe operation.

In this work, we tackle the above-mentioned issue by introducing a lightweight, post-hoc confidence fusion framework that extends open-set detection to a three-class setting: distinguishing between in-distribution objects, out-of-distribution objects, and background clutter. Our method leverages multiple complementary uncertainty signals, including softmax confidence, entropy, and Gaussian mixture model (GMM) densities [7], which are then fused into a single calibrated confidence score using a shallow multilayer perceptron (MLP). The MLP is deliberately small, enabling training in under two minutes on a CPU at deployment time, making it no more expensive than the thresholding strategies widely used today. Moreover, our framework does not modify the detector itself, ensuring full model-agnostic compatibility whenever multiple confidence signals are available.

We validate our approach using RT-DETR [8] with a spectrally normalized backbone, incorporating GMMs in the embedding space and an MLP postprocessor. Our method achieves real-time performance without degrading closed-set detection accuracy, while introducing the ability to explicitly separate background from OOD targets. This capability is particularly valuable for UAV navigation, where avoiding OOD objects and ignoring background must be treated as fundamentally different decisions.

The main contributions of this work are summarized as

Spyridon Loukovitis is with the School of Electrical & Computer Engineering, National Technical University Athens, Polytechnioupoli, Zografou, 15780, Greece (el20120@mail.ntua.gr).

Vasileios Karampinis is with the School of Electrical & Computer Engineering, National Technical University Athens, Polytechnioupoli, Zografou, 15780, Greece (vkarampinis@ails.ece.ntua.gr).

Athanasios Voulodimos is with the School of Electrical & Computer Engineering, National Technical University Athens, Polytechnioupoli, Zografou, 15780, Greece (thanosv@mail.ntua.gr).

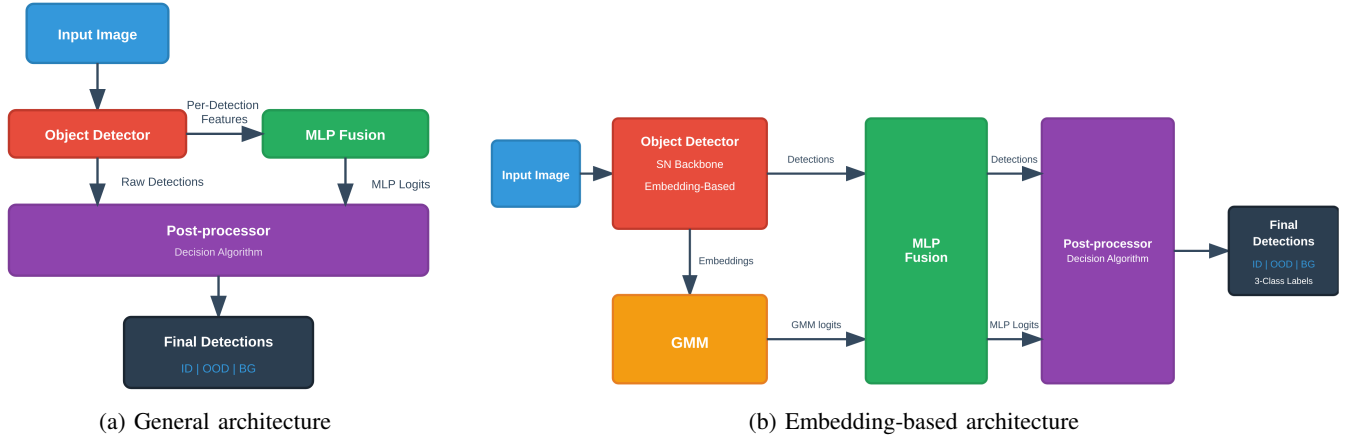(a) General architecture       (b) Embedding-based architecture

Fig. 1: Comparison of general and embedding-based feature fusion architectures.

follows:

- **Three-Class Open-Set Detection:** We extend open-set aerial detection beyond the binary ID/OOD setting to a three-class problem that explicitly distinguishes between in-distribution objects, out-of-distribution objects, and background clutter, addressing a critical gap in UAV perception.
- **Improved Two-Class Performance:** In the standard ID/OOD setting, our approach achieves marginal but consistent improvements over existing single-score and threshold-based baselines.
- **Model-Agnostic Post-Hoc Fusion:** Our method does not alter detector training and can be applied to any architecture that provides multiple confidence signals, ensuring broad applicability and ease of deployment.
- **Effective Real-Time Implementation:** We provide a complete implementation combining RT-DETR with a spectrally normalized backbone, Gaussian mixture models, and an MLP postprocessor, demonstrating strong performance under realistic conditions while preserving real-time throughput.

## II. RELATED WORK

### A. Domain Generalization

Domain Generalization (DG) seeks to learn domain-invariant representations that transfer to unseen domains without any access to target-domain data during training.

**Domain Generalization in Object Detection**. DG has been used for various vision tasks, including classification, semantic segmentation, and object detection. Object detection emphasizes learning domain-invariant representations and producing robust proposal features. Several approaches have been examined for maintaining robust feature representations under domain shift, with one important direction involving feature alignment through adversarial training [9]. Other studies [10] introduced epistemic and aleatoric uncertainty as a means of evaluating the generalization capability of model features in unseen domains, and proposed a metric to assess feature robustness, an important addition to mitigating the risk of navigation.

**Domain Generalization in Aerial Object Detection**. The importance of generalization in unseen domains is equally critical in air-to-air object detection. Early works on aerial object detection highlighted the challenge of constantly detecting other UAVs, primarily due to their small size, which occupies roughly 5% of the observed region [11]. To address this challenge, methods such as AirTrack [12] proposed a framework based on frame alignment and cascading detection. Other studies [13] presented a combination of detection and tracking for enhancing autonomy in Advanced Air Mobility systems. Building upon the framework in [13], [14] introduced an auxiliary branch based on depth estimation, providing a relative estimate to the distance of the approaching object. While the above works demonstrated promising results in aerial object detection, they provided no evaluation on their framework's ability to generalize. Adverse weather conditions are common in the real-world, such phenomena introduce significant domain shift, requiring robust frameworks that are capable of generalizing in unseen domains. To evaluate detector generalization under such conditions, [15] introduced the AOT-C dataset, which consists of corrupted data based on different weather conditions and sensor corruptions like blur, noise, and color quantization.

Despite substantial progress in aerial object detection, the discussed methods assume a closed-set of pre-established classes (e.g., aircraft, drones). In real-world deployments, however, UAVs frequently encounter novel objects. Closed-set detectors tend to misclassify such instances or fail to detect them, undermining robustness in dynamic environments. This limitation underscores the need for advancing current methods toward open-set aerial object detection.

### B. Open-Set Aerial Object Detection

Out-of-distribution (OOD) [3] detection has been extensively studied as a mechanism for identifying samples outside the training distribution and mitigating overconfidence in neural networks. The most common methods in OOD estimate ID density [16] and exclude test samples that deviate from the estimated distribution. More recent approaches have introduced teacher-student learning architectures with

masked image modeling training scheme [17], while others explored meta-learning techniques [18] or even Vision-Language models [19], achieving substantial performance uplifts and offering alternative strategies to address the critical challenge of OOD detection in navigation. Open-set [3] extends the OOD detection to the object detection task by localizing known objects while simultaneously ignoring unknown ones. Approaches to this problem generally fall into three categories. Post hoc scoring methods [4], [20], [21] add an unknown rejection score to a pretrained detector, using logits, features, or sensitivity signals, allowing integration without retraining. Training time unknown-aware detection [5] modifies the detector or its loss functions to better separate objectness from class evidence (e.g., by introducing unknown/objectness heads, background re-labeling, or regularization). While this typically improves recall for unknown objects, it requires retraining. Open-world detection [22], [23] integrates unknown detection with discovery and incremental learning by clustering unknown proposals, expanding the label space, and addressing the challenge of forgetting previously seen objects.

In robotics, methods such as Open-set RCNN [24] demonstrated promising results by introducing a training time approach based on prototype and instance-level contrastive learning for separating the known objects from the background clutter. Other approaches, such as SAFE [21], proposed a post-hoc strategy by incorporating a light MLP-head into a pretrained detector to differentiate between adversarial and in-distribution examples, classifying the adversarial ones as unseen objects. Additional methods have focused on uncertainty quantification, incorporating models of epistemic and aleatoric uncertainty, enhancing robustness in safety-critical tasks. For instance, [4] introduced Gaussian Mixture Models (GMMs) for post-hoc uncertainty estimation, while others applied spectral normalization and temperature scaling. [25] built upon the introduced uncertainty quantification methods, introducing a model-agnostic framework capable of integrating GMMs, spectral normalization, and temperature scaling into any embedding-based object detector, significantly improving the AUROC a metric showcasing the effectiveness of the separation between known and unknown objects.

Although these methods achieved promising results in distinguishing known objects from novel ones, none have addressed the more nuanced challenge of separating OOD instances into previously unseen objects and background clutter. This distinction is particularly critical in aerial navigation, where the majority of OOD detections come from background clutter. To address this limitation, our work introduces a post-hoc confidence fusion framework that integrates multiple uncertainty signals into a single calibrated confidence score through an optimized MLP head.

## III. METHODOLOGY

In this section, we introduce a general, detector-agnostic algorithm that fuses multiple confidence estimates and per-detection features through a lightweight multilayer percep-

TABLE I: The benchmarking results of 13 object detectors on AOT and AOT-C in terms of Average Precision (AP), inference speed (fps) and model size (M)

| Object detector | $AP_{clean}$ ↑ | $AP_{cor}$ ↑ | fps ↑ | Model Size (M) ↓ |
|---|---|---|---|---|
| YOLOv5 [26] | 64.6 | 53.5 | 99 | 46.5 |
| YOLOv8 [27] | 56.4 | 41.2 | **110** | 43.7 |
| YOLOX [28] | **69.3** | 43.8 | 68 | 54.2 |
| RetinaNet [29], [30] | 35.7 | 20.0 | 17 | **37.9** |
| FasterR-CNN [31], [32] | 52.9 | 29.7 | 15 | 41.3 |
| DiffusionDet [33] | 63.8 | 35.7 | 30 | 110.5 |
| DETR [34] | 58.7 | 26.1 | 27 | 41.2 |
| CenterNet2 [35] | 66.2 | 35.9 | 24 | 71.6 |
| GMM-DET (FasterR-CNN) [4] | 64.2 | 48.0 | 15 | 41.3 |
| RT-DETR-R50 [8] | 66.2 | 49.6 | 28 | 40.1 |
| Joint Thresholding | 66.8 | 49.3 | 28 | 40.1 |
| **Our Method 2 class** | 65.0 | 49.3 | 27 | 40.2 |
| **Our Method 3 class** | 69.2 | **58.7** | 27 | 40.2 |

tron (MLP), as illustrated in Fig.1a. This formulation provides a flexible framework for improving the area under the ROC curve (AUROC) by learning to combine complementary uncertainty cues. Building on this approach, we propose a model-agnostic embedding-based variant (Fig.1b), which leverages the intermediate feature representations produced by modern detectors to achieve enhanced uncertainty calibration. Most importantly, our framework extends beyond standard binary in-/out-of-distribution classification and enables explicit three-class discrimination between in-distribution objects, out-of-distribution objects, and background clutter. This capability is particularly important for reliable autonomous navigation in both terrestrial and aerial environments, where safety-critical decisions must depend on robust uncertainty estimation.

### A. General Fusion Algorithm

We first describe the general model-agnostic algorithm for open-set aerial object detection, illustrated in Algorithm 1. The approach constructs a new training set of per-detection features and labels by running a pretrained detector on data containing both in-distribution (ID) and out-of-distribution (OOD) samples. Each detection yields a feature vector that may include raw detector confidences, uncertainty scores, logits, or embeddings, along with a label indicating whether the detection corresponds to an ID object, an OOD object, or background clutter.

Formally, given a detector and a dataset $\mathcal{D} = \mathcal{D}_{ID} \cup \mathcal{D}_{OOD}$, the detector is applied to all images. Each prediction is matched to its ground-truth label, producing pairs $(X_i, Y_i)$ where $X_i$ is the feature vector of the detection and $Y_i \in \{ID, OOD, BG\}$ is the class label. This collection of pairs constitutes a new dataset tailored for uncertainty calibration.

From this dataset, a subset of features is selected to serve as input to a lightweight multilayer perceptron (MLP). The desired output configuration is also chosen: a binary classifier ($K = 2$) for ID vs. OOD, or a three-way classifier ($K = 3$) for ID, OOD, and background. The MLP is then trained on the constructed dataset. Lastly, thresholds are tuned on the MLP logits to satisfy desired open-set recognition guarantees (e.g., controlling the false acceptance rate), which need not correspond to a simple $\arg\max$ decision rule.

At deployment time, for each new detection, the same

**Algorithm 1** Model-Agnostic Open-Set Detection with MLP Fusion

1: **Definitions:**
    - *Detector Output*: per detection features $f$
    - *Scores*: confidences and uncertainties derived from features $c$
    - *OOD Dataset*: (X,Y) $\mathcal{D}$

2: **procedure** TRAIN_FUSION_MLP($\mathcal{D}, f, c$)
3:     **Select features**: From $f \cup c$ select $f_s$
4:     **Select outputs**: choose $K \in \{2, 3\}$
5:     $S \leftarrow \emptyset$
6:     **for all** images $u \in \mathcal{D}_{\text{ID}} \cup \mathcal{D}_{\text{OOD}}$ **do**
7:         Run detector $\rightarrow$ predictions $x_i = f_s$
8:         Match with labels for $y_i$
9:         $S \leftarrow S \cup \{(x_i, y_i)\}$
10:    **end for**
11:    Train a $K$-way MLP classifier $g(\cdot)$ on $S$
12:    $thresholding(\cdot) \leftarrow$ decision boundaries
13: **end procedure**

14: **function** CLASSIFY_DETECTION($f_s$)
15:    $logits \leftarrow g(f_s)$
16:    decision $\leftarrow thresholding(logits)$
17:    **return** decision
18: **end function**

---

**Algorithm 2** Embedding Based Algorithm

1: **Definitions:**
    - *Detector output*: class logits $l$, bounding boxes $b$, embeddings $e$
    - *GMM logits*: $\ell_{\text{gmm}}$
    - *GMM dataset*: $(X, Y)$ with ID class labels for GMMs
    - *OOD dataset*: $\mathcal{D} = \mathcal{D}_{\text{ID}} \cup \mathcal{D}_{\text{OOD}}$

2: **procedure** TRAIN_GMM($X, Y$)
3:    **for all** images $x \in X$ **do**
4:         Run detector $\rightarrow$ predictions $(b_i, l_i, e_i)$
5:         Match predictions to GT via Hungarian matcher
6:         Assign $e_i$ to its GT label
7:    **end for**
8:    **for all** class $c$ with samples $x_c \subset X$ **do**
9:         $\mu_c \leftarrow \frac{1}{|x_c|} \sum_{x_c} f_\theta(x_c)$
10:        $\Sigma_c \leftarrow \frac{1}{|x_c|-1} \sum_{x_c} (f_\theta(x_c) - \mu_c)(f_\theta(x_c) - \mu_c)^T$
11:        $\pi_c \leftarrow \frac{|x_c|}{|X|}$
12:    **end for**
13: **end procedure**

14: **procedure** TRAIN_FUSION_MLP($\mathcal{D}, f, c$)
15:    **Select features**: From $f \cup c$ select $f_s$
16:    **Select outputs**: choose $K \in \{2, 3\}$
17:    $S \leftarrow \emptyset$
18:    **for all** images $u \in \mathcal{D}_{\text{ID}} \cup \mathcal{D}_{\text{OOD}}$ **do**
19:         Run detector $\rightarrow$ predictions $x_i = f_s$
20:         Match with labels for $y_i$
21:         $S \leftarrow S \cup \{(x_i, y_i)\}$
22:    **end for**
23:    Train a $K$-way MLP classifier $g(\cdot)$ on $S$
24:    $thresholding(\cdot) \leftarrow$ decision boundaries
25: **end procedure**

26: **function** CLASSIFY_DETECTION($f_s$)
27:    $logits \leftarrow g(f_s)$
28:    decision $\leftarrow thresholding(logits)$
29:    **return** decision
30: **end function**

---

set of features is extracted, the trained MLP is applied to obtain fused logits, and these are passed through the calibrated decision function. The output is a classification of each detection as ID, OOD, or background, enabling reliable open-set detection in real time.

*B. Embedding-Based Fusion Algorithm*

Building on the general framework described in Section 1, we present a more specific embedding-based implementation tailored to modern detectors that produce per-detection feature embeddings. The methodology follows prior work on embedding-space density modeling [4], [36] and extends it with calibration and pruning strategies, as well as fusion through our MLP.

1) **Detector training with spectral normalization:** The base detector is trained with spectral normalization applied to convolutional layers to enforce bi-Lipschitz continuity and produce well-behaved embeddings.
2) **Temperature calibration of logits:** On a held-out calibration set, scalar temperature parameters are learned to rescale both detector logits and GMM log-likelihoods by minimizing negative log-likelihood. This improves comparability across different uncertainty scores.
3) **Gaussian mixture modeling:** Using the training set, Gaussian mixture models (GMMs) are fitted to the embeddings of each class. Each detection embedding is then mapped to a vector of per-class GMM log-

likelihoods, which serve as additional uncertainty signals.
4) **Logit calibration:** The raw GMM log-likelihoods are rescaled using temperature scaling, ensuring that their magnitudes are consistent with detector-derived confidences.
5) **Score pruning:** Detections with low raw confidence scores (sigmoid $< 0.2$) are discarded, reducing redundancy and eliminating spurious predictions that otherwise dominate AUROC errors.

This procedure provides, for every detection, both calibrated detector scores and GMM-derived logits and confidences. These signals are then used as input features for the fusion MLP described in the previous subsection. The overall

embedding-based pipeline is summarized in Algorithm 2, which combines GMM training, fusion MLP training, and the final OOD decision rule.

## C. Detection Classification and Ground Truth Matching

To establish a consistent evaluation framework, we define how detector outputs are categorized relative to ground truth annotations. When comparing detector outputs with ground truth labels, four types of detections emerge:

1) **True Positive ID (TP-ID):** Detections that match with a known ground truth object and predict the correct class label
2) **False Positive ID (FP-ID):** Detections that match with a known ground truth object but predict an incorrect class label
3) **Out-of-Distribution (OOD):** Detections that match with ground truth objects whose class is not present in the training set
4) **Background (BG):** Detections that do not match with any ground truth objects

For the purpose of open-set detection, we classify both TP-ID and FP-ID detections as in-distribution (ID) detections. This design choice separates the problem of ID/OOD/background categorization from the problem of correct class prediction within the ID set. This definition differs from some approaches in the literature that consider only correctly classified detections as ID. When comparing against prior methods, we recompute their results according to our definition to ensure fair evaluation.

## D. Fusion MLP Training

The fusion MLP is deliberately lightweight, containing only a few hidden units. Training can be completed on a standard CPU in under two minutes, which allows the model to be recalibrated "on the fly" in real-world scenarios without significant computational overhead.

## E. Evaluation Protocol

We evaluate our method under three classification settings: (i) binary ID vs. OOD, (ii) binary ID vs. OOD+background, and (iii) three-class ID vs. OOD vs. background. For all settings, we also track mean average precision (mAP) and frames per second (FPS) to ensure that open-set calibration does not degrade closed-set accuracy or real-time performance.

## F. Domain Shift in MLP Training

Finally, we study the impact of training the fusion MLP with OOD data from sources different from the deployment domain. We find that training on unrelated datasets or synthetic features significantly degrades performance, underscoring the importance of either accessing representative OOD data from the target domain or generating realistic image-domain OOD examples that produce detector features aligned with deployment conditions.

## IV. EXPERIMENTS AND RESULTS

In this section we present the evaluation of our proposed method. We describe the experimental setup, followed by results on the binary (two-class) and three-class settings, and conclude with an ablation study on OOD data and domain shift.

## A. Experimental Setting

Our experiments are conducted using the same detector as in the Joint Thresholding baseline, namely RT-DETR with a spectrally normalized backbone [25]. We evaluate on two benchmark datasets. The first is **AOT-C**, which includes out-of-distribution samples from real flight data provided in [15]. The second is **COCO-OS**, where the first 50 COCO [37] classes are treated as in-distribution and the remaining 30 as out-of-distribution. For each dataset, we construct training, validation, and calibration splits containing only ID data, together with an OOD split that is further divided into training and validation subsets when training the MLP fusion model.

The primary evaluation metric is the **area under the ROC curve (AUROC)**, which directly measures OOD discrimination. We also report **TPR at fixed OSR levels** (5%, 10%, 20%), while ensuring that closed- and open-set **mean average precision (mAP)** and **inference speed (FPS)** are preserved, with only minor fluctuations and in some cases major improvements. The fusion MLPs are deliberately kept small, allowing training to complete on a standard CPU within a few minutes, thereby enabling practical recalibration without heavy computational overhead. Experiments using a GPU are conducted on an NVIDIA A10G.

## B. OOD Dataset Preparation

Between the training/validation splits and the OOD test set, two types of domain shift arise. The first stems from the OOD data itself, which the MLP is designed to target. The second arises because even ID and background detections can originate from different datasets, introducing additional distribution shift, mirroring deployment conditions. To address both, we prepare the OOD training data as follows.

- **Train split:**
  - ID samples: both from the detector's training set, and from the OOD test set.
  - OOD samples: from the OOD test set.
  - Background samples: both from the detector's training set, and from the OOD test set.
- **Validation split:** only samples from the OOD test set, in order to measure open-set performance under the first type of domain shift.

The ratio between ID training and ID test samples is adjusted depending on the severity of the domain shift. A higher degree of ID domain shift requires more training samples to avoid degradation in mAP. Concretely, we adopt a 1:1 ratio for AOT-C and COCO-OS, while for AOT [38] (without corruptions) we use a 3:1 ratio.

TABLE II: Ablation study for MLP inputs in the two-class setting. Each row indicates which input features are included (✓/✗). We report AUROC and TPR at fixed OSR levels (5%, 10%, 20%).

| Dataset | Score | Entropy | Density | GMM Entr. | GMM Dens. | Logits | GMM Logits | AUROC | TPR@5% | TPR@10% | TPR@20% |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Real Flights** | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | 0.891 | **0.717** | 0.754 | 0.821 |
| | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | 0.889 | 0.629 | **0.758** | 0.819 |
| | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | 0.885 | 0.681 | 0.724 | 0.795 |
| | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | **0.897** | 0.687 | 0.719 | **0.835** |
| **COCO** | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | 0.788 | 0.390 | 0.493 | 0.633 |
| | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | 0.788 | 0.390 | 0.493 | 0.633 |
| | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | **0.894** | **0.636** | **0.739** | 0.823 |
| | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | **0.894** | 0.624 | 0.702 | **0.829** |

TABLE III: Comparison of algorithms on Real Flights and COCO datasets. We report mAP, $AUROC_{bd}$, and AUROC.

**Real Flights**

| Model | Method | mAP | $AUROC_{bd}$ | AUROC |
|---|---|---|---|---|
| YOLOv5 | Standard | 39.3 | 0.800 | 0.789 |
| Faster R-CNN | GMM-DET | 35.9 | 0.775 | 0.723 |
| RT-DETR + SN | Joint | **41.0** | 0.887 | 0.874 |
| RT-DETR + SN | MLP | 39.0 | **0.887** | **0.897** |

**COCO**

| Model | Method | mAP | $AUROC_{bd}$ | AUROC |
|---|---|---|---|---|
| YOLOv5 | Standard | **44.4** | 0.839 | 0.685 |
| Faster R-CNN | GMM-DET | 41.6 | 0.836 | 0.872 |
| RT-DETR + SN | Joint | 42.0 | 0.701 | 0.756 |
| RT-DETR + SN | MLP | 42.1 | **0.845** | **0.894** |

TABLE IV: Three-class results: macro AUROC and Open-Set mAP (higher is better). An asterisk (*) in the result means that all detections in the dataset got pruned.

| Algorithm | Real Flights | | COCO | |
|---|---|---|---|---|
| | AUROC | OS mAP | AUROC | OS mAP |
| Score | 0.86 | 41.5 | 0.64 | 36.1 |
| Entropy | 0.75 | 42.7 | 0.57 | 39.7 |
| Density | 0.79 | 40.9 | 0.57 | 39.0 |
| GMM Entropy | 0.78 | **45.9** | 0.60 | 32.7 |
| GMM Density | 0.81 | * | 0.46 | 23.8 |
| **MLP** | **0.91** | 39.1 | **0.89** | **41.0** |

**mAP** was calculated with the same 20% restriction from OOD to ID miss-classification.

### C. Two-Class Setting (ID vs. OOD)

We first conduct an ablation study on the choice of input features to the fusion MLP. Table II reports AUROC and TPR at fixed OSR levels for different feature combinations. We observe that including **both detector-derived scores and GMM-based signals** results in both the highest AUROC and the best TPR at 20% OSR, in both datasets, without overfitting. This configuration is therefore adopted for all subsequent comparisons.

Table III compares our method against common open-set methods. We present the open-set **mAP** along with two definitions of the binary AUROC. **AUROC** follows the common definition where background detections are ignored, while $AUROC_{bd}$ treats background detections as OOD detections. We observe that our **Fusion** algorithm achieves the best AUROC in both definitions and datasets, with an average improvement of **2.7%** over the second best, while maintaining comparable or better open-set mAP performance. While the improvement seems incremental, our method maintains **robust performance** across both datasets in contrast to other methods.

Lastly, Table I shows that our two-class algorithm maintains a robust closed-set mAP performance, both in the AOT and in the AOT-C dataset, compared to the base model, while maintaining a **real-time** inference speed of **27 fps**. For direct comparison with the other open-set algorithms in Table I the

### D. Three-Class Setting (ID vs. OOD vs. Background)

We next evaluate our method in the three-class classification setting, where detections must be explicitly assigned to in-distribution, out-of-distribution, or background. To the best of our knowledge, no prior work has addressed this problem directly. We therefore present our results primarily as a new benchmark, while comparing against a simple two-threshold heuristic applied to standard confidence scores. This baseline partitions the score range into three intervals, providing a straightforward extension of binary open-set detection.

To ensure comparability in Tables I and IV, all methods are evaluated under the same decision protocol, with thresholds selected such that the probability of misclassifying OOD samples as either ID or background remains below 20%. We report macro pairwise AUROC as a threshold-independent measure of three-class separability, and open-set mAP, which reflects task-level detection performance. As shown, our fusion MLP achieves the highest AUROC on both Real Flights and COCO, with a substantial margin over all competing methods. While simple score-based baselines retain moderate discriminative power in the simpler Real Flights dataset, they consistently fail to generalize to the more complex COCO setting, where their AUROC scores collapse close to random chance.

This trend is also reflected in detection performance. Even the best competing technique, the score of the detector, only

(a) ID object classification: Airplane      (b) OOD object classification: Drone      (c) Background Classification
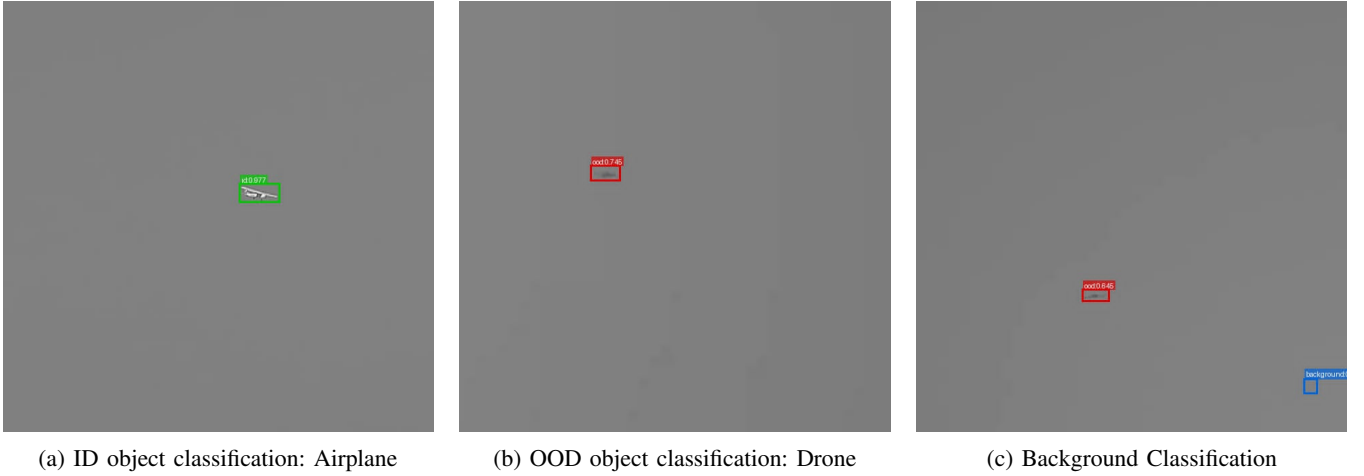
Fig. 2: Qualitative Results on Real Flights Dataset. ID classifications in green, OOD in red and background in blue. The UAV separates ood objects from background detections improving both safety and efficiency.

maintains a **49.3 mAP** in the closed setting, highlighting the impact of false positives from OOD and background confusion. By contrast, Table I demonstrates that our three-class fusion algorithm eliminates many of these errors, thereby improving mAP drastically to **58.7** while preserving robust open-set **mAP of 39.1** and a real-time inference speed of 27 FPS. Examples of the 3 class classification can be seen in Fig. 2.

### E. OOD Ablation Study

Training a detector such as RT-DETR requires powerful hardware and considerable computational resources. Moreover, all open-set recognition algorithms rely on access to OOD data that is representative of the deployment domain in order to calibrate thresholds. In practice, however, the point at which such OOD samples become available may not coincide with access to high-performance hardware. For this reason, we deliberately design our fusion MLP to be small enough to train on a standard CPU in just a few minutes.

An alternative to lightweight on-site training would be to avoid using in-domain OOD data altogether. To explore this possibility, we perform experiments on the Real Flights dataset in which the MLP is trained without access to its native OOD samples. Instead, we construct OOD training sets using three different strategies. We present them below, along with the AUROC achieved in the 2 class open-set problem:

1) Random detections from the COCO dataset, treated as OOD samples. ($AUROC = 0.867$)
2) Samples from an unrelated drone dataset. ($AUROC = 0.823$)
3) Artificially simulated features sampled outside the distributions of ID and background detections. ($AUROC = 0.835$)

We observe that none of the datasets achieves results better than Joint Thresholding [25], which achieved the second best results on the dataset. To further analyze the shortcomings of the method we use the best resulting dataset ($COCO$)

and analyze the 3 class open-set performance. The pairwise AUROCs between the 2 classes are ($id/ood, id/bg, ood/bg$) : $(0.722, 0.957, 0.951)$.

These results highlight that the distribution shift between proxy OOD datasets and deployment OOD data is too large for the MLP to generalize effectively, underscoring the need for OOD samples that better reflect the deployment domain. As future work, we therefore advocate the use of synthetic image datasets to create realistic OOD training samples and close the domain gap.

## V. Conclusions

We presented a lightweight post-hoc confidence fusion framework that enhances the robustness of aerial object detection in open-set conditions. In the two-class ID/OOD setting, our method achieves measurable improvements over thresholding baselines, yielding more reliable predictions under distribution shift.

Most importantly, we extend open-set detection to a three-class problem, explicitly distinguishing between in-distribution objects, out-of-distribution targets, and background clutter. Background detections, which are commonly ignored in other settings, are prevalent in autonomous navigation because the background is both complex and dominant in aerial imagery. For this reason, three-class classification is essential: unknown objects represent potential hazards that must be avoided, while background regions can be safely disregarded. By enabling this separation without compromising accuracy or real-time operation, our approach offers a practical step toward safer and more reliable robotic navigation in open and uncertain environments.

## References

[1] K. Wang, X. Fu, Y. Huang, C. Cao, G. Shi, and Z.-J. Zha, "Generalized uav object detection via frequency domain disentanglement," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 1064–1073.

[2] A. Wu and C. Deng, "Single-domain generalized object detection in urban scene via cyclic-disentangled self-distillation," in *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, 2022, pp. 847–856.

[3] J. Yang, K. Zhou, Y. Li, and Z. Liu, "Generalized out-of-distribution detection: A survey," *International Journal of Computer Vision*, vol. 132, no. 12, pp. 5635–5662, 2024.

[4] D. Miller, N. Sünderhauf, M. Milford, and F. Dayoub, "Uncertainty for identifying open-set errors in visual object detection," *IEEE Robotics and Automation Letters*, vol. 7, no. 1, pp. 215–222, 2021.

[5] Z. Zhou, Y. Yang, Y. Wang, and R. Xiong, "Open-set object detection using classification-free object proposal and instance-level contrastive learning," *IEEE Robotics and Automation Letters*, vol. 8, no. 3, pp. 1691–1698, 2023.

[6] K. Wang, C. Shen, X. Li, and J. Lu, "Uncertainty quantification for safe and reliable autonomous vehicles: A review of methods and applications," *IEEE Transactions on Intelligent Transportation Systems*, 2025.

[7] D. Reynolds, "Gaussian mixture models," in *Encyclopedia of biometrics*. Springer, 2015, pp. 827–832.

[8] Y. Zhao, W. Lv, S. Xu, J. Wei, G. Wang, Q. Dang, Y. Liu, and J. Chen, "DETRs Beat YOLOs on Real-time Object Detection," Seattle, Washington, USA, pp. 16 965–16 974, June 2024.

[9] X. Fan, Q. Wang, J. Ke, F. Yang, B. Gong, and M. Zhou, "Adversarially adaptive normalization for single domain generalization," in *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 8204–8213.

[10] S. Gasperini, J. Haug, M.-A. N. Mahani, A. Marcos-Ramiro, N. Navab, B. Busam, and F. Tombari, "Certainnet: Sampling-free uncertainty estimation for object detection," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 698–705, 2021.

[11] Y. Zheng, Z. Chen, D. Lv, Z. Li, Z. Lan, and S. Zhao, "Air-to-air visual detection of micro-uavs: An experimental evaluation of deep learning," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1020–1027, 2021.

[12] S. Ghosh, J. Patrikar, B. Moon, M. M. Hamidi, and S. Scherer, "Airtrack: Onboard deep learning framework for long-range aircraft detection and tracking," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 1277–1283.

[13] A. Arsenos, E. Petrongonas, O. Filippopoulos, C. Skliros, D. Kollias, and S. Kollias, "Nefeli: A deep-learning detection and tracking pipeline for enhancing autonomy in advanced air mobility," *Aerospace Science and Technology*, vol. 155, p. 109613, 2024.

[14] V. Karampinis, A. Arsenos, O. Filippopoulos, E. Petrongonas, C. Skliros, D. Kollias, S. Kollias, and A. Voulodimos, "Ensuring uav safety: A vision-only and real-time framework for collision avoidance through object detection, tracking, and distance estimation," in *2024 International Conference on Unmanned Aircraft Systems (ICUAS)*, 2024, pp. 1072–1079.

[15] A. Arsenos, V. Karampinis, E. Petrongonas, C. Skliros, D. Kollias, S. Kollias, and A. Voulodimos, "Common corruptions for evaluating and enhancing robustness in air-to-air visual object detection," *IEEE Robotics and Automation Letters*, vol. 9, no. 7, pp. 6688–6695, 2024.

[16] X. Du, Z. Wang, M. Cai, and Y. Li, "Vos: Learning what you don't know by virtual outlier synthesis," *arXiv preprint arXiv:2202.01197*, 2022.

[17] J. Li, P. Chen, Z. He, S. Yu, S. Liu, and J. Jia, "Rethinking out-of-distribution (ood) detection: Masked image modeling is all you need," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 11 578–11 589.

[18] X. Wu, J. Lu, Z. Fang, and G. Zhang, "Meta ood learning for continuously adaptive ood detection," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 19 353–19 364.

[19] A. Miyai, J. Yang, J. Zhang, Y. Ming, Y. Lin, Q. Yu, G. Irie, S. Joty, Y. Li, H. Li *et al.*, "Generalized out-of-distribution detection and beyond in vision language model era: A survey," *arXiv preprint arXiv:2407.21794*, 2024.

[20] R. Li, C. Zhang, H. Zhou, C. Shi, and Y. Luo, "Out-of-distribution identification: Let detector tell which i am not sure," in *European Conference on Computer Vision*. Springer, 2022, pp. 638–654.

[21] S. Wilson, T. Fischer, F. Dayoub, D. Miller, and N. Sünderhauf, "Safe: Sensitivity-aware features for out-of-distribution object detection," in *Proceedings of the ieee/cvf international conference on computer vision*, 2023, pp. 23 565–23 576.

[22] O. Zohar, K.-C. Wang, and S. Yeung, "Prob: Probabilistic objectness for open world object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 11 444–11 453.

[23] K. Joseph, S. Khan, F. S. Khan, and V. N. Balasubramanian, "Towards open world object detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 5830–5840.

[24] Z. Zhou, Y. Yang, Y. Wang, and R. Xiong, "Open-set object detection using classification-free object proposal and instance-level contrastive learning," *IEEE Robotics and Automation Letters*, vol. 8, no. 3, pp. 1691–1698, 2023.

[25] S. Loukovitis, A. Arsenos, V. Karampinis, and A. Voulodimos, "Model-agnostic open-set air-to-air visual object detection for reliable uav perception," 2025. [Online]. Available: https://arxiv.org/abs/2509.09297

[26] G. Jocher, A. Stoken, J. Borovec, NanoCode012, Christopher-STAN, L. Changyu, Laughing, tkianai, A. Hogan, lorenzomammana, yxNONG, AlexWang1900, L. Diaconu, Marc, wanghaoyang0106, ml5ah, Doug, F. Ingham, Frederik, Guilhen, Hatovix, J. Poznanski, J. Fang, L. Yu, changyu98, M. Wang, N. Gupta, O. Akhtar, PetrDvoracek, and P. Rai, "ultralytics/yolov5: v3.1 - bug fixes and performance improvements," https://doi.org/10.5281/zenodo.4154370, Oct 2020.

[27] G. Jocher, A. Chaurasia, and J. Qiu, "Ultralytics yolo," https://github.com/ultralytics/ultralytics, Jan 2023, [Online; accessed August 2, 2025].

[28] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, "Yolox: Exceeding yolo series in 2021," *arXiv preprint arXiv:2107.08430*, 2021.

[29] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.

[30] Y. Henon, "Pytorch-retinanet: Pytorch implementation of retinanet," 2020. [Online]. Available: https://github.com/yhenon/pytorch-retinanet

[31] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 28, 2015.

[32] sovit-123, "Faster r-cnn pytorch training pipeline," 2025. [Online]. Available: https://github.com/sovit-123/fasterrcnn-pytorch-training-pipeline

[33] S. Chen, P. Sun, Y. Song, and P. Luo, "Diffusiondet: Diffusion model for object detection," in *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023, pp. 19 773–19 786.

[34] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European conference on computer vision*. Springer, 2020, pp. 213–229.

[35] X. Zhou, V. Koltun, and P. Krähenbühl, "Probabilistic two-stage detection," in *arXiv preprint arXiv:2103.07461*, 2021.

[36] J. Mukhoti, A. Kirsch, J. Van Amersfoort, P. H. Torr, and Y. Gal, "Deep deterministic uncertainty: A new simple baseline," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 24 384–24 394.

[37] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.

[38] "Airborne object tracking dataset," https://registry.opendata.aws/airborne-object-tracking, accessed: 2023-07-23.