# Per-frame mAP Prediction for Continuous Performance Monitoring of Object Detection During Deployment

Quazi Marufur Rahman, Niko Sünderhauf and Feras Dayoub

*Abstract*— **Performance monitoring of object detection is crucial for safety-critical applications such as autonomous vehicles that operate under varying and complex environmental conditions. Currently, object detectors are evaluated using summary metrics based on a single dataset that is assumed to be representative of all future deployment conditions. In practice, this assumption does not hold, and the performance fluctuates as a function of the deployment conditions. To address this issue, we propose an introspection approach to performance monitoring during deployment without the need for ground truth data. We do so by predicting when the per-frame mean average precision drops below a critical threshold using the detector's internal features. We quantitatively evaluate and demonstrate our method's ability to reduce risk by trading off making an incorrect decision by raising the alarm and absenting from detection.**

## I. INTRODUCTION

Object detection is a crucial part of many safety-critical applications such as robotics and autonomous systems. For safe operation, an autonomous vehicle (AV), for example, needs to accurately locate and identify critical objects like other vehicles and pedestrians on the road. To achieve this goal, there is ongoing research [1]–[10] to improve the speed and accuracy of object detection models. However, due to the discrepancy between training data and deployment environments (i.e., dataset shift [11]) and many other unavoidable factors like sensor failure or degraded image quality, a consistent deployment performance can not be guaranteed. Hence, object detection accuracy can fluctuate without any prior notification while deployed on an autonomous vehicle. A silent failure like this in the object detection model is a significant concern. Due to this failure, the AV can cause catastrophic damage if it operates based on erroneous object detection. Undetected performance drops are a significant bottleneck for the widespread deployment of autonomous vehicles in our everyday lives. Hence, for safety, robustness, and reliability, the importance of performance monitoring of a deployed object detection model is paramount.

The standard practice to prepare an object detection model for deployment is to train and evaluate the model using training and evaluation split of some dataset to measure the accuracy and generalization capacity. Here, the assumption is the training and evaluation data are representative of the real operating environment. However, this assumption does not hold in the context of autonomous vehicles where the operating environment is continuously evolving and might

change unexpectedly. Consequently, object detection performance fluctuates without any prior notification. Moreover, the performance might drop below any critical threshold, which can cause a fatal incident. See Figure 1 for an overview.

One possible solution is to develop an exceptionally accurate and domain adaptive object detection system for autonomous vehicles. However, it is impossible in most practical circumstances to account for all imaginable future deployment conditions during training. Another approach is to identify when the performance of the deployed object detector drops below a critical threshold. So without the need to increase the detection accuracy directly, a performance drop identifier can protect the autonomous vehicle by providing crucial alerts during periods of silent failure. However, measuring the performance drop directly during deployment is impractical due to the absence of ground-truth data in this phase. Therefore, we advocate equipping object detectors with self-assessment capability to detect instances of performance drop during deployment.

Self-assessment is becoming a prerequisite for any vision-based efficient, safe, and robust robotic system. This capability is often referred to as introspective perception [12], [13]. For autonomous driving, an introspective object detection system can hand over the control to human drivers when it can predict inconsistency in its operation. There are several works [14]–[17] towards addressing the requirement of providing self-assessment in a deep learning based robotic vision system. However, there are very few works towards introspective systems for object detection. To this end, our paper makes the following contributions:

1) We propose an introspective approach to performance monitoring of object detection during deployment without access to ground truth labels.
2) We propose an internal integrated feature based on the mean, max and statistics pooling techniques for performance monitoring.
3) We introduce the use of per-frame mAP prediction for continuous performance monitoring of object detectors.

The rest of the paper is organized as follows: In Section II, we review the related works on introspective perception systems. In Section III, we introduce our framework to find the performance drop for an object detection system. Section IV outlines our experimental setup. Section V presents the results and finally in Section VI we draw conclusions and suggest areas for future work.

The authors are with the Australian Centre for Robotic Vision at Queensland University of Technology (QUT), Brisbane, QLD 4001, Australia. Contact: quazi.rahman@qut.edu.au
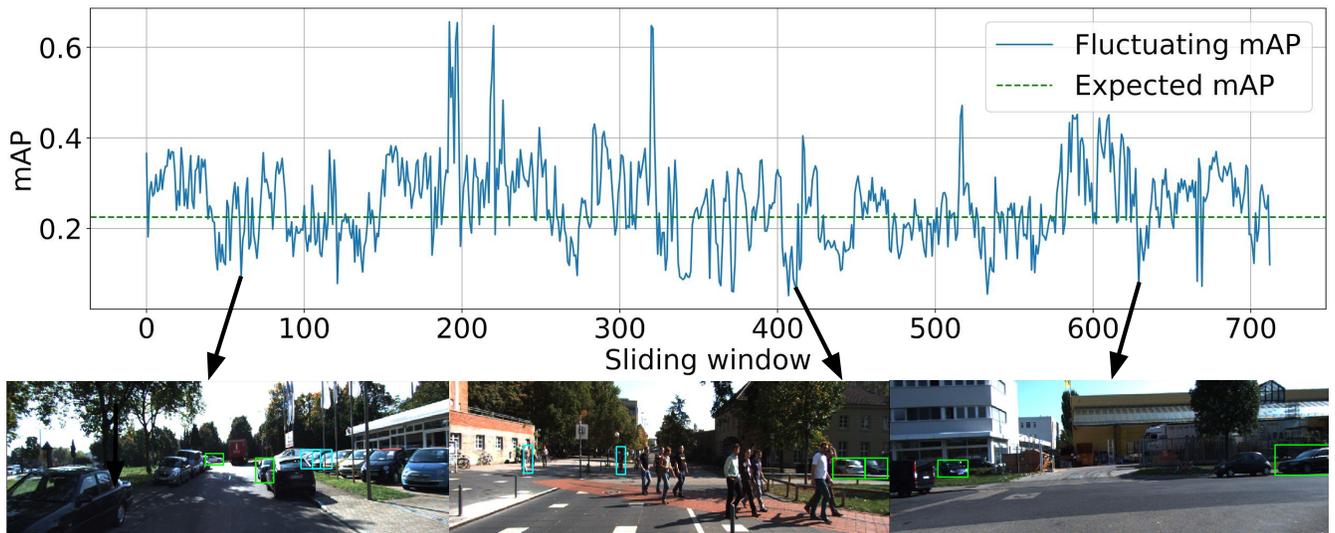
Fig. 1: First row, an example of object detection performance fluctuation during the deployment, tracked using a sliding window of ten frames. mAP is computed for the ten frames at a time. The dashed line represents a predefined critical threshold. We can see that mAP drops below this threshold from time to time. The second row shows some samples from the low mAP regions. Green and Cyan boxes represent false negative and false positive errors made by the object detector.

## II. RELATED WORK

In robotics, the idea of self-assessment was introduced by [13] to achieve reliable performance in a real environment. They described this self-assessment as the introspection capability of a mobile robot while operating in an unknown environment. Later [14] and [18] adopted the idea of introspection for classification and semantic mapping respectively in the context of robotics. These works examine the output of the underlying system to predict the likelihood of failure on any given input.

Another line of research is to predict the perception system performance from the input itself. In this paradigm, [19] introduced an evaluator algorithm to predict the failure of a human pose estimation model. Zhang *et al.* [20] introduced the terminology *basesys* and *alert* in failure prediction context. They proposed a general framework where *alert* is used to raise a warning when the underlying system *basesys* fails to make a correct decision from an input. Daftry *et al.* [12] proposed an introspective framework to predict an image classifier model failure deployed on a micro aerial vehicle. Following a similar methodology, [21] proposed a model to predict how hard an image is for an underlying classifier. Using a probabilistic model, [22] predicted the probable performance of a robot's perception system based on past experience in the same location.

Recently, confidence estimation and Bayesian approaches for uncertainty estimation have gained popularity to detect how well the underlying model has performed on the input. TrustScore [23], Maximum Class Probability [24] and True Class Probability [25] are some of the works that measure the confidence of the underlying model for a given task using the confidence estimation paradigm. Using a Bayesian approach, [26] proposed to use dropout as a Bayesian approximation

technique to represent model uncertainty. Following their work, [27], [28] have used dropout sampling to identify the quality of image and video segmentation network. Here, all of these works focus on predicting model failure using different approaches for classification and segmentation tasks.

In the context of object detection, [29], [30] have used different approaches to identify the failure of an object detection system. Both of these works are beneficial to identify false positive errors made by an object detector. Whereas, [31] and [32] have proposed different approaches to detect false negative errors made by an object detection model. Our proposed approach differs from these methods in that we can detect images with low per-frame mAP, which covers both false positive and false negative errors as well as poor object localization.

## III. APPROACH OVERVIEW

In this section we describe our proposed framework to predict the performance drop of an object detection system during deployment without using any ground-truth data. We assume that the deployed object detection model weight remains frozen during this phase. First of all, we will define the problem.

Assuming we have an object detector $O$ with backbone deep neural network $B$, $O$ is trained to detect a set of objects $T$ from a training dataset, $D_t$. We also have an evaluation dataset $D_e$, similarly distributed as $D_t$. $D_e$ contains a set of images $I = \{I_1, I_2 \ldots I_n\}$ and corresponding annotations $L = \{L_1, L_2 \ldots L_n\}$ per image. After the object detection training phase, $O$ is applied on $D_e$ to detect all the objects from $T$. Thus, we get a set of predictions per image, $P = \{P_1, P_2 \ldots P_n\}$. Using the pairs of annotations and predictions per image $(L_i, P_i)$, we compute the per-frame mAP, $M = \{M_1, M_2 \ldots M_n\}$ following the procedure at

[33]. Here, per-frame mAP quantifies $O$'s performance to detect all the existing objects in each image.

We assign each image of $D_e$ into *success* and *failure* classes using the Equation 1. Here $\lambda$ is chosen to be the $k^{th}$ percentile of $M$. The *failure* class contains the $k\%$ image frames from the $D_e$ where $O$ was not accurate enough to detect the available objects. The choice of $k$ here is application specific. We want to train the introspective perception system *alert* to predict the images similar to the *failure* class where per-frame mAP will be lower than $\lambda$.

$$\mathcal{L}(I) = \begin{cases} failure, & mAP_{\text{per-frame}} < \lambda \\ success, & \text{otherwise} \end{cases} \quad (1)$$

To train the *alert* we use features $F = \{F_1, F_2 \ldots F_n\}$ for each image from $D_e$. Following the failure prediction network proposed by [21] and [25], the final convolutional layer of backbone $B$ is used to extract all the necessary features. Assuming that, there are $N$ channels at the last layer of $B$ and each activation map is of size $W \times H$, we apply Equation 2 on the last layer to extract the mean pooling feature $F_{mean}$. Here, $f(x, y)$ represents the spatial unit of each activation map.

$$F_{mean} = \frac{\sum_{x=1}^{H} \sum_{y=1}^{W} f(x, y)}{W \cdot H} \quad (2)$$

Applying Equation 3 on the last layer of $B$, we generate the max pooling feature $F_{max}$.

$$F_{max} = \max_{x \in [1, H]} \max_{y \in [1, W]} f(x, y) \quad (3)$$

Inspired by [34], we calculate the standard deviation from each activation map to generate the statistics pooling feature $F_{std}$ following the Equation 4. Here $std(f_i)$ calculates the standard deviation of $i^{th}$ feature map.

$$F_{std} = std(f_1) \oplus std(f_2) \oplus \ldots std(f_N) \quad (4)$$

All the features described above are concatenated together to generate the feature $F_{mean\_max\_std}$ for the *alert* system. Equation 5 formulates this process.

$$F_{mean\_max\_std} = F_{mean} \oplus F_{max} \oplus F_{std} \quad (5)$$

We train a binary classifier using the $F_{mean\_max\_std}$ feature and the corresponding labels from Equation 5 and Equation 1 respectively. The classifier is trained to predict the probability of an image feature to be in the *failure* class. Following [20], we will refer to the object detection model and its corresponding binary classifier as *basesys* and *alert* respectively. Figure 2 shows the incorporation between the *basesys* and *alert* system.

## IV. Experimental Setup

In this section we will describe the settings that we have used to train the *basesys* and *alert* system.
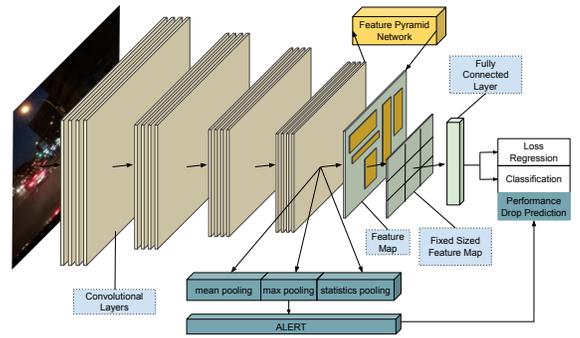


Fig. 2: The architecture of our proposed *basesys* and *alert* system together. The last convolutional feature of the backbone is pooled using the mean, max and statistical pooling layer to generate the feature for *alert*. *Alert* consists of a binary classifier that predict the performance drop of the *basesys*.

*a) Datasets::* We used all images from *kitti* [35] dataset and one frame per video from *bdd* [36] dataset to train both *basesys* and corresponding *alert* system. Randomly selected 60%, 20% and 20% images from both dataset have been used for *basesys* training, evaluation and testing purpose. As person and car classes are available in both of these dataset, we used these two objects as *basesys* target class. After the object detection training, *basesys* is used to detect person and car from the 20% evaluation split. Based on *basesys* performance on the evaluation split, we collect image features and labels for the *alert* system following the procedure described in Section III. Thus the features and labels collected from *basesys* evaluation split works as training dataset for the corresponding *alert* system. To test the *alert*, we first apply *basesys* on the testing split and measure its per-frame performance drop, which works as the testing data for the *alert* system. In some of our experiments, we will train and test *basesys* and *alert* using training and testing split of a single dataset. We will refer this settings as *similar dataset*. In rest of the experiments, *basesys* and *alert* will be trained using training split of one dataset and tested using testing split of another dataset. This arrangements will be referred as *cross dataset* settings.

TABLE I: *Basesys* mean average precision (mAP) using ResNet18 and ResNet50 backbone. Here *basesys* is trained and tested using *similar dataset* and *cross dataset* settings. *Basesys* accuracy drops when trained and tested on different dataset.

| ResNet18 | | testing dataset | | ResNet50 | | testing dataset | |
|---|---|---|---|---|---|---|---|
| | | kitti | bdd | | | kitti | bdd |
| training dataset | kitti | 0.292 | 0.130 | training dataset | kitti | 0.377 | 0.182 |
| | bdd | 0.200 | 0.331 | | bdd | 0.259 | 0.499 |

*b) Basesys Training::* We have used Faster RCNN object detection network [2] as the *basesys* in all of our experiments. *Basesys* has been trained using transfer learning to detect person and car object from both *kitti* and *bdd* dataset. Two different versions of Residual Neural Network

TABLE II: Area Under the Precision Recall Curve (AUPRC) and Area Under the ROC Curve (AUROC) score for *alert* system in the similar dataset settings. Here *alert* is used to identify *basesys* performance drop in a known environment. The notation *A/B/C* denotes that *basesys* and *alert* is trained on dataset *A* using backbone *C* and *alert* is used to identify *basesys* performance drop on dataset*B*.

| | kitti/kitti/18 | | kitti/kitti/50 | | bdd/bdd/18 | | bdd/bdd/50 | |
|---|---|---|---|---|---|---|---|---|
| Feature | AUPRC | AUROC | AUPRC | AUROC | AUPRC | AUROC | AUPRC | AUROC |
| n_proposals | 0.180 | 0.128 | 0.186 | 0.110 | 0.205 | 0.368 | 0.197 | 0.363 |
| mean_conf_score | 0.205 | 0.320 | 0.192 | 0.358 | 0.452 | 0.653 | 0.463 | 0.665 |
| classifier | 0.728 | 0.851 | 0.689 | 0.831 | 0.498 | 0.744 | 0.488 | 0.734 |
| places365 | 0.654 | 0.799 | 0.670 | 0.823 | 0.516 | 0.753 | 0.507 | 0.744 |
| layer | 0.760 | 0.890 | 0.480 | 0.764 | 0.622 | 0.814 | 0.587 | **0.798** |
| mean | 0.738 | 0.876 | 0.602 | 0.822 | 0.587 | 0.800 | 0.549 | 0.777 |
| max | 0.756 | 0.887 | 0.673 | 0.819 | 0.621 | 0.811 | 0.587 | 0.790 |
| mean_std | 0.747 | 0.879 | 0.689 | **0.855** | 0.609 | 0.815 | 0.577 | 0.791 |
| mean_max | 0.777 | 0.898 | 0.708 | 0.841 | 0.627 | 0.818 | 0.587 | 0.793 |
| mean_max_std | **0.781** | **0.902** | **0.712** | 0.846 | **0.633** | **0.820** | **0.595** | 0.795 |

TABLE III: Area Under the Precision Recall Curve (AUPRC) and Area Under the ROC Curve (AUROC) for *alert* in the cross dataset settings. Here *alert* is identifying *basesys* performance drop in an unknown environment. The notation *A/B/C* denotes that *basesys* and *alert* is trained on dataset *A* using backbone *C* and *alert* is used to identify *basesys* performance drop on dataset *B*.

| | bdd/kitti/18 | | bdd/kitti/50 | | kitti/bdd/18 | | kitti/bdd/50 | |
|---|---|---|---|---|---|---|---|---|
| Feature | AUPRC | AUROC | AUPRC | AUROC | AUPRC | AUROC | AUPRC | AUROC |
| n_proposals | 0.558 | 0.381 | 0.736 | 0.447 | 0.531 | 0.465 | 0.538 | 0.497 |
| mean_conf_score | 0.641 | 0.507 | 0.786 | 0.566 | 0.675 | 0.623 | 0.730 | 0.710 |
| classifier | 0.783 | 0.629 | 0.818 | 0.483 | 0.672 | 0.663 | 0.557 | 0.578 |
| places365 | 0.781 | 0.624 | 0.821 | 0.493 | 0.857 | 0.837 | 0.630 | 0.637 |
| layer | 0.780 | 0.684 | 0.815 | 0.580 | 0.868 | 0.836 | 0.662 | 0.670 |
| mean | 0.754 | 0.665 | 0.809 | 0.563 | 0.858 | 0.831 | 0.733 | 0.753 |
| max | 0.778 | 0.682 | 0.813 | 0.582 | 0.647 | 0.605 | 0.661 | 0.686 |
| mean_std | 0.759 | 0.672 | 0.815 | 0.569 | 0.855 | 0.823 | 0.751 | 0.768 |
| mean_max | 0.786 | 0.692 | 0.822 | 0.586 | 0.826 | 0.786 | 0.701 | 0.726 |
| mean_max_std | **0.790** | **0.696** | **0.822** | **0.587** | **0.883** | **0.856** | **0.833** | **0.825** |

[37], ResNet18 and ResNet50 have been used as the *basesys* backbone. In our experiments, the *basesys*, trained using RestNet50 backbone has performed better than the ResNet18 backbone. Table I shows comparative performance using the mean average precision (mAP) for all different *basesys* and dataset combinations.

*c) Feature Collection::* We experimented with multiple features to find the most suitable one for the proposed *alert* system. The first set of features are collected from *basesys* bounding box proposals.

- *mean_conf_score*: This feature exploits object proposal confidence score to determine *basesys* performance drop. As *basesys* proposes multiple bounding boxes with corresponding confidence scores and labels during object detection, we use the mean of confidence scores which are greater than $0.5$ to build the first feature. Here, a lower mean confidence score indicates a potential performance drop in the *basesys*.
- *n_proposals*: We assume that a crowded environment might be a factor for *basesys* performance drop. To evaluate this assumption, we used the number of proposals having a confidence score greater than $0.5$ as a performance drop indicator.

The second set of features are collected from two external deep convolutional neural networks.

- *classifier:* Two different versions of Residual neural

network, Resnet18 and ResNet50 have been used to extract image features to train the *alert* system. Both of these networks are pre-trained on ImageNet [38] dataset.
- *places365:* We used ResNet18 and ResNet50 network pre-trained on Places365 [39] dataset to extract features to train the *alert* system.

In both cases, average pooling has been used at the final convolutional layer to extract the necessary image features.

We use the *basesys* backbone to extract the third set of features. These will be referred as the internal features.

- *layer:* We applied the mean-pooling operation in all of the convolutional layers of the backbone and concatenated them to create this feature.
- *mean*, *max* and *std*: Applying the mean, max and statistics pooling technique described in Section III at the last convolutional layer of *basesys* backbone, we extracted the *mean*, *max* and *std* features.
- *mean_std* and *mean_max:* Using the concatenation operation and following the feature generation technique proposed in [34] and [40], we generate two new features *mean_std* and *mean_max* using the *mean*, *max* and *std* feature.
- *mean_max_std:* This feature is generated by applying the Equation 5 at the last convolutional layer of *basesys*
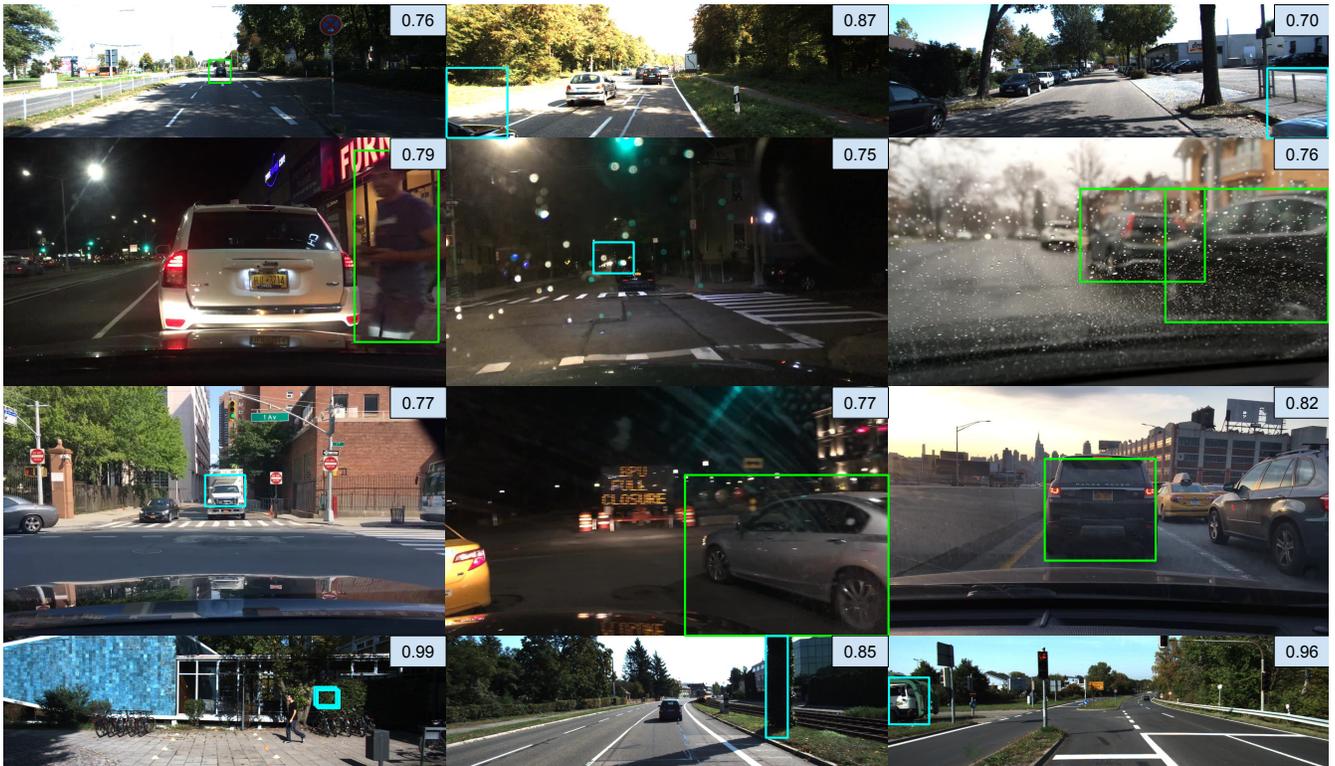
Fig. 3: Examples of *alert* prediction to identify *basesys* performance drop. Here the Green and Cyan bounding boxes show the false negative and false positive errors respectively made by an object detector. *Alert* prediction is showed at the upper right corner of each image. The first row shows samples from the *kitti/kitti/50* experimental settings. The second, third and fourth row show samples from the *bdd/bdd/18*, *kitti/bdd/18* and *bdd/kitti/50* experiments.

backbone.

*d)* ***Alert Training:*** We used a multi layer fully connected binary classifier with $50\%$ dropout rate to train all the *alert* systems. Besides, we used binary cross entropy loss with balanced sampling to train the *alert* network.

## V. EVALUATION AND RESULTS

### A. AUPRC and AUROC Metrics

This section summarizes the *alert* accuracy using Area Under the Precision Recall Curve (AUPRC) and Area Under the ROC Curve (AUROC) metric. Here, we will refer all our experimental settings using the notation *A/B/C*. It means the *basesys* and *alert* are trained on dataset *A* using backbone *C* and *alert* is used to identify *basesys* performance drop on dataset *B*. Here *C* can be *18* or *50*, resembling the ResNet18 and ResNet50 backbone for the *basesys*.

Table II summarizes the *alert* accuracy for similar dataset settings. Our proposed *mean_max_std* feature achieves $0.781$ and $0.902$ as AUPRC and AUROC score, and outperforms all other features in the case of *kitti/kitti/18*. For *kitti/kitti/50*, *bdd/bdd/18* and *bdd/bdd/50* experimental settings, features collected from the *basesys* performs better than all other features in terms of AUPRC and AUROC score.

The proposed *alert* system is beneficial for cross dataset settings too. Table III shows the AUPRC and AUROC

scores for *alert* when it is used to identify *basesys* performance when deployed on an unknown environment. For *bdd/kitti/18* settings, *alert* achieves $0.790$ and $0.696$ as AUPRC and AUROC score respectively when used with the *mean_max_std* feature. In all cross dataset experimental settings, *mean_max_std* features outperforms all other features for identifying *basesys* performance drop.

TABLE IV: The true warning rate of the *alert* system to identify *basesys* performance drop.

| ResNet18 | | testing dataset | | ResNet50 | | testing dataset | |
|---|---|---|---|---|---|---|---|
| | | kitti | bdd | | | kitti | bdd |
| training | kitti | 59.1% | 81.8% | training | kitti | 66.0% | 78.7% |
| dataset | bdd | 79.3% | 55.4% | dataset | bdd | 81.4% | 52.4% |

### B. True Warning Rate

Using the best performing feature, *mean_max_std*, we use the true warning rate metric to determine the quality of the *alert* system in raising a warning against *basesys* performance drop. Here, warning rate is the ratio of correctly raised warning vs the total number of frames with per-frame mAP below the critical threshold. Table IV shows the true warning rate raised by *alert* system.

The results in Table IV show that in *cross dataset* settings the true warning rate is higher than the *similar dataset* settings. As *basesys* accuracy drops in cross dataset settings
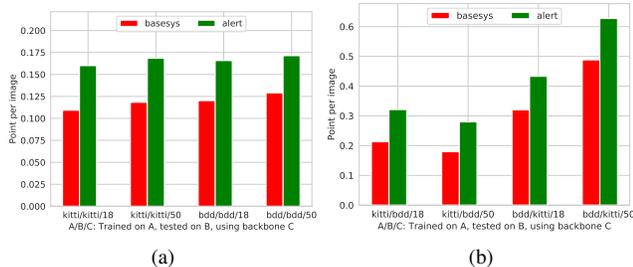
Fig. 4: Risk-Averse Metric for the proposed *alert* system. (a) Point per image earned by *basesys* with and without considering *alert* warning when trained and tested on similar dataset. (b) Point per image for *basesys* with and without *alert* system when trained and tested on different dataset. In both cases, *basesys* earns more point per image when associated with *alert*.
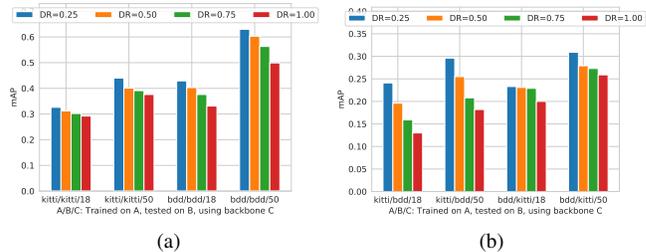


Fig. 5: mAP vs declaration rate metric for the proposed *alert* system. We used four different declaration rate to calculate the corresponding mAP metric. An increasing declaration rate shows a gradual drop in the mean average precision. (a) shows the mAP vs DR metric for the similar dataset settings (b) mAP vs DR metric for cross dataset settings.

(Table I), *alert* becomes more useful in these cases by identifying the critical cases. When the detector with ResNet50 backbone is trained on BDD and tested on Kitti, *alert* can identify $81.4\%$ of the frames where *basesys* per-frame mAP is lower than the critical threshold. Figure 3 displays multiple samples of *alert* raising the alarm and flagging frames where *basesys* performance drop below a critical threshold of $0.5$. The frames show conditions such as night, rain, cluttered environments.

### C. Risk-Averse Metric

In Risk-Averse Metric (RAM) [20] we evaluate *alert*'s capability to trade-off the risk of making an incorrect decision with not making a decision at all. RAM gives *basesys* $+1.0$ and $-0.5$ respectively for a correct and incorrect prediction. *basesys* will get $0$ point if it does not make any decision considering the warning raised by *alert*. For crucial system like self-driving car's object detection we expect *basesys* to trade-off incorrect decision for no decision. In such case, *basesys* can handover its control to some more competent systems. Figure 4a shows the point per image earned by *basesys* when it operates with and without considering the warning raised by *alert* in similar dataset settings. In all cases, *basesys* earns more point per image if it abstains from making an incorrect decision taking *alert*'s warning in consideration. Figure 4b shows the RAM metric for corss dataset settings.

### D. mAP vs Declaration Rate Metric

In this section we evaluate the *basesys* accuracy score for different declaration rate (DR) [20]. Here, declaration rate is the proportion of images on which *basesys* operates. The rest of the images are discard assuming that *basesys* per-frame mAP will be lower than the critical threshold on those images. To calculate this metric we first sort the images in the ascending order of *alert* confidence. Next, mAP of top DR percentage of images are computed to plot the mAP vs DR metric. For a perfect *alert* the mAP for low DR

images would be very high and decrease gracefully as DR approaches to $1.0$. In Figure 5a we show the mAP score for four different declaration rate in similar dataset settings. The mAP score drops gradually with the increasing declaration rate. Figure 5b shows the mAP vs DR metric for cross dataset settings. In both cases, we use *mean_max_std* features in *alert* to identify *basesys* performance drop.

## VI. CONCLUSION

Deep learning-based object detection is a critical component of a wide variety of robotic applications, from autonomous vehicle to warehouse automation due to its accuracy and efficiency. However, its performance is a function of the deployment conditions and could drop below a critical threshold leading to increased risk. Although there is always room to improve accuracy and speed, safety is still a significant concern that should not be overlooked. To this end, we presented an introspection approach to performance monitoring of deep learning based object detection. We showed that our approach can improve safety by raising an alarm when per-frame mean average precision is detected to drop below a critical. We also showed that internal features from the detector itself could be used to predict when per-frame mAP degrade. Our results showed quantitatively the ability of our method to reduce risk by trading off making an incorrect detection with raising the alarm and absenting from detection.

## REFERENCES

[1] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *Proceedings of the IEEE international conference on computer vision*, 2019, pp. 9627–9636.

[2] 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*, 2015, pp. 91–99.

[3] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.

[4] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian, "Centernet: Keypoint triplets for object detection," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 6568–6577, 2019.

[5] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *ArXiv*, vol. abs/2004.10934, 2020.

[6] K. He, R. B. Girshick, and P. Dollár, "Rethinking imagenet pre-training," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4917–4926, 2019.

[7] Z. Cai and N. Vasconcelos, "Cascade r-cnn: Delving into high quality object detection," *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6154–6162, 2018.

[8] Z. Li, C. Peng, G. Yu, X. Zhang, Y. Deng, and J. Sun, "Detnet: A backbone network for object detection," *ArXiv*, vol. abs/1804.06215, 2018.

[9] T.-Y. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2999–3007, 2017.

[10] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. M. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," *ArXiv*, vol. abs/2005.12872, 2020.

[11] J. Quionero-Candela, M. Sugiyama, A. Schwaighofer, and N. D. Lawrence, *Dataset shift in machine learning*. The MIT Press, 2009.

[12] S. Daftry, S. Zeng, J. A. Bagnell, and M. Hebert, "Introspective perception: Learning to predict failures in vision systems," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 1743–1750.

[13] A. C. Morris, "Robotic introspection for exploration and mapping of subterranean environments," Ph.D. dissertation, Carnegie Mellon University, The Robotics Institute, 2007.

[14] H. Grimmett, R. Triebel, R. Paul, and I. Posner, "Introspective classification for robot perception," *The International Journal of Robotics Research*, vol. 35, no. 7, pp. 743–762, 2016.

[15] H. Hu and G. Kantor, "Introspective evaluation of perception performance for parameter tuning without ground truth," in *Robotics: Science and Systems*, 2017.

[16] S. Rabiee and J. Biswas, "Ivoa: Introspective vision for obstacle avoidance," *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 1230–1235, 2019.

[17] S. Mohseni, M. Pitale, V. Singh, and Z. Wang, "Practical solutions for machine learning safety in autonomous vehicles," in *SafeAI@AAAI*, 2020.

[18] R. Triebel, H. Grimmett, R. Paul, and I. Posner, "Driven learning for driving: How introspection improves semantic mapping," in *ISRR*, 2013.

[19] N. Jammalamadaka, A. Zisserman, M. Eichner, V. Ferrari, and C. V. Jawahar, "Has my algorithm succeeded? an evaluator for human pose estimators," in *ECCV*, 2012.

[20] P. Zhang, J. Wang, A. Farhadi, M. Hebert, and D. Parikh, "Predicting failures of vision systems," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 3566–3573.

[21] P. Wang and N. Vasconcelos, "Towards realistic predictors," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 36–51.

[22] C. Gurau, D. Rao, C. H. Tong, and I. Posner, "Learn from experience: Probabilistic prediction of perception performance to avoid failure," *The International Journal of Robotics Research*, vol. 37, pp. 981 – 995, 2018.

[30] B. Cheng, Y. Wei, H. Shi, R. S. Feris, J. Xiong, and T. S. Huang, "Decoupled classification refinement: Hard false positive suppression for object detection," *ArXiv*, vol. abs/1810.04002, 2018.

[23] H. Jiang, B. Kim, and M. R. Gupta, "To trust or not to trust a classifier," in *NeurIPS*, 2018.

[24] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," *ArXiv*, vol. abs/1610.02136, 2017.

[25] C. Corbière, N. Thome, A. Bar-Hen, M. Cord, and P. Pérez, "Addressing failure prediction by learning model confidence," in *Advances in Neural Information Processing Systems*, 2019, pp. 2902–2913.

[26] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *ICML*, 2016.

[27] T. Devries and G. W. Taylor, "Leveraging uncertainty estimates for predicting segmentation quality," *ArXiv*, vol. 1807.00502, 2018.

[28] P.-Y. Huang, W. T. Hsu, C.-Y. Chiu, T.-F. Wu, and M. Sun, "Efficient uncertainty estimation for semantic segmentation in videos," in *ECCV*, 2018.

[29] D. Miller, L. Nicholson, F. Dayoub, and N. Sünderhauf, "Dropout sampling for robust object detection in open-set conditions," *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1–7, 2018.

[31] Q. M. Rahman, N. Sünderhauf, and F. Dayoub, "Did you miss the sign? a false negative alarm system for traffic sign detectors," *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 3748–3753, 2019.

[32] M. S. Ramanagopal, C. Anderson, R. Vasudevan, and M. Johnson-Roberson, "Failing to learn: Autonomously identifying perception failures for self-driving cars," *IEEE Robotics and Automation Letters*, vol. 3, pp. 3860–3867, 2018.

[33] 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.

[34] D. Snyder, D. Garcia-Romero, D. Povey, and S. Khudanpur, "Deep neural network embeddings for text-independent speaker verification." in *Interspeech*, 2017, pp. 999–1003.

[35] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.

[36] F. Yu, W. Xian, Y. Chen, F. Liu, M. Liao, V. Madhavan, and T. Darrell, "Bdd100k: A diverse driving video database with scalable annotation tooling," *arXiv preprint arXiv:1805.04687*, vol. 2, no. 5, p. 6, 2018.

[37] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[38] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[39] B. Zhou, À. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, pp. 1452–1464, 2018.

[40] L. Zhang and Y. Guo, "Delving into fully convolutional networks activations for visual recognition," in *ICMIP 2018*, 2018.