

# Probabilistic Object Detection: Definition and Evaluation

David Hall<sup>1,2</sup> Feras Dayoub<sup>1,2</sup> John Skinner<sup>1,2</sup> Haoyang Zhang<sup>1,2</sup> Dimity Miller<sup>1,2</sup>  
Peter Corke<sup>1,2</sup> Gustavo Carneiro<sup>1,3</sup> Anelia Angelova<sup>4</sup> Niko Sünderhauf<sup>1,2</sup>

<sup>1</sup>Australian Centre for Robotic Vision

<sup>2</sup>Queensland University of Technology <sup>3</sup>University of Adelaide <sup>4</sup>Google Brain

<sup>2</sup>{d20.hall, feras.dayoub, j6.skinner, haoyang.zhang.acrv,  
d24.miller, peter.corke, niko.suenderhauf}@qut.edu.au

<sup>3</sup>gustavo.carneiro@adelaide.edu.au <sup>4</sup>anelia@google.com

## Abstract

We introduce *Probabilistic Object Detection*, the task of detecting objects in images and accurately quantifying the spatial and semantic uncertainties of the detections. Given the lack of methods capable of assessing such probabilistic object detections, we present the new *Probability-based Detection Quality measure (PDQ)*. Unlike AP-based measures, PDQ has no arbitrary thresholds and rewards spatial and label quality, and foreground/background separation quality while explicitly penalising false positive and false negative detections. We contrast PDQ with existing mAP and moLRP measures by evaluating state-of-the-art detectors and a Bayesian object detector based on Monte Carlo Dropout. Our experiments indicate that conventional object detectors tend to be spatially overconfident and thus perform poorly on the task of probabilistic object detection. Our paper aims to encourage the development of new object detection approaches that provide detections with accurately estimated spatial and label uncertainties and are of critical importance for deployment on robots and embodied AI systems in the real world.

## 1. Introduction

Visual object detection provides answers to two questions: *what* is in an image and *where* is it? State-of-the-art approaches that address this problem are based on deep convolutional neural networks (CNNs) that localise objects by predicting a bounding box, and providing a class label with a confidence score, or a full label distribution, for every detected object in the image [27, 37, 38]. The ability of deep CNNs to quantify epistemic and aleatoric uncertainty [19] has recently been identified as paramount for deployment in safety critical applications, where the perception and decision making of an agent has to be trusted [1, 19, 43, 49].

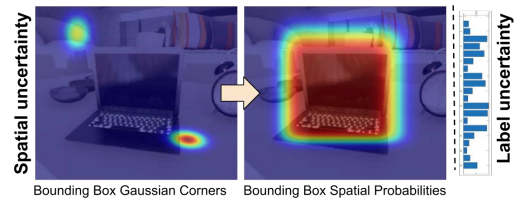


Figure 1: In contrast to conventional object detection, *probabilistic* object detections express semantic *and* spatial uncertainty. Our probabilistic object detections represent object locations as probabilistic bounding boxes where corners are modelled as 2D Gaussians (left) used to express a spatial uncertainty over the pixels (centre). Semantic uncertainty is represented as full label probability distributions (right).

While state-of-the-art object detectors have limited capability to express epistemic and aleatoric uncertainty about the class label through the confidence score or label distribution [14, 15, 33, 44, 48], uncertainty about the spatial aspects of the detection is currently not at all quantified. Furthermore, none of the existing benchmarks using average precision (AP) as the basis for their evaluation [7, 8, 20, 23, 26, 40] can evaluate how well detectors quantify spatial and semantic uncertainties.

We introduce **Probabilistic Object Detection**, the task of detecting objects in images while accurately quantifying the spatial and semantic uncertainties of the detections. Probabilistic Object Detection poses a key challenges that goes beyond the established conventional object detection: the detector must quantify its *spatial uncertainty* by reporting *probabilistic* bounding boxes, where the box corners are modelled as normally distributed. As illustrated in Figure 1, this induces a spatial probability distribution over the image for each detection. The detector must also reliably quantify its *semantic uncertainty* by providing a full probability distribution over the known classes for each detection.

To evaluate how well detectors perform on this chal-

lenging task, we introduce a new evaluation measure, **Probability-based Detection Quality** (PDQ). In contrast to AP-based measures, PDQ explicitly evaluates the reported probability of the true class via its *Label Quality* component. Furthermore, PDQ contains a *Spatial Quality* term that evaluates how well a detection’s spatial probability distribution matches the true object.

Unlike existing measures such as mAP [26] and moLRP [34], PDQ jointly evaluates spatial and label uncertainty quality, foreground and background separation quality, and the number of true positive (correct), false positive (spurious), and false negative (missed) detections. Importantly, PDQ does not rely on fixed thresholds or tuneable parameters, and provides optimal assignments of detections to ground-truth objects. Although PDQ has been primarily developed to evaluate *new* types of probabilistic object detectors that are designed to quantify spatial and semantic uncertainties, PDQ can also evaluate conventional state-of-the-art, non-probabilistic detectors.

As we show in Section 7, current conventional detection methods perform poorly on the task of probabilistic object detection due to spatial over-confidence and are outperformed by a recently proposed probabilistic object detector that incorporates Monte Carlo Dropout into a VGG16-based Single Shot MultiBox Detector (SSD) [29].

In summary, our contributions include defining the challenging new task of probabilistic object detection, introducing the new evaluation measure PDQ, evaluating current object detectors, and showing for the first time that novel probabilistic object detectors achieve better performance on this new task, that is highly relevant for applications such as robotics or embodied AI.

## 2. Motivation

Object detection embedded in a robot or autonomous system, such as a self-driving car, is part of a complex, active, goal-driven system. In such a scenario, object detection provides crucial perception information that ultimately determines the performance of the robot in its environment. Mistakes in object detection could lead to catastrophic outcomes that not only risk the success of the robot’s mission, but potentially endanger human lives [1, 2, 24, 32, 35, 49].

For safe and trusted operation in robots or autonomous systems, CNNs must express meaningful *uncertainty* information [1, 14, 15, 19, 43, 49]. Object detectors will have to quantify uncertainty for both the reported labels and bounding boxes, which would enable them to be treated as yet another sensor within the established and trusted framework of Bayesian information fusion [39, 47]. However, while state-of-the-art object detectors report at least an *uncalibrated* indicator of label uncertainty via label distributions or label scores [14, 15, 33, 44, 48], they currently do *not* report spatial uncertainty. As a result, eval-

uating the quality of the label or spatial uncertainties is not within the scope of typical benchmark measures and competitions [7, 8, 20, 23, 26, 40].

We argue in favour of accurate quantification of spatial and semantic uncertainties for object detectors in computer vision and robotics applications. Our work builds on this idea by creating a measure that will guide research towards developing detection systems that can operate effectively within a robot’s sensor fusion framework.

## 3. Related Work

**Conventional Object Detection:** Object detection is a fundamental task in computer vision and aims to localise each instance of certain object classes in an image using a bounding box. The typical output from an object detection system is a set of bounding boxes with a class label score [5, 9, 50]. Since the advent of convolutional neural networks (CNNs) [21], object detection has experienced impressive progress in terms of accuracy and speed [4, 12, 13, 25, 27, 37, 38]. Nonetheless, current overconfident object detection systems fail to provide spatial and semantic uncertainties, and as a result, can be a source of risk in various vision and robotics applications. The probabilistic object detection task introduced by this paper requires that object detectors estimate the spatial and semantic uncertainty of their detections.

**Uncertainty Estimation:** To improve system robustness and accuracy or avoid risks, quantifying uncertainty has become popular in many vision tasks. Kendall et al. [18] propose a Bayesian model that outputs a pixel-wise semantic segmentation with a measure of model uncertainty for each class. In [19], the authors propose to model the aleatoric and epistemic uncertainties for the pixel-wise semantic segmentation and depth regression tasks, and argue that epistemic uncertainty is important for safety-critical applications and training with small data sets. Kampffmeyer et al. [17] propose a model that estimates pixel-wise classification uncertainty in urban remote sensing images – they argue that the estimated uncertainty can indicate the correctness of pixel labelling. Miller et al. [29, 30] estimate both spatial and classification uncertainties for object detection and use the uncertainty to accept or reject detections under open-set conditions. Nair et al. [31] provide four different voxel-based uncertainty measures for their 3D lesion segmentation system to enable a more complete revision by clinicians. In [45] an uncertainty map for super-resolution of diffusion MR brain images is generated to enable a risk assessment for the clinical use of the super-resolved images. In [42], the authors build an ensemble of predictors to estimate the uncertainty of the centre of nuclei in order to produce more accurate classification results. All the methods above, except the last one [42], estimate uncertainty based on the Monte Carlo (MC) dropout technique [10, 11]. The

papers above provide evidence that it is important to estimate uncertainty for various vision tasks. Most of the proposed methods, except [29, 42], deal with pixel-wise classification. We argue that it is essential to capture the uncertainty of object detectors as motivated in Section 2.

**Performance Measures:** For the past decade, detection algorithms have predominantly been evaluated using average precision (AP) or variants thereof. Average precision was introduced for the PASCAL VOC challenge [8] in 2007 to replace measuring the area under the ROC curve. It is the average of the maximum precision values at different recall values. These use a pre-defined threshold for the intersection over union (IoU), typically 0.5, defining a true positive detection. This is calculated and averaged across all classes. Since then, AP has become the standard evaluation measure in the PASCAL VOC challenge and is the basis for many other works examining object detection [4, 25, 26, 27, 38, 40]. Most recently, a variation of AP was created which averages AP over multiple IoU thresholds (varying from 0.5 to 0.95 in intervals of 0.05) [26]. This averaging over IoUs rewards detectors with better localisation accuracy. In this work we refer to this measure as mean average precision (mAP) to distinguish it from AP despite mAP typically referring to averaging AP over all classes.

AP-based measures have biased the community to develop object detectors with high recall rate and localisation precision, but these measures have several weaknesses [3, 16, 36]. They rely on fixed IoU thresholds which can lead to overfitting for certain IoU thresholds – the negative consequence is that a small change in the thresholds can cause abrupt score changes. Additionally, these measures use the label score as the detection ranking evidence, without considering the spatial quality, which can lead to sub-optimal detection assignment. In our work, we propose the new evaluation measure PDQ to evaluate both label and spatial qualities of object detections, without using any fixed thresholds and relying on an optimal assignment of detection to ground-truth objects based on both spatial and label qualities.

Oksuz et al. [34] propose the Localisation Recall Precision (LRP) metric to overcome two main deficiencies of mAP: the inability to distinguish different precision-recall (PR) curves, and the lack of a direct way to measure bounding box localisation accuracy. When used for analysing multi-class detectors, the mean optimal LRP (moLRP) is used. Comparing to mAP, moLRP is also based on PR curves but measures localisation quality, false positive rate and false negative rate at some optimal label threshold for each class. The localisation quality is represented by the IoU between the detection and the ground-truth object, scaled by the IoU threshold being used to plot the PR curves. In contrast, our PDQ measure estimates the spatial

uncertainty through probabilistic bounding boxes and evaluates how well the detection bounding box’s spatial probability distribution coincides with the true object.

## 4. Probabilistic Object Detection

Probabilistic Object Detection is the task of detecting objects in an image, while accurately quantifying the spatial and semantic uncertainties of the detections. Probabilistic Object Detection thus extends conventional object detection, and makes the quantification of uncertainty an essential part of the task and its evaluation.

Probabilistic Object Detection requires a detector to provide for each known object in an image:

- a categorical distribution over all class labels, and
- a bounding box represented as  $\mathcal{B} = (\mathcal{N}_0, \mathcal{N}_1) = (\mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0), \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1))$  such that  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\Sigma}_i$  are the mean and covariances for the multivariate Gaussians describing the top-left and bottom-right corner of the box.

From this probabilistic box representation  $\mathcal{B}$ , we can calculate a probability distribution  $P$  over all pixels  $(u', v')$ , such that  $P(u', v')$  is the probability that the pixel is contained in the box:

$$P(u', v') = \iint_{0,0}^{v',u'} \mathcal{N}_0(u, v) du dv \iint_{v',u'}^{H,W} \mathcal{N}_1(u, v) du dv,$$

where  $H, W$  is the height and width of the image. This is illustrated in Fig. 1, with Gaussians over two corners illustrated on the left, and the resulting distribution  $P(u', v')$  in the centre.

The evaluation of each detection focuses on the probability value assigned to the true class label, and the spatial probability mass from  $P(u', v')$  assigned to the ground truth object *vs.* the probability mass assigned to the background. Since existing measures for conventional object detection such as mAP [26] or moLRP [34] are not equipped to evaluate the probabilistic aspects of a detection, we introduce a novel evaluation measure for Probabilistic Object Detection in the following section.

## 5. Probability-based Detection Quality (PDQ)

This section introduces the major technical contribution of our paper: the probability-based detection quality (PDQ) measure which evaluates the quality of detections based on spatial and label probabilities. Unlike AP-based measures, our approach penalises low spatial uncertainty when detecting background as foreground, or when detecting foreground as background, and explicitly evaluates the label probability in calculating detection quality. PDQ has no thresholds or tuneable parameters that can redefine the conditions of success. Furthermore, PDQ is based on an approach that provides optimal assignment of detections to

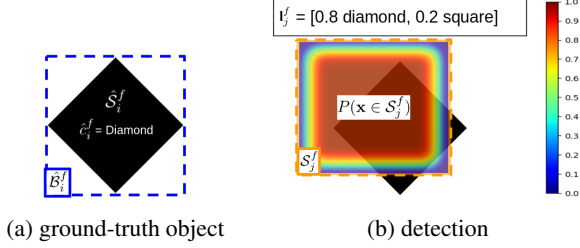


Figure 2: In our notation, a ground-truth object (a) consists of a segmentation mask  $\hat{S}_i^f$  (black), a bounding box  $\hat{B}_i^f$  (blue box), and a class label  $\hat{c}_i^f$  which here is *diamond*. A detection (b) consists of a probability density function  $P(\mathbf{x} \in \mathcal{S}_j^f)$  (illustrated as a heatmap), a segmentation mask  $\mathcal{S}_j^f$  (all pixels within the orange box), and a probability distribution across all classes  $\mathbf{I}_j^f$ , which here provides probabilities for diamond and square classes.

ground-truth objects, incorporating both the label and spatial attributes of the detections in this assignment.

A reference implementation of PDQ will be made available on github (link withheld for double-blind review).

**Notation** We write the  $i$ -th ground-truth object in the  $f$ -th frame (image) as the set  $\mathcal{G}_i^f = \{\hat{S}_i^f, \hat{B}_i^f, \hat{c}_i^f\}$ , comprising a segmentation mask defined by a set of pixels  $\hat{S}_i^f$ , a set of bounding box corners  $\hat{B}_i^f$  fully encapsulating all pixels in  $\hat{S}_i^f$ , and a class label  $\hat{c}_i^f$ .

We define the  $j$ -th detection in the  $f$ -th frame as the set  $\mathcal{D}_j^f = \{P(\mathbf{x} \in \mathcal{S}_j^f), \mathcal{S}_j^f, \mathbf{I}_j^f\}$ , comprising a probability function that returns the spatial probability that a given pixel is a part of the detection (regardless of class prediction)  $P(\mathbf{x} \in \mathcal{S}_j^f)$ , a set of pixels with a non-zero  $P(\mathbf{x} \in \mathcal{S}_j^f)$  which we refer to as the detection segmentation mask  $\mathcal{S}_j^f$ , and a label probability distribution across all possible class labels  $\mathbf{I}_j^f$ . A visualisation of both ground-truth objects and detections is provided in Figure 2.

**Requirements** PDQ requires pixel-accurate ground-truth annotations for the segmentation mask  $\hat{S}_i^f$ . Such annotations can be easily obtained from simulated environments [6, 41] and also from datasets containing only bounding box annotations by considering all pixels within a box part of the segmentation mask. PDQ can evaluate probabilistic detectors that provide bounding boxes with Gaussian corners as defined in Section 4, or conventional detectors by assuming  $P(\mathbf{x} \in \mathcal{S}_j^f) = 1 - \epsilon$  for all pixels inside the respective bounding box and  $\epsilon$  outside, for a small  $\epsilon > 0$ .

**Overview** PDQ evaluates both the *spatial* and *label* quality of a detector. It is therefore based on a combination of a spatial quality measure  $Q_S$  and a label quality measure  $Q_L$ . Both are calculated between all possible pairs of detections and ground-truth objects within a single frame. We

**Data:** a dataset of  $f = 1 \dots N_F$  frames with detections  $\mathcal{D}_j^f$  and ground-truths  $\mathcal{G}_i^f$

**forall** frames in the dataset **do**

**forall** pairs  $(\mathcal{G}_i^f, \mathcal{D}_j^f)$  **do**

calculate spatial quality  $Q_S(\mathcal{G}_i^f, \mathcal{D}_j^f)$

calculate label quality  $Q_L(\mathcal{G}_i^f, \mathcal{D}_j^f)$

calculate pPDQ( $\mathcal{G}_i^f, \mathcal{D}_j^f$ ) =  $\sqrt{Q_S \cdot Q_L}$

**end**

Based on the pPDQ(.) computed between all pairs, find optimal assignment between detections and ground-truth objects, yielding optimal pPDQ for frame  $f$ .

**end**

Combine frame-wise optimal pPDQs into an overall PDQ measure.

**Algorithm 1:** PDQ Evaluation Process

define the geometric mean between these two quality measures as the pairwise PDQ (pPDQ), and use it to find the optimal assignment between all detections and ground-truth objects within an image. The optimal pPDQ measures are then combined into an overall PDQ measure for the whole dataset. However, many of these intermediate results can also be recorded and analysed for a more detailed breakdown of performance. Algorithm 1 summarises the overall PDQ calculation. In the following, we detail each of the involved steps and both quality measures.

## 5.1. Spatial Quality

The spatial quality  $Q_S$  measures how well a detection  $\mathcal{D}_j^f$  captures the spatial extent of a ground-truth object  $\mathcal{G}_i^f$ , and takes into account the spatial probabilities for individual pixels as expressed by the detector.

Spatial quality  $Q_S$  comprises two loss terms, the foreground loss  $L_{FG}$  and the background loss  $L_{BG}$ . Spatial quality is defined as the exponentiated negative sum of the two loss terms, as follows:

$$Q_S(\mathcal{G}_i^f, \mathcal{D}_j^f) = \exp(-(L_{FG}(\mathcal{G}_i^f, \mathcal{D}_j^f) + L_{BG}(\mathcal{G}_i^f, \mathcal{D}_j^f))), \quad (1)$$

where  $Q_S(\mathcal{G}_i^f, \mathcal{D}_j^f) \in [0, 1]$ . The spatial quality in (1) is equal to 1 if the detector assigns a spatial probability of 1 to all ground-truth pixels, while not assigning any probability mass to pixels outside the ground-truth segment. This behaviour is governed by the two loss terms explained below.

**Foreground Loss** The foreground loss  $L_{FG}$  is defined as the average negative log-probability the detector assigns to the pixels of a ground-truth segment.

$$L_{FG}(\mathcal{G}_i^f, \mathcal{D}_j^f) = -\frac{1}{|\hat{S}_i^f|} \sum_{\mathbf{x} \in \hat{S}_i^f} \log(P(\mathbf{x} \in \mathcal{S}_j^f)), \quad (2)$$



where, as defined above,  $\hat{S}_i^f$  is the set of all pixels belonging to the  $i$ -th ground-truth segment in frame  $f$ , and  $P(\cdot)$  is the spatial probability function that assigns a probability value to every pixel of the  $j$ -th detection. The foreground loss is minimised if the detector assigns a probability value of one to every pixel of the ground-truth segment, in which case  $L_{FG} = 0$ . It grows without bounds otherwise.

Notice that  $L_{FG}$  intentionally ignores pixels that are inside the ground-truth bounding box  $\hat{B}_i^f$  but are *not* part of the ground-truth segment  $\hat{S}_i^f$ . This avoids treating the detection of background pixels as critically important in the case of irregularly shaped objects when pixel-level annotations are available, unlike AP-based methods using bounding-box IoUs, as illustrated in Figure 3.

**Background Loss** The background loss term  $L_{BG}$  penalises any probability mass that the detector incorrectly assigned to pixels outside the ground-truth bounding box. It is formally defined as

$$L_{BG}(\mathcal{G}_i^f, \mathcal{D}_j^f) = -\frac{1}{|\hat{S}_i^f|} \sum_{\mathbf{x} \in \mathcal{V}_{i,j}^f} \log((1 - P(\mathbf{x} \in \mathcal{S}_j^f))), \quad (3)$$

which is the sum of negative log-probabilities of all pixels in the set  $\mathcal{V}_{i,j}^f = \{\mathcal{S}_j^f - \hat{B}_i^f\}$ , i.e. pixels that are part of the detection, but not of the true bounding box. A visualisation of this evaluation region  $\mathcal{V}_{i,j}^f$  is shown in Figure 4.

Note that we average over  $|\hat{S}_i^f|$  rather than  $|\mathcal{V}_{i,j}^f|$  to ensure that foreground and background losses are scaled equivalently, measuring the loss incurred per ground-truth pixel the detection aims to describe. The background loss term is minimised if all pixels outside the ground-truth bounding box are assigned a spatial probability of zero.

## 5.2. Label Quality

While spatial quality measures how well the detection describes *where* the object is within the image, label quality  $Q_L$  measures how effectively a detection identifies *what* the object is. We define  $Q_L$  as the probability estimated by the detector for the object’s ground-truth class. Note that this is irrespective of whether this class is the highest ranked in the detector’s probability distribution. Unlike with mAP, this value is explicitly used to influence detection quality rather than just for ranking detections regardless of actual label probability. We define label quality as:

$$Q_L(\mathcal{G}_i^f, \mathcal{D}_j^f) = \mathbf{I}_j^f(\hat{c}_i^f). \quad (4)$$

## 5.3. Pairwise PDQ (pPDQ)

The pairwise PDQ (pPDQ) between a detection  $\mathcal{D}_j^f$  and a ground-truth object  $\mathcal{G}_i^f$  in frame  $f$  is the geometric mean of the spatial quality and label quality measures  $Q_S$  and  $Q_L$ :

$$\text{pPDQ}(\mathcal{G}_i^f, \mathcal{D}_j^f) = \sqrt{Q_S(\mathcal{G}_i^f, \mathcal{D}_j^f) \cdot Q_L(\mathcal{G}_i^f, \mathcal{D}_j^f)}. \quad (5)$$

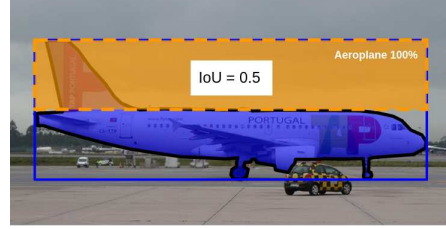


Figure 3: Example of a detection of an aeroplane (orange box), a ground-truth box (blue line), and a ground-truth segmentation mask, (blue-coloured region with black border). At an IoU threshold of 0.5, AP-based methods consider the orange detection entirely correct, despite covering only 16% of the plane’s pixels. There is no correlation between the bounding box overlap analysed and the content within the bounding box. By comparison, PDQ penalises this detection heavily for only detecting this small portion without any spatial uncertainty. The pPDQ for this detection containing no spatial uncertainty is  $3.64 \times 10^{-6}$ .

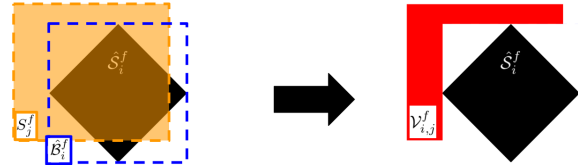


Figure 4: PDQ defines the background evaluation region  $\mathcal{V}_{i,j}^f$  (red) as the set of pixels that are part of the detection  $\mathcal{S}_j^f$  (orange), but not of the true bounding box  $\hat{B}_i^f$  (blue).

Using the geometric mean requires both components to have high values for a high pPDQ score, and is zero if either component reaches zero. Notice that it is also possible to use a weighted geometric mean for applications where the spatial or label quality component is more important.

## 5.4. Assignment of Optimal Detection-Object Pairs

It is important that, for every frame, each detection is matched to, at most, one ground-truth object and vice versa. This is also done for mAP, but it utilises a greedy assignment process based upon label confidence ranking, rather than ensuring that the optimal assignment takes into account both the spatial and label aspects of the detection. To mitigate this problem, we use our proposed pPDQ score in (5) between possible detection-object pairings to determine the optimal assignment through the Hungarian algorithm [22]. This provides the optimal assignment between two sets of information which produce the best total pPDQ score.

Using assignments from the Hungarian algorithm, we store the pPDQs for all non-zero assignments in the  $f$ -th frame in a vector  $\mathbf{q}^f = [q_1^f, q_2^f, q_3^f, \dots, q_{N_{TP}^f}^f]$  where  $N_{TP}^f$  is the number of non-zero (true positive) assignments within

the  $f$ -th frame. Note that these “true positive” detections are not ones which are considered 100% accurate as is done for AP-based measures. Instead these are detections which, even marginally, describe the ground-truth object they are matched with and provide a non-zero pPDQ. If the pPDQ from an optimal assignment is zero, there is no association between the ground-truth object and detection. This occurs when either a ground-truth object is undetected (false negative) or a detection does not describe an object (false positive). We also record the number of false negatives and false positives for each frame, expressed formally as  $N_{FN}^f$  and  $N_{FP}^f$  respectively, to be used in our final evaluation. After obtaining  $\mathbf{q}^f$ ,  $N_{TP}^f$ ,  $N_{FN}^f$ , and  $N_{FP}^f$  for each frame, the PDQ score can be calculated.

### 5.5. PDQ Score

The final PDQ score across a set of ground-truth objects  $\mathcal{G}$  and detections  $\mathcal{D}$  is the total pPDQ for each frame divided by the total number of TPs, FNs and FPs assignments across all frames. This can be seen as the average pPDQ across all TPs, FNs and FPs observed, which is calculated as follows:

$$PDQ(\mathcal{G}, \mathcal{D}) = \frac{1}{\sum_{f=1}^{N_F} N_{TP}^f + N_{FN}^f + N_{FP}^f} \sum_{f=1}^{N_F} \sum_{i=1}^{N_{TP}^f} \mathbf{q}^f(i), \quad (6)$$

where  $\mathbf{q}^f(i)$  is the pPDQ score for the  $i$ -th assigned detection-object pair in the  $f$ -th frame. This final PDQ score provides a consistent, probability-based measure, evaluating both label and spatial probabilities, that can determine how well a set of detections has described a set of ground-truth objects without the need for thresholds to determine complete success or failure of any given detection.

### 6. Evaluation of PDQ Traits

The previous section introduced PDQ, a new measure to evaluate the performance of detectors for *probabilistic* object detection. PDQ has been designed with one main goal in mind: it should reward detectors that can accurately quantify both their spatial and label uncertainty. In this section, we are going to demonstrate that this goal has been met, by showing PDQ’s behaviour in controlled experiments. We show the most critical experiments here and more are provided in supplementary material.

**PDQ Rewards Accurate Spatial Uncertainty** We perform a controlled experiment on the COCO 2017 validation dataset with a simulated object detector. For every ground truth object with true bounding box corners  $\hat{\mathbf{x}}_0$  and  $\hat{\mathbf{x}}_1$ , the detector generates a detection with bounding box corners sampled as  $\mathbf{x}_0 \sim \mathcal{N}(\hat{\mathbf{x}}_0, \hat{\Sigma})$  and  $\mathbf{x}_1 \sim \mathcal{N}(\hat{\mathbf{x}}_1, \hat{\Sigma})$ , with  $\hat{\Sigma} = \text{diag}(\hat{s}^2, \hat{s}^2)$ . We vary the value of  $\hat{s}^2$  throughout the experiments and refer to  $\hat{s}^2$  as the detector’s *true* variance.

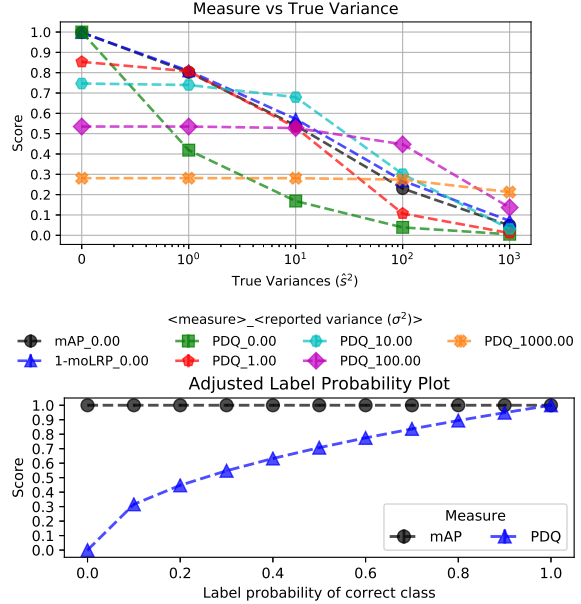


Figure 5: Top: PDQ rewards detectors that accurately evaluate their true spatial uncertainty. Bottom: PDQ explicitly evaluates label uncertainty, in contrast to mAP. See Section 6 for explanation of the experiments.

Independent of the value of  $\hat{s}^2$ , the simulated detector expresses spatial uncertainty for each probabilistic detection with a different variance  $\sigma^2$ , which we refer to as the *reported* variance. Each detection is assigned probability 1.0 for the *true* label, corresponding to perfect classification.

When varying the values of  $\hat{s}^2$  and  $\sigma^2$  and evaluating the resulting detections, PDQ should reward when  $\hat{s}^2$  is similar to  $\sigma^2$ , i.e. when the *reported* spatial uncertainty is close to the *true* spatial uncertainty that was used to sample the detection corners. When both *reported* and *true* spatial uncertainty are equal, PDQ should reach its peak performance. This would indicate that PDQ does indeed reward the accurate estimation of spatial uncertainty.

Figure 5 confirms this conjecture. We repeated the experiment described above 20 times, evaluating on all objects in the 5,000 images of the COCO 2017 validation set. Every line corresponds to a detector with a different *reported* variance  $\sigma^2$ . The *true* variance  $\hat{s}^2$  varies along the x axis. We can see that for each value of  $\hat{s}^2$ , simulated detectors with  $\sigma^2 = \hat{s}^2$  give the best performance.

**PDQ Explicitly Evaluates Label Uncertainty** We perform a controlled experiment in a simulated scenario where a single object is detected by a single detection with perfect spatial quality. We vary the detection’s reported label probability for the correct class and ensure that it always remains the dominant class in the probability distribution. The resulting PDQ and mAP scores are compared in Figure 5. We observe that PDQ is affected by the label probability of the

correct class via its label quality term. This is in contrast to mAP which uses label probability only to determine the dominant class and for ranking detection matches.

## 7. Evaluation of Object Detectors

In this section we evaluate a number of state-of-the-art conventional detectors and the recently proposed probabilistic object detector MC-Dropout SSD [30] that is based on Monte Carlo Dropout. We compare the ranking of all tested detectors using PDQ and its components, as well as the established measures mAP and moLRP [34], and discuss our most important observations and gained insights.

### 7.1. Experimental Set-up

**Evaluated Detectors** The state-of-the-art conventional object detectors evaluated were SSD-300 [27], YOLOv3 [37], FasterRCNN with ResNet backbone (FRCNN R) [51], FasterRCNN with ResNet backbone and feature pyramid network (FPN) (FRCNN R+FPN) [28], and FasterRCNN with ResNeXt backbone and FPN (FRCNN X+FPN) [28]. To evaluate these conventional detectors with PDQ, we set  $P(\mathbf{x} \in \mathcal{S}_j^f) = 1 - \epsilon$  for all pixels  $\mathbf{x}$  inside the provided standard bounding box, and  $\epsilon$  for all pixels outside, when performing the calculations in equations (2) and (3), with  $\epsilon = 10^{-14}$  to avoid infinite loss.

In addition to conventional object detectors, we evaluate a probabilistic MC-Dropout object detector based on the work by Miller et al. [29, 30]. We follow the established implementation [29], where Monte Carlo Dropout [10] is utilised in a SSD-300 object detector [27] with two dropout layers inserted and activated during both training and testing. Each image is tested with 20 forward passes through the network with randomised dropout masks to obtain samples of detections. The recommended merging strategy was used to cluster these samples [29], namely a BSAS clustering method [46] with spatial affinity IoU and label affinity ‘same label’ (we found an IoU threshold of 0.7 performed better than the 0.95 threshold recommended in [29]). Final probabilistic detections were obtained by averaging sample label probability distributions and estimating  $\mathcal{N}_0$  and  $\mathcal{N}_1$  from the average and covariance of sample bounding boxes.

We furthermore modify a FasterRCNN with ResNeXt backbone and feature pyramid network to approximate probabilistic detections. We achieve this by the following process: for every detection surviving the normal non-maximum suppression, we find all of the suppressed detections that have an IoU of above 0.75 with the surviving detections and cluster them (including the survivors). We then calculate the Gaussian corner mean and covariances of each cluster, weighted by the detection’s winning label confidences. We denote this method as probFRCNN in Table 1.

**Evaluation Protocol and Datasets** Evaluation was performed on the 5,000 images of the MS COCO 2017 validation set [26], after all detectors have been trained or fine-tuned on the 2017 training dataset. During the evaluation, we ignored all detections with the winning class label probability below a threshold  $\tau$ . We compare the effect of this process for  $\tau = 0.5$  and 0.05.

### 7.2. Insights

Table 1 presents the results of our evaluation, comparing PDQ and its components with mAP and moLRP. From these results we observe the following:

#### 1. PDQ exposes the performance differences between probabilistic and non-probabilistic object detectors.

When evaluating using mAP or moLRP, both SSD-300 [27] and FasterRCNN with ResNeXt and FPN [28], and their respective probabilistic variants (MC-DropoutSSD [29] and probFRCNN) show very similar performance. However, evaluating with PDQ reveals their performance differences in terms of probabilistic object detection: both probabilistic variants perform significantly better than their non-probabilistic counterparts. This is especially true for their overall spatial quality and its foreground and background quality components. Comparing probFRCNN with MC-DropoutSSD, we found that probFRCNN achieved a higher PDQ score, benefiting from its more accurate base network.

#### 2. PDQ reveals differences in spatial and label quality.

Since PDQ comprises meaningful components, it allows a detailed analysis of how well detectors perform in terms of spatial and label quality. For example, in Table 1 we observe that the YOLOv3 detector achieves the highest label quality (95.8%/92.8% for  $\tau = 0.5/0.05$ ), but the worst spatial quality (6.2%/5.1%) of all tested detectors. This gives important insights into worthwhile directions of future research, suggesting YOLO can be more trusted to understand *what* an object is than other detectors but is less reliable in determining precisely *where* that object is.

#### 3. Probabilistic localisation performance of existing object detectors is weak.

Spatial quality in PDQ measures how well detectors probabilistically localise objects in an image. Conventional object detectors assume full confidence in their bounding box location and achieve low spatial qualities between 5.1% and 17.6%, indicating they are spatially overconfident. Since conventional object detectors have comparatively high label qualities, we conclude that for probabilistic object detection tasks where spatial uncertainty estimation is important, improving the localisation performance and the estimation of spatial uncertainty has the biggest potential of improving performance.

#### 4. PDQ does not obscure false positive errors.

Unlike mAP and moLRP, PDQ explicitly penalises a detector for spurious (false positive) detections, as well as for missed

Table 1: PDQ-based Evaluation of Probabilistic and Non-Probabilistic Object Detectors. Legend: mLRP = 1 – moLRP, Sp = Spatial Quality, Lbl = Label Quality, FG = Foreground Quality ( $\exp(-L_{FG})$ ), BG = Background Quality ( $\exp(-L_{BG})$ ), TP = True Positives, FP = False Positives, FN = False Negatives. pPDQ, Sp, Lbl, FG and BG averaged over all TP.

Approach ( $\tau$ )	mAP (%)	mLRP (%)	PDQ (%)	pPDQ (%)	Sp (%)	Lbl (%)	FG (%)	BG (%)	TP	FP	FN
probFRCNN (0.5)	35.5	32.2	<b>28.4</b>	<b>56.7</b>	<b>45.0</b>	90.7	<b>77.8</b>	<b>60.7</b>	23,434	10,016	13,347
MC-Dropout SSD (0.5) [29]	15.8	15.6	12.8	47.3	39.9	74.0	73.1	57.3	10,510	<b>2,165</b>	26,271
MC-Dropout SSD (0.05) [29]	19.5	16.6	1.3	26.1	27.3	35.9	60.1	46.2	24,843	461,074	11,938
SSD-300 (0.5) [27]	15.0	14.3	3.9	18.1	9.7	80.2	57.5	25.1	8,999	4,746	27,782
SSD-300 (0.05) [27]	19.3	16.0	0.6	9.7	6.4	40.2	38.1	32.3	21,961	324,067	14,820
YOLOv3 (0.5) [37]	29.7	30.8	5.7	14.6	6.2	<b>95.8</b>	52.2	20.4	17,390	7,728	19,391
YOLOv3 (0.05) [37]	30.1	27.7	3.3	12.2	5.1	92.8	44.6	22.9	23,447	50,074	13,334
FRCNN R (0.5) [51]	32.8	29.1	6.7	19.1	10.3	88.8	62.2	23.6	19,930	20,044	16,851
FRCNN R (0.05) [51]	34.3	29.1	3.0	17.1	9.5	78.5	57.8	25.1	23,081	93,141	13,700
FRCNN R+FPN (0.5) [28]	34.6	31.2	11.8	27.1	16.9	86.5	60.6	35.7	22,537	14,706	14,244
FRCNN R+FPN (0.05) [28]	37.0	30.4	4.2	23.1	15.8	69.5	54.4	38.7	29,326	123,511	7,455
FRCNN X+FPN (0.5) [28]	37.4	<b>32.7</b>	11.9	27.9	17.6	88.2	60.8	36.8	24,523	20,444	12,258
FRCNN X+FPN (0.05) [28]	<b>39.0</b>	32.1	4.4	24.8	16.7	74.4	55.6	39.1	<b>29,922</b>	130,009	<b>6,859</b>



(a) FRCNN X+FPN (0.05) (b) FRCNN X+FPN (0.5)

Figure 6: Visualisation of all TPs (blue segmentation mask and corresponding BBox), FPs (orange BBox), and FNs (orange segmentation mask) as defined by PDQ for FRCNN X+FPN with  $\tau = 0.5$  (a) and  $\tau = 0.05$  (b). We see here that a lower  $\tau$  leads to far more FPs that are strongly penalised by PDQ but are largely ignored under mAP.

(false negative) detections. We observe that decreasing the label threshold  $\tau$  and consequently massively increasing the number of false positive detections (see Fig. 6 for an example) actually increases mAP, and does not tend to affect moLRP much. In contrast, PDQ scores decrease significantly. PDQ is designed to evaluate systems for application in real-world systems and does not filter detections based on label ranking or calculating the optimal threshold  $\tau$ . It involves *all* reported detections in its analysis.

## 8. Conclusions and Future Work

We introduced Probabilistic Object Detection, a challenging new task that is highly relevant for domains where accurately estimating the spatial and semantic uncertainties of the detections is of high importance such as embodied AI (such as robotics, autonomous systems, driverless cars), and medical imaging. To foster further research in this direction, we introduced the probability-based detection quality (PDQ) measure which explicitly evaluates both spatial and label uncertainty.

PDQ is not meant to *replace* mAP, but to *complement* it. Both evaluation measures are designed for two *different* problems. While mAP has been the established performance measure for conventional object detection, we developed PDQ specifically for the new task of *probabilistic* object detection.

After evaluating a range of object detectors, including the first emerging probabilistic object detector in Section 7, we are confident that PDQ is a useful performance measure that can guide and inform the research of even better probabilistic object detectors in the future. In future work we will investigate how to train object detectors to directly optimise for PDQ by incorporating it into the training loss function. The concept of probabilistic object detection can also be easily extended to Probabilistic *Instance Segmentation* where each pixel would contain a probability of belonging to a certain object instance, along with a label distribution.

## References

- [1] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané. Concrete problems in ai safety. *arXiv*



- preprint *arXiv:1606.06565*, 2016. 1, 2
- [2] D. Cireşan, U. Meier, J. Masci, and J. Schmidhuber. Multi-column deep neural network for traffic sign classification. *Neural networks*, 32:333–338, 2012. 2
- [3] G. Csurka, D. Larlus, F. Perronnin, and F. Meylan. What is a good evaluation measure for semantic segmentation?. In *BMVC*, volume 27, page 2013. Citeseer, 2013. 3
- [4] J. Dai, Y. Li, K. He, and J. Sun. R-fcn: Object detection via region-based fully convolutional networks. In *NIPS*, pages 379–387, 2016. 2, 3
- [5] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *international Conference on computer vision & Pattern Recognition (CVPR’05)*, volume 1, pages 886–893. IEEE Computer Society, 2005. 2
- [6] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun. CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 1–16, 2017. 4
- [7] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015. 1, 2
- [8] M. Everingham, L. V. Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The Pascal Visual Object Classes (VOC) Challenge. *International Journal of Computer Vision*, 88(2):303–338, June 2010. 1, 2, 3
- [9] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2010. 2
- [10] Y. Gal and Z. Ghahramani. Bayesian convolutional neural networks with bernoulli approximate variational inference. *arXiv preprint arXiv:1506.02158*, 2015. 2, 7
- [11] Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016. 2
- [12] R. Girshick. Fast r-cnn. In *ICCV*, pages 1440–1448, 2015. 2
- [13] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Region-based convolutional networks for accurate object detection and segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 2015. 2
- [14] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. In *ICML*, 2017. 1, 2
- [15] D. Hendrycks and K. Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *Proceedings of International Conference on Learning Representations*, 2017. 1, 2
- [16] D. Hoiem, Y. Chodpathumwan, and Q. Dai. Diagnosing error in object detectors. In *European conference on computer vision*, pages 340–353. Springer, 2012. 3
- [17] M. Kampffmeyer, A.-B. Salberg, and R. Jenssen. Semantic segmentation of small objects and modeling of uncertainty in urban remote sensing images using deep convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 1–9, 2016. 2
- [18] A. Kendall, V. Badrinarayanan, and R. Cipolla. Bayesian segnet: Model uncertainty in deep convolutional encoder-decoder architectures for scene understanding. *arXiv preprint arXiv:1511.02680*, 2015. 2
- [19] A. Kendall and Y. Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pages 5574–5584, 2017. 1, 2
- [20] I. Krasin, T. Duerig, N. Alldrin, V. Ferrari, S. Abu-El-Haija, A. Kuznetsova, H. Rom, J. Uijlings, S. Popov, S. Kamali, M. Mallocci, J. Pont-Tuset, A. Veit, S. Belongie, V. Gomes, A. Gupta, C. Sun, G. Chechik, D. Cai, Z. Feng, D. Narayanan, and K. Murphy. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from <https://storage.googleapis.com/openimages/web/index.html>*, 2017. 1, 2
- [21] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *NIPS*, volume 1, page 4, 2012. 00312. 2
- [22] H. W. Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2:83–97, 1955. 5
- [23] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Mallocci, T. Duerig, and V. Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *arXiv:1811.00982*, 2018. 1, 2
- [24] I. Lenz, H. Lee, and A. Saxena. Deep learning for detecting robotic grasps. *The International Journal of Robotics Research*, 34(4-5):705–724, 2015. 2
- [25] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. *IEEE transactions on pattern analysis and machine intelligence*, 2018. 2, 3
- [26] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollr, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1, 2, 3, 7
- [27] 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*, pages 21–37. Springer, 2016. 1, 2, 3, 7, 8
- [28] F. Massa and R. Girshick. maskrcnn-benchmark: Fast, modular reference implementation of Instance Segmentation and Object Detection algorithms in PyTorch. <https://github.com/facebookresearch/maskrcnn-benchmark>, 2018. 7, 8
- [29] D. Miller, F. Dayoub, M. Milford, and N. Sünderhauf. Evaluating Merging Strategies for Sampling-based Uncertainty Techniques in Object Detection. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2019. 2, 3, 7, 8
- [30] D. Miller, L. Nicholson, F. Dayoub, M. Milford, and N. Sünderhauf. Dropout Sampling for Robust Object Detection in Open-Set Conditions. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2018. 2, 7
- [31] T. Nair, D. Precup, D. L. Arnold, and T. Arbel. Exploring uncertainty measures in deep networks for multiple sclerosis lesion detection and segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 655–663. Springer, 2018. 2
- [32] T. Namba and Y. Yamada. Risks of deep reinforcement learn-

- ing applied to fall prevention assist by autonomous mobile robots in the hospital. *Big Data and Cognitive Computing*, 2(2):13, 2018. [2](#)
- [33] A. Nguyen, J. Yosinski, and J. 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](#), [2](#)
- [34] K. Oksuz, B. Can Cam, E. Akbas, and S. Kalkan. Localization recall precision (lrp): A new performance metric for object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 504–519, 2018. [2](#), [3](#), [7](#)
- [35] C. Otte. Safe and interpretable machine learning: A methodological review. In C. Moewes and A. Nürnberger, editors, *Computational Intelligence in Intelligent Data Analysis*, pages 111–122, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg. [2](#)
- [36] J. Pont-Tuset and F. Marques. Supervised evaluation of image segmentation and object proposal techniques. *IEEE transactions on pattern analysis and machine intelligence*, 38(7):1465–1478, 2016. [3](#)
- [37] J. Redmon and A. Farhadi. YOLOv3: An Incremental Improvement. *arXiv:1804.02767 [cs]*, Apr. 2018. [1](#), [2](#), [7](#), [8](#)
- [38] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, pages 91–99, 2015. [1](#), [2](#), [3](#)
- [39] C. Richter, W. Vega-Brown, and N. Roy. Bayesian learning for safe high-speed navigation in unknown environments. In *Robotics Research*, pages 325–341. Springer, 2018. [2](#)
- [40] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, Dec. 2015. [1](#), [2](#), [3](#)
- [41] M. Savva, A. X. Chang, A. Dosovitskiy, T. Funkhouser, and V. Koltun. MINOS: Multimodal indoor simulator for navigation in complex environments. *arXiv:1712.03931*, 2017. [4](#)
- [42] K. Sirinukunwattana, S. e. A. Raza, Y. Tsang, D. Snead, I. Cree, and N. Rajpoot. Locality sensitive deep learning for detection and classification of nuclei in routine colon cancer histology images. *IEEE Transactions on Medical Imaging*, 35:1–1, 02 2016. [2](#), [3](#)
- [43] N. Sünderhauf, O. Brock, W. Scheirer, R. Hadsell, D. Fox, J. Leitner, B. Upcroft, P. Abbeel, W. Burgard, M. Milford, et al. The limits and potentials of deep learning for robotics. *The International Journal of Robotics Research*, 37(4-5):405–420, 2018. [1](#), [2](#)
- [44] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *Proceedings of International Conference on Learning Representations*, 2014. [1](#), [2](#)
- [45] R. Tanno, D. E. Worrall, A. Ghosh, E. Kaden, S. N. Sotiropoulos, A. Criminisi, and D. C. Alexander. Bayesian image quality transfer with cnns: exploring uncertainty in dmri super-resolution. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 611–619. Springer, 2017. [2](#)
- [46] S. Theodoridis and K. Koutroumbas. *Pattern Recognition, Second Edition*. Academic Press, Inc., Orlando, FL, USA, 2nd edition, 2003. [7](#)
- [47] S. Thrun, W. Burgard, and D. Fox. *Probabilistic robotics*. MIT press, 2005. [2](#)
- [48] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In *CVPR 2011*, pages 1521–1528, June 2011. [1](#), [2](#)
- [49] K. R. Varshney and H. Alemzadeh. On the safety of machine learning: Cyber-physical systems, decision sciences, and data products. *Big data*, 5 3:246–255, 2017. [1](#), [2](#)
- [50] P. Viola and M. J. Jones. Robust real-time face detection. *International journal of computer vision*, 57(2):137–154, 2004. [2](#)
- [51] J. Yang, J. Lu, D. Batra, and D. Parikh. A Faster Pytorch Implementation of Faster R-CNN. <https://github.com/jwyang/faster-rcnn.pytorch>, 2017. [7](#), [8](#)