

Object Detection as Probabilistic Set Prediction

Georg Hess^{1,2}, Christoffer Petersson^{1,2}, and Lennart Svensson¹

¹ Chalmers University of Technology, Gothenburg, Sweden, georghe@chalmers.se
² Zenseact, Gothenburg, Sweden

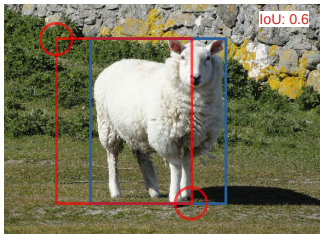
Abstract. Accurate uncertainty estimates are essential for deploying deep object detectors in safety-critical systems. The development and evaluation of probabilistic object detectors have been hindered by shortcomings in existing performance measures, which tend to involve arbitrary thresholds or limit the detector’s choice of distributions. In this work, we propose to view object detection as a set prediction task where detectors predict the distribution over the set of objects. Using the negative log-likelihood for random finite sets, we present a proper scoring rule for evaluating and training probabilistic object detectors. The proposed method can be applied to existing probabilistic detectors, is free from thresholds, and enables fair comparison between architectures. Three different types of detectors are evaluated on the COCO dataset. Our results indicate that the training of existing detectors is optimized toward non-probabilistic metrics. We hope to encourage the development of new object detectors that can accurately estimate their own uncertainty. Code available at <https://github.com/georghess/pmb-nll>.

Keywords: Probabilistic object detection, random finite sets, proper scoring rules, uncertainty estimation.

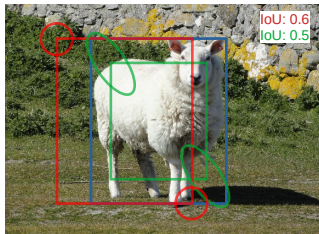
1 Introduction

Accurately locating and classifying a set of objects has a range of applications, such as autonomous driving, transportation, surveillance, scene analysis, and image captioning. Common approaches for solving this rely on a deep object detector which provides a set of detections containing bounding box parameters, semantic class and classification confidence. However, as pointed out in previous works [3,5,11,15,28] most state-of-the-art networks lack the ability to assess their own regression confidence and fail to provide a complete uncertainty description. As an effect, this can limit the performance in downstream tasks such as multi-object tracking, sensor fusion, or decision making, ultimately hindering humans to establish trust in the deep learning agent.

There are many strategies to evaluate predictive uncertainties in the deep learning regime. Broadly speaking, a distribution should perform well on two criteria: calibration and sharpness. For a distribution to be well calibrated, it should not be over- or under-confident, but reflect the true confidence in its predictions. Sharpness instead promotes concentrated and, consequently, informative distributions [9]. Both these properties can be measured simultaneously



(a) Assigning predictions to ground truth is non-trivial in probabilistic OD. The prediction can either be considered a true positive with bad uncertainty estimates or a false positive.



(b) mAP and [13] prefer the red prediction with larger IoU. Our method considers uncertainties and multiple assignments, and finds the green prediction a more probable match to the blue ground truth.

Fig. 1: Predictions (red and green) and ground truth (blue), highlighting the object detection assignment problem. Ellipses represent spatial uncertainty.

by using a proper scoring rule such as negative log-likelihood [10]. Proper scoring rules assess the quality of predictive uncertainties and are minimized only when the prediction is equivalent to the distribution that generated the ground truth observations [10]. Besides measuring calibration and sharpness, proper scoring rules enable a theoretically sound ranking of different predictive distributions.

Evaluating the quality of predictive uncertainties in object detection (OD) is a non-trivial task. First, any measure has to jointly consider the performance in terms of ability to detect, correctly classify and accurately locate objects. Second, as we do not know the correspondence between predictions and ground truth objects, any analysis is colored by the selected assignment rules. As an example, the prediction in Fig. 1a can be considered either a correct detection with bad uncertainty predictions or a false positive. Having multiple predictions makes the assignment even harder, as shown in Fig. 1b. The most common measure in OD, mAP, uses handcrafted assignment rules based on IoU and class confidence and fails to consider predicted uncertainties. The probability-based detection quality (PDQ) [11] tries to address these issues, but is limited to Gaussian distributions for regression. More recently, the lack of proper scoring rules for evaluating probabilistic object detection was pointed out by [13], also proving that PDQ is not a proper scoring rule. However, while they use proper scoring rules for the different subtasks, such as the energy score for regression and the Brier score for classification, predictions are assigned to targets using ad hoc IoU-based rules which ignore regression uncertainties. As highlighted earlier, these types of assignment rules have a large influence on the reported performance, do not yield proper scoring rules, and make it harder to draw conclusions about model performance.

In this paper, we propose to use random finite sets (RFS) to model the probabilistic object detection task. Object detection is often seen as a set prediction task, and we extend this perspective to probabilistic object detection (PrOD).

We describe the set of objects in a given image by a single random variable, and the task of our object detection networks is to describe the distribution of that variable. This simple change of perspective enables us to use the negative log-likelihood to evaluate the uncertainty estimates of our detections, which gives rise to the first proper scoring rule for object detection. Our framework explicitly models the assignment problem, is general enough to be applied to any type of distribution, enables easy ranking between different algorithms, and can be decomposed to highlight different types of errors (detection, regression and classification). Our key contributions are the following.

- We propose to view the set of objects in an image as a single stochastic variable. By applying the negative log-likelihood (NLL) to a distribution over sets, we present the first proper scoring rule for object detection.
- We show how to apply the random finite set framework to object detection by interpreting the detector output as parameters of multi-Bernoulli (MB) and Poisson multi-Bernoulli (PMB) densities. Further, we present how to efficiently calculate and interpret the NLL of the MB and PMB densities.
- Using our proposed scoring rule, we evaluate one-stage, two-stage, and set-based detectors on the popular MS COCO dataset, and showcase their strengths and shortcomings using the decomposability of PMB-NLL.
- Further, we leverage the fact that the proposed method is differentiable and fine-tune the detectors to optimize PMB-NLL directly. Our results show that this helps detectors to reduce the number of false and duplicate detections.

2 Related Work

Quantifying uncertainties with deep neural networks has been a long-standing challenge. We aim to provide a brief overview here, as to motivate the importance of our work. Interested readers are referred to [1,5,8] for details.

Types of Uncertainties. In computer vision, uncertainties are generally divided into two categories: aleatoric and epistemic [15]. The first category refers to noise inherent to the data, which can originate from sensor noise, class ambiguities, label noise and such, and cannot be reduced with more data. Epistemic uncertainties are due to uncertainties in model parameters, and can, in principle, be eliminated given enough data. In this work, we do not aim to disentangle the two types, but consider overall predictive uncertainty [26].

Uncertainty Estimation. Most approaches for quantifying uncertainties in object detection either apply Monte Carlo dropout [12,21,33], deep ensembles [6,19] or direct modeling [14,16,36]. Unfortunately, uncertainty estimates are often overlooked when evaluating probabilistic detectors, while methods that do evaluate their uncertainties use a range of different performance measures, making comparison challenging. The lack in standard performance measures has also been pointed out as a main obstacle for uncertainty estimation [1,5,28].

Evaluating Uncertainty. As the commonly used performance measure mAP fails to consider spatial uncertainties and is insensitive to badly calibrated classification, several methods trying to address these issues have been suggested. The Probability-based Detection Quality (PDQ) [11] evaluates both spatial and semantic uncertainties, but is limited to Gaussian spatial uncertainties, requires practitioners to select confidence thresholds, and has been shown by [13] to not be a proper scoring rule, thereby introducing biases into its ranking of detectors. The authors of [13] promote the use of proper scoring rules for object detection. However, their approach disregards the spatial uncertainty information when assigning predictions to targets, requires confidence thresholds, and does not provide clear recommendations on model ranking.

Set Prediction. While object detection inherently can be seen as a set prediction task, this has been made more explicit by a range of set-based detectors [2,27,35,37]. These detectors highlight the assignment problem, i.e., how to assign predictions to ground truth elements when calculating losses or metrics. In this work, we extend this perspective to probabilistic object detection by modeling the problem using distributions over random finite sets. This paradigm is applicable to any type of detector, set-based or not, and naturally models and solves the assignment problem.

Random Finite Sets. Random finite sets have been used extensively in the model-based multi-object tracking community [7,20,30,32]. The RFS framework has proven useful for modeling potentially detected and undetected objects as it captures uncertainties in the cardinality of present objects and their individual properties. However, these algorithms are often evaluated without taking their uncertainties into account. Recently, the authors of [24] suggested the use of negative log-likelihood for probabilistic evaluation of model-based multi-object trackers and presented an efficient approximation of the NLL. Our work shows how to interpret parameters of existing deep object detectors as RFSs and uses [24] to calculate our proper scoring rule. Unlike the custom designed and low-dimensional regression problems explored in [24], we apply this method to a large scale dataset, jointly evaluating detection, classification, and regression.

3 Probabilistic Modeling for Object Detection

Object detection is a set prediction task, where, given an image \mathbf{X} , the aim is to predict the set of corresponding objects \mathbb{Y} present in said image. Here, the number of objects n in the set $\mathbb{Y} = \{y_1, y_2, \dots, y_n\}$ is unknown beforehand. Further, for each object $y_i = (c_i, b_i)$, we do not know which class $c_i \in \{1, \dots, C\}$ it belongs to, nor where its bounding box $b_i \in \mathbb{R}^4$ is located in the image. In supervised learning, we aim to learn a model that, given the image \mathbf{X} , predicts a set of \hat{n} objects $\hat{\mathbb{Y}} = \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_{\hat{n}}\}$ which is close to the ground truth label \mathbb{Y} in some sense. For probabilistic object detection, we further want an uncertainty description for the number of objects and their individual properties.

In this work, we evaluate probabilistic object detectors by seeing the set of objects \mathbb{Y} as a single random variable. The task for our networks is to predict the distribution of this set $f(\mathbb{Y}|\mathbf{X})$. This is a natural and general probabilistic extension to the set prediction perspective, as a distribution over sets can capture the varying cardinality and uncertainty in properties for individual objects. Using this novel perspective, all predictions for a single image are evaluated together by applying the negative log-likelihood

$$\text{NLL}((\mathbb{Y}, \mathbf{X}), f) = -\log(f(\mathbb{Y}|\mathbf{X})). \quad (1)$$

This can be compared to existing methods where classification and regression are treated separately and evaluated conditioned on an ad hoc assignment rule [13], or network performance is measured using non-proper scoring rules [11].

To use the negative log-likelihood in practice, we need our deep object detectors to predict distributions $f(\mathbb{Y}|\mathbf{X})$. We propose to use random finite sets (RFSs) and the Poisson multi-Bernoulli (PMB) distribution and demonstrate how the PMB parameters are naturally obtained from the output of standard probabilistic deep object detectors. Further, using the results of [24], we show how to efficiently calculate and decompose the negative log-likelihood of $f(\mathbb{Y}|\mathbf{X})$ into detection, classification and regression errors.

Notation: Scalars and vectors are denoted by lowercase or uppercase letters with no special typesetting x , matrices by uppercase boldface letters \mathbf{X} , and sets by uppercase blackboard-bold letters \mathbb{X} . We define $\mathbb{N}_a = \{i \in \mathbb{N} | i \leq a\}$, $a \in \mathbb{N}$.

3.1 Modeling Detections with Random Finite Sets

We need a way to describe the distribution over \mathbb{Y} using deep neural networks. Interestingly, existing probabilistic detectors already contain the parameters needed. To this end, we propose to model \mathbb{Y} with random finite sets. Random finite sets are described using a multi-object density $f(\mathbb{Y})$, which means that sampling from $f(\mathbb{Y})$ yields finite sets of objects with varying cardinality, where objects consist of a class and a bounding box. We should note that RFSs are not the only way to describe a distribution over \mathbb{Y} . However, we will show that our method has multiple properties suitable for object detection and advantages such as being compatible with existing architectures.

Bernoulli RFS. One of the simplest RFSs is the Bernoulli RFS, commonly used for modeling single potential objects in the multi-target tracking community [7,29]. Here, we use it to model each individual detected object, and its density is

$$f_{\text{B}}(\mathbb{Y}) = \begin{cases} 1 - r & \text{if } \mathbb{Y} = \emptyset, \\ rp(y) & \text{if } \mathbb{Y} = \{y\}, \\ 0 & \text{if } |\mathbb{Y}| > 1, \end{cases} \quad (2)$$

where $p(y)$ is the single-object density. For instance, assuming the class and bounding box to be independent, $p(y) = p_{\text{cls}}(c)p_{\text{reg}}(b)$ contains the class distribution $p_{\text{cls}}(c)$ and some density describing the object's spatial distribution

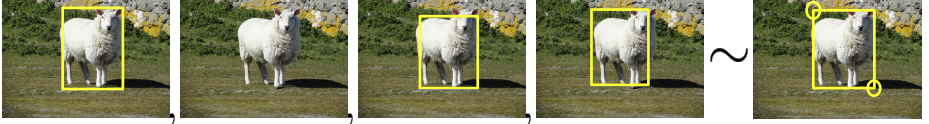


Fig. 2: Four sampled sets (left) from a Bernoulli RFS (right) with existence probability $r = 0.75$. The RFS can model the absence of objects, as well as semantic and spatial uncertainties. The image is only included for context.

$p_{\text{reg}}(b)$. Further, $r \in [0, 1]$ is the probability of existence, which is the probability that the Bernoulli RFS yields an object when sampling from it. Note that a Bernoulli RFS can account for at most one object since the likelihood is zero for any set with cardinality greater than one.

The parameters of $p(y)$ are already present in probabilistic detectors. Depending on the architecture, we can interpret r as objectness and predicted it directly, or, find it as the sum of probabilities assigned to foreground classes and let $p_{\text{cls}}(c)$ be the class distribution conditioned on existence. Note that $f_{\text{B}}(\emptyset) = 1 - r$ is the probability that the object is not present, which we may think of as the event where the prediction is background. An example of a Bernoulli RFS prediction and corresponding samples is shown in Fig. 2.

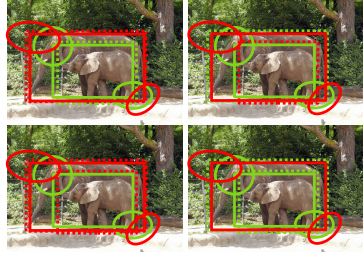
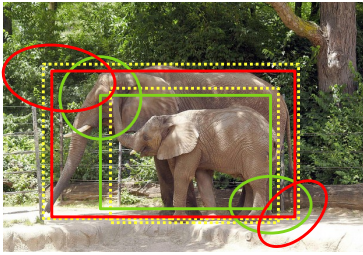
Multi-Bernoulli RFS. Generally, the number of objects in an image can vary greatly. Modeling many potential objects can be achieved by taking the union of multiple Bernoulli RFSs [7], resulting in a multi-Bernoulli (MB) RFS. In other words, individual predictions made by a detector are interpreted as parameters of individual Bernoulli RFSs, and by taking their union we combine them into a single random variable. Unlike a Bernoulli RFS, an MB RFS can be used to model the set of potentially detected objects for an entire image.

Formally, let $\mathbb{X}_1, \dots, \mathbb{X}_m$ be m independent Bernoulli RFSs with the densities $f_{\text{B}_1}(\mathbb{X}_1), \dots, f_{\text{B}_m}(\mathbb{X}_m)$, existence probabilities r_1, \dots, r_m , and single-object densities $p_1(x), \dots, p_m(x)$. Then $\mathbb{X} = \cup_{i=1}^m \mathbb{X}_i$ is an MB RFS with multi-object density

$$f_{\text{MB}}(\mathbb{X}) = \sum_{\cup_{i=1}^m \mathbb{X}_i = \mathbb{X}} \prod_{j=1}^m f_{\text{B}_j}(\mathbb{X}_j), \quad (3)$$

where $\sum_{\cup_{i=1}^m \mathbb{X}_i = \mathbb{X}}$ denotes the sum over all disjoint sets whose union is \mathbb{X} . In other words, when evaluating the multi-object density $f_{\text{MB}}(\mathbb{Y})$ of a set \mathbb{Y} we sum the multi-object densities of all possible assignments between elements in \mathbb{Y} and Bernoulli components in f_{MB} .

We illustrate this concept with an example. Consider an image containing two objects $\mathbb{Y} = \{y_1, y_2\}, y_i = (c_i, b_i)$ and two predictions, as shown in Fig. 3. Each prediction consists of a class distribution and a spatial pdf. We let these parameterize the densities $f_{\text{B}_1}(\cdot), f_{\text{B}_2}(\cdot)$ of two separate Bernoulli RFS, whereas the multi-object density $f_{\text{MB}}(\mathbb{Y})$ of their union is the MB RFS used to describe



(a) Set \mathbb{Y} (yellow) with two ground truth objects and a multi-Bernoulli with two Bernoullis \mathbb{X}_1 and \mathbb{X}_2 with densities $f_{B_1}(\cdot)$ and $f_{B_2}(\cdot)$. Ellipses represent uncertainties in bounding box location and shape. The large spatial uncertainties in $f_{B_1}(\cdot)$ make it a decent description of both true objects.

(b) Visualization of the four potential assignments, ordered in decreasing likelihood. In the bottom row, both ground truth objects have been assigned to \mathbb{X}_1 and \mathbb{X}_2 respectively. As $f_{B_1}(\mathbb{Y}) = f_{B_2}(\mathbb{Y}) = 0$ for sets with cardinality larger than one, both these assignments have a likelihood of zero.

Fig. 3: Visualization of likelihood evaluation for a multi-Bernoulli RFS.

all objects in the image. When evaluating the likelihood $f_{MB}(\mathbb{Y})$ using (3), we sum the four ways to assign ground truth objects to the Bernoulli RFSs

$$f_{MB}(\mathbb{Y}) = f_{B_1}(\{y_1\})f_{B_2}(\{y_2\}) + f_{B_1}(\{y_2\})f_{B_2}(\{y_1\}) + f_{B_1}(\{y_1, y_2\})f_{B_2}(\emptyset) + f_{B_1}(\emptyset)f_{B_2}(\{y_1, y_2\}), \quad (4)$$

where each individual assignment is visualized in Fig. 3b. As $f_{B_1}(\cdot)$ and $f_{B_2}(\cdot)$ both evaluate to zero for sets with more than one element, the last two assignments have a likelihood of zero, and we are left with two terms

$$f_{MB}(\mathbb{Y}) = r_1 p_{1,cls}(c_1) p_{1,reg}(b_1) \cdot r_2 p_{2,cls}(c_2) p_{2,reg}(b_2) + r_1 p_{1,cls}(c_2) p_{1,reg}(b_2) \cdot r_2 p_{2,cls}(c_1) p_{2,reg}(b_1). \quad (5)$$

In contrast to existing methods with handcrafted assignment rules [11,13,18], the assignment problem is modeled explicitly and rigorously by the RFS framework. The intuition behind considering all possible assignments is that we cannot know the correspondence between ground truths and predictions. In cases with overlapping boxes, predictions may have large IoU with multiple objects, making the assignment highly ambiguous. Further, for PrOD, large uncertainties can make it even harder to pair predictions to true objects.

3.2 Proper Scoring Rule for Object Detection

A scoring rule measures the quality of predictive uncertainty in terms of sharpness and calibration [10]. It does so by assigning a numerical value $S(p_\theta, (\mathbf{x}, \mathbf{y}))$ to a predicted distribution $p_\theta(\mathbf{y}|\mathbf{x})$, given that some event $(\mathbf{x}, \mathbf{y}) \sim p^*(\mathbf{y}|\mathbf{x})p(\mathbf{x})$ materialized, where a lower number indicates better quality. A scoring rule is further known to be strictly proper if it is minimized only when p_θ is equal to the

distribution p^* that generated the observed event. For OD this translates to the predictive distribution being the same as the distribution from which the annotations have been generated. The noise present in a perfect prediction p_θ should in other words be equal to the noise in the annotations. These properties make proper scoring rules suitable for evaluating and ranking different predictions.

Negative log-likelihood (NLL) is a local proper scoring rule used to evaluate the quality of predictive distributions for both regression and classification. Local refers to the fact that the predicted distribution is only evaluated at the event that materialized [4]. If we let \mathbb{Y} denote the set of ground truth objects present in the current image and let $f_{\text{MB}}(\mathbb{X})$ be the multi-object density of an MB RFS produced by some model, then

$$\text{NLL}(\mathbb{Y}, f_{\text{MB}}) = -\log f_{\text{MB}}(\mathbb{Y}) = -\log \left(\sum_{\substack{\Psi_{i=1}^m \mathbb{Y}_i = \mathbb{Y}}} \prod_{j=1}^m f_{B_j}(\mathbb{Y}_j) \right). \quad (6)$$

As discussed in the previous section, to evaluate the likelihood of an MB RFS density $f_{\text{MB}}(\mathbb{Y})$, we consider all possible assignments. As the number of predictions m , or the cardinality of the ground truth set $|\mathbb{Y}|$ grows, the number of assignments grows super-exponentially, making the NLL computation intractable. However, recently it was shown how to efficiently approximate the NLL of certain RFS densities, including the MB density [24]. Assuming that the ground truth objects, as well as the individual Bernoulli components, are somewhat separated, only a few assignments have a substantial contribution to the overall likelihood. Referring back to the example from Fig. 3, we can see that mainly the first assignment contributed to the sum of likelihoods. Thus, we approximate the NLL by only considering the most likely assignments. We find these assignments efficiently by solving an optimal assignment problem

$$\min_{\mathbf{A}} \sum_k \sum_l C_{k,l} A_{k,l} \quad (7a)$$

$$\text{s.t.} \quad \sum_{k=1}^{m+|\mathbb{Y}|} A_{k,l} = 1, \quad \sum_{l=1}^{|\mathbb{Y}|} A_{k,l} \leq 1, \quad (7b)$$

$$C_{k,l} = -\log \left(\frac{p_k(y_l)}{1 - r_k} r_k \right), \quad (7c)$$

where \mathbf{C} is a cost matrix and its derivation can be found in Appendix A. In (7c), both the cost of assigning the object l to prediction k , $p_k(y_l)r_k$, and the alternative of not assigning the prediction to anything, $1 - r_k$, are considered jointly. The assignment matrix \mathbf{A} describes the pairing between predictions and ground-truth objects, where ground truth object y_l is assigned to the k -th component of the MB i.f.f. $[\mathbf{A}]_{k,l} = 1$. Murty's algorithm [22,23] efficiently computes the Q lowest cost associations $\mathbf{A}_1^*, \dots, \mathbf{A}_Q^*$ to this assignment problem. We obtain

$$\text{NLL}(\mathbb{Y}, f_{\text{MB}}) \approx -\log \left(\sum_{q=1}^Q \prod_{k=1}^m f_{B_k}(\mathbb{Y}_k(\mathbf{A}_q^*)) \right), \quad (8)$$

where $\mathbb{Y}_k(\mathbf{A}_q^*) = \{y_j \in \mathbb{Y} | [\mathbf{A}_q^*]_{k,j} = 1\}$, i.e., \mathbb{Y}_k contains the ground truth y_j if y_j was assigned to Bernoulli component k , otherwise it is the empty set. Comparing this expression to (6), only Q terms have to be calculated. During our experiments, we use $Q = 25$ as we find that the approximation does not change considerably when using additional assignments.

3.3 Modeling All Objects

Using only a MB RFS to describe the objects in an image can be problematic as it assumes that the number of predictions is greater than or equal to the number of objects present. For an algorithm providing too few detections, multiple objects are assigned to the same Bernoulli in (3), resulting in the MB likelihood being zero and an infinite NLL. Fortunately, there are RFSs that can model an arbitrary number of objects. Within model-based multi-object tracking, the Poisson Point Process (PPP) is used to model undetected objects [7], and we show here how to use it for OD to ensure a finite NLL. The PPP is then combined with the detections, yielding the Poisson multi-Bernoulli (PMB) RFS. Importantly, we also establish a technique to obtain the PPP directly from the output of our deep object detectors.

Poisson Point Process. Intuitively speaking, the PPP is intended to capture objects that are not properly detected. By complementing the detections in the MB, we model both the detected and undetected objects in an image. In contrast to the MB, the cardinality of a PPP is Poisson distributed which gives it a non-zero probability for any set cardinality. Thus we avoid the issue of infinite NLL due to lack of detections. The multi-object density of a PPP is

$$f_{\text{PPP}}(\mathbb{X}) = \exp(-\bar{\lambda}) \prod_{x \in \mathbb{X}} \lambda(x), \quad (9)$$

where $\lambda(\cdot)$ is the intensity function and $\bar{\lambda} = \int \lambda(x') dx'$ is the expected cardinality of the set. The intensity function is expected to describe the properties of poorly detected objects, e.g., partially occluded objects, far-away objects, or even classes of objects that are inherently harder to detect. The intensity function is similar to a density function, but its integral does not have to sum to one.

Poisson multi-Bernoulli RFS. With models for both detected and undetected objects, we have to combine them to a single model for all objects. To this end, we propose to use a Poisson multi-Bernoulli (PMB) RFS, which is the union of a PPP and an MB RFS. The PMB RFS also arises naturally as the posterior density of all objects after a single measurement update, when using standard models in model-based target tracking [7,32].

The multi-object density of a PMB is

$$f_{\text{PMB}}(\mathbb{X}) = \sum_{\mathbb{X}^{\text{U}} \uplus \mathbb{X}^{\text{D}} = \mathbb{X}} f_{\text{PPP}}(\mathbb{X}^{\text{U}}) f_{\text{MB}}(\mathbb{X}^{\text{D}}), \quad (10)$$

were $\mathbb{X}^U \uplus \mathbb{X}^D$ refers to summing over all possible ways of partitioning \mathbb{X} into two disjoint sets, one being the set of undetected objects \mathbb{X}^U and the other one being the set of detected objects \mathbb{X}^D . When evaluating the likelihood of a set \mathbb{Y} this translates to, for each object in \mathbb{Y} , considering it to be detected and assigning it to a Bernoulli following (3), or it being undetected and assigning it to the PPP.

Selecting the PPP Intensity. To use the PMB in object detection, we must describe the PPP intensity function $\lambda(\cdot)$. During this work, we explored various ways of learning $\lambda(\cdot)$ from data, e.g., estimating the parameters of a uniform intensity function or describing $\lambda(\cdot)$ as a constant mixture model. However, the method we found to work best for the detectors considered in our experiments, is to create $\lambda(\cdot)$ from low confidence predictions. In practice, we parameterize the intensity function as the unnormalized mixture of low confidence predictions where the mixture weights are the existence probabilities

$$\lambda(x) = \sum_i r_i p_i(x). \quad (11)$$

Specifically, we remove all predictions from the MB RFS, whose existence probabilities are $r < 0.1$, and instead use them to construct the intensity function using (11). The theoretical motivation for this change is that the Kullback-Leibler divergence between a Bernoulli RFS with existence probability $r < 0.1$ and a PPP with intensity function $\lambda(x) = rp(x)$ is small [31]. The proposed PMB density should therefore be a good approximation to the MB density that we had before, but this minor adjustment is sufficient to avoid issues with infinite NLL.

NLL Evaluation. For evaluating the NLL of a PMB RFS, we use the same approach as for the MB and consider only the Q most likely assignments. The cost matrix from (7c) used in the optimization is extended to

$$C_{k,l} = \begin{cases} -\log\left(\frac{p_k(y_l)}{1-r_k}r_k\right), & \text{if } k \leq |\mathbb{Y}| \\ -\log\lambda(y_l), & \text{if } k = l + |\mathbb{Y}| \\ \infty, & \text{otherwise,} \end{cases} \quad (12)$$

which translates to appending a diagonal matrix of size $|\mathbb{Y}| \times |\mathbb{Y}|$ to the original cost matrix. The NLL from (8) is extended as

$$\text{NLL}(\mathbb{Y}, f_{\text{PMB}}) \approx \int \lambda(y') dy' - \log\left(\sum_{q=1}^Q \prod_{y \in \mathbb{Y}^U(\mathbf{A}_q^*)} \lambda(y) \prod_{k=1}^m f_{B_k}(\mathbb{Y}_k(\mathbf{A}_q^*))\right), \quad (13)$$

where we define $\mathbb{Y}^U(\mathbf{A}_q^*) = \mathbb{Y} \setminus \cup_{i=1}^m \mathbb{Y}_i(\mathbf{A}_q^*)$, i.e., \mathbb{Y}^U contains all the ground truth elements matched to the PPP.

NLL Decomposition. Often the most likely assignment yields a good approximation to the NLL. For $Q = 1$, the NLL can be decomposed into four parts and expressed in terms of assignments

$$\begin{aligned} \text{NLL}(\mathbb{Y}, f_{\text{PMB}}) \approx \min_{\gamma \in \Gamma} & - \underbrace{\sum_{(i,j) \in \gamma} \log(r_i p_{i,\text{cls}}(c_j))}_{\text{Classification}} - \underbrace{\sum_{(i,j) \in \gamma} \log(p_{i,\text{reg}}(b_j))}_{\text{Regression}} \quad (14) \\ & - \underbrace{\sum_{i \in \mathbb{F}(\gamma)} \log(1 - r_i)}_{\text{False detections}} + \underbrace{\int \lambda(y') dy' - \sum_{j \in \mathbb{M}(\gamma)} \log \lambda(y_j)}_{\text{Missed objects}}, \end{aligned}$$

where Γ is the set of all possible assignment sets, $(i, j) \in \gamma$ means that prediction i has been assigned to ground truth j , and $\mathbb{F}(\gamma) = \{i \in \mathbb{N}_m \mid \nexists j : (i, j) \in \gamma\}$ is the set of indices of the Bernoullis not matched to any ground-truth, i.e. false positives. Note that we assume the classification and regression distributions are independent $p_i(x) = p_{i,\text{cls}}(\cdot)p_{i,\text{reg}}(\cdot)$ for this decomposition. Further, we define $\mathbb{M}(\gamma) = \{j \in \mathbb{N}_{|\mathbb{Y}|} \mid \nexists i : (i, j) \in \gamma\}$ as the set of indices of ground-truths not matched to any Bernoulli component, i.e., missed objects. This decomposition enables further insight into the types of errors made by an algorithm, e.g., instead of treating all false positives equally as in [11], we take their existence probability into account for deciding how much to penalize an algorithm.

4 Experiments

For our experiments, we evaluate three existing object detection models: DETR [2], RetinaNet [17], and Faster-RCNN [25], all using ResNet50 backbones. These are chosen to represent a set-based, one-stage, and two-stage detector, which highlights that the RFS framework is applicable regardless of architecture. All these models are publicly available through the Detectron2 [34] object detection framework, with hyperparameters³ optimized to produce competitive detection results for the COCO dataset [18]. Further, the models have previously been retrofitted with variance networks to estimate their spatial uncertainty [13]. Due to hardware limitations, models in [13] used a smaller batch size and adjusted learning rates, resulting in decreased mAP compared to numbers reported by Detectron2. For fair comparison, we use the same hyperparameters as [13], but note that increasing the batch size can improve mAP for all models.

The models are also fine-tuned with MB-NLL (8) as loss function. During training, the aim is to detect all objects, hence the PPP for undetected objects is ignored. We also found training to be more stable when the number of assignments Q is set to one. Further, when calculating the assignment costs for matching, ignoring spatial uncertainties improved training stability. That is, in (7c), we use the L2 distance instead of $\log(p_{k,\text{reg}})$. This can be thought of as learning the spatial uncertainty given the predicted mean of the bounding box.

³ Hyperparameters are used as is unless stated otherwise.

Table 1: mAP, PMB-NLL and PDQ with/without threshold for three detectors on the COCO validation set. * detections excluded due to ∞ NLL.

Detector	Loss	PMB-NLL ↓	mAP ↑	PDQ@F1 ↑	PDQ@0.0 ↑
DETR	ES	120.33	0.407	0.262	0.033
	NLL	152.13	0.376	0.113	0.014
	MB-NLL	124.20	0.389	0.271	0.023
RetinaNet	ES	127.66	0.362	0.228	0.028
	NLL	126.86	0.351	0.185	0.021
	MB-NLL	121.02	0.361	0.251	0.023
Faster-RCNN	ES*	140.53	0.373	0.281	0.087
	NLL*	139.08	0.371	0.282	0.087
	MB-NLL	117.77	0.326	0.199	0.024

For evaluation, the $Q = 25$ assignments with the highest likelihood are used to approximate the PMB-NLL, as larger values for Q do not affect the approximation for any of the models considered. In contrast to training, the matching cost is used as described in (12). Further, following the COCO standard, models are limited to 100 predictions and no confidence threshold is used. DETR is designed to provide exactly 100 predictions, while we apply NMS and keep the 100 top-scoring predictions for RetinaNet and Faster-RCNN. For all models, bounding boxes are parameterized by their top-left and bottom-right coordinates $[x_1, y_1, x_2, y_2]$. While the pre-trained models from [13] used a Gaussian distribution for regression, we found that using a Laplace distribution results in considerably lower NLL for both training and evaluation, across all models.

Evaluating Object Detection with Proper Scoring Rule. We report mAP, PDQ [11] and PMB-NLL in Table 1 and the decomposed results following (14) are shown in Table 2, with additional analysis in the supplementary material. For models with loss ES (energy score) and NLL, please refer to [13] for their details. PDQ is reported both when thresholding prediction confidences at the detectors’ optimal F1-score and without any thresholding. The decomposition in Table 2 is calculated for the assignment with the lowest NLL and shown averaged per image and per assigned objects. For instance, the DETR ES regression term 82.3/12.3 is read as, on average the regression distribution of *all* matched predictions contributes with 82.3 to the overall NLL. For a *single* matched prediction, it contributes with 12.3 to the total NLL on average.

We can see from Table 1 that optimizing the networks toward MB-NLL rather than the NLL formulation used in [13] consistently gives lower PMB-NLL at evaluation. With the exception of Faster-RCNN, this lower PMB-NLL is achieved without sacrificing mAP performance. We can also note that mAP does not indicate quality of predictive uncertainty. For instance, Faster-RCNN trained with the energy score achieves competitive performance in terms of mAP, but its uncertainty estimates result in the second worse PMB-NLL among the models. Further, although PDQ is described as a threshold-free performance measure,

Table 2: Decomposed PMB-NLL on COCO validation set. Numbers are given as [mean per image]/[mean per prediction]. FP=NLL of unmatched predictions. PPP match+PPP rate=missed objects. *detections excluded due to ∞ NLL.

Detector	Loss	Regression ↓	Classification ↓	FP ↓	PPP match ↓	PPP rate ↓
DETR	ES	82.3 / 12.3	3.79/0.57	17.6/0.71	15.6/ 23.3	1.46
	NLL	124.7/17.5	3.57/ 0.50	18.4/ 0.59	5.0 /23.8	1.52
	MB-NLL	93.3 /14.6	3.26/0.51	3.1 /0.99	24.8/25.4	0.18
RetinaNet	ES	103.1/14.4	7.95/1.11	10.1/ 0.22	4.3 /23.8	2.87
	NLL	105.0/14.5	7.87/1.09	9.8 / 0.22	2.3 / 20.7	2.91
	MB-NLL	79.9 / 13.9	4.00/ 0.70	2.6 /0.44	34.1/21.1	0.98
F-RCNN	ES*	105.5/15.1	8.23/1.18	12.9/0.57	14.0/37.1	0.43
	NLL*	104.7/15.0	8.23/1.18	13.2/0.58	13.5/36.1	0.43
	MB-NLL	62.0 / 12.3	4.94/ 0.98	3.3 / 0.36	46.8/ 20.4	1.30

it is sensitive to false positive detections regardless of their confidence, as predictions with low and high confidence receive the same penalty by PDQ. When including low confidence predictions (PDQ@0.0), the reported PDQ results become hard to distinguish between detectors. For PMB-NLL, FP penalties are instead proportional to the predicted existence probabilities.

Inspecting the decomposition of PMB-NLL in Table 2 gives further insights into the strengths and weaknesses of the detectors. Models that have not been trained with MB-NLL, show high penalties for producing many false positives. This is exemplified in Fig. 4, where the model trained with MB-NLL produces a single prediction per object, and fewer false detections. Rather than producing multiple plausible predictions per object, where each prediction has low spatial uncertainty, they are compiled into a single hypothesis with slightly larger uncertainty. More examples of this are available in the Appendix. We theorize that the ES and NLL training is optimized toward mAP evaluation, where low confidence predictions are not penalized as heavily, and that the MB-NLL loss instead encourages models to produce plausible set predictions. Comparing across architectures, we can see from FP penalties that RetinaNet generally assigns lower existence probabilities to its incorrect predictions compared to DETR. For the matched predictions, DETR instead has the strongest classification performance, indicating that DETR generally has higher existence probabilities for its predictions, regardless of them being assigned to a ground truth or not.

The example in Fig. 4 also underlines important advantages with our evaluation. For the person in the image, DETR ES has one prediction with good regression but low confidence (in turquoise with 0.4), and one confident (in white with 0.96) with bad regression. Depending on which prediction is assigned to the true object, the error is either related to classification or regression. As highlighted previously, confidence thresholds are in practice needed by PDQ and used explicitly in other methods [13]. Thus, existing methods only consider the confident prediction and report large regression errors. In contrast, PMB-NLL evaluates both possibilities and seamlessly weighs their contribution based on



Fig. 4: Ground truth (left), predictions from DETR ES (middle) and DETR MB-NLL (right). Predictions with a score less than 0.1 not shown. Models not trained with MB-NLL tend to produce many false positive detections.

their individual likelihood, where the most likely assignment is in fact the one with lower existence probability.

Further, it is interesting to study the balance between matched predictions, and ground truth objects matched to the PPP. For Faster-RCNN trained with MB-NLL, many objects are assigned to the PPP. However, its PPP is a reasonably good description of the missed objects, resulting in a total PMB-NLL which is lower than the other models. For an application where a high recall level is desirable, the RetinaNet ES model might be a better choice, at the cost of worse regression and classification performance.

5 Conclusions

We propose the use of random finite sets for probabilistic object detection. Instead of predicting a set of objects \mathbb{X} , we ask our models to predict the distribution over the set of objects $f(\mathbb{X})$. Using a distribution over sets enables us to evaluate model performance for the true set of objects \mathbb{Y} by applying the proper scoring rule negative log-likelihood $-\log(f(\mathbb{Y}))$. Our proposed method is general enough to be applied to detectors with any type of regression or classification distribution. It handles the assignments between predictions and objects automatically and can be decomposed into different error types. We evaluate three types of detectors using our new scoring rule and highlight their strengths and weaknesses. Our method enables fair comparison between probabilistic object detectors and we hope this will encourage the creation of novel architectures that aim for accurate uncertainty estimates rather than just accurate means. Future directions include how to better optimize networks toward our scoring rule and exploring further scoring rules within the random finite set framework.

Acknowledgements This work was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation. Computational resources were provided by the Swedish National Infrastructure for Computing at C3SE and NSC, partially funded by the Swedish Research Council, grant agreement no. 2018-05973.

References

1. Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U.R., et al.: A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion* **76**, 243–297 (2021)
2. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: *European Conference on Computer Vision*. pp. 213–229. Springer (2020)
3. Choi, J., Chun, D., Kim, H., Lee, H.J.: Gaussian yolov3: An accurate and fast object detector using localization uncertainty for autonomous driving. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 502–511 (2019)
4. Dawid, A.P., Musio, M.: Theory and applications of proper scoring rules. *Metron* **72**(2), 169–183 (2014)
5. Feng, D., Harakeh, A., Waslander, S.L., Dietmayer, K.: A review and comparative study on probabilistic object detection in autonomous driving. *IEEE Transactions on Intelligent Transportation Systems* (2021)
6. Feng, D., Wei, X., Rosenbaum, L., Maki, A., Dietmayer, K.: Deep active learning for efficient training of a lidar 3d object detector. In: *2019 IEEE Intelligent Vehicles Symposium (IV)*. pp. 667–674. IEEE (2019)
7. García-Fernández, Á.F., Williams, J.L., Granström, K., Svensson, L.: Poisson multi-Bernoulli mixture filter: direct derivation and implementation. *IEEE Transactions on Aerospace and Electronic Systems* **54**(4), 1883–1901 (2018)
8. Gawlikowski, J., Tassi, C.R.N., Ali, M., Lee, J., Humt, M., Feng, J., Kruspe, A., Triebel, R., Jung, P., Roscher, R., et al.: A survey of uncertainty in deep neural networks. *arXiv preprint arXiv:2107.03342* (2021)
9. Gneiting, T., Balabdaoui, F., Raftery, A.E.: Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **69**(2), 243–268 (2007)
10. Gneiting, T., Raftery, A.E.: Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association* **102**(477), 359–378 (2007)
11. Hall, D., Dayoub, F., Skinner, J., Zhang, H., Miller, D., Corke, P., Carneiro, G., Angelova, A., Sünderhauf, N.: Probabilistic object detection: Definition and evaluation. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. pp. 1031–1040 (2020)
12. Harakeh, A., Smart, M., Waslander, S.L.: Bayesod: A Bayesian approach for uncertainty estimation in deep object detectors. In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*. pp. 87–93. IEEE (2020)
13. Harakeh, A., Waslander, S.L.: Estimating and evaluating regression predictive uncertainty in deep object detectors. In: *International Conference on Learning Representations* (2021)
14. He, Y., Wang, J.: Deep mixture density network for probabilistic object detection. In: *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. pp. 10550–10555. IEEE (2020)
15. Kendall, A., Gal, Y.: What uncertainties do we need in bayesian deep learning for computer vision? In: *Advances in Neural Information Processing Systems*. vol. 30 (2017)
16. Lee, Y., Hwang, J.w., Kim, H.I., Yun, K., Kwon, Y.: Localization uncertainty estimation for anchor-free object detection. *arXiv preprint arXiv:2006.15607* (2020)

17. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: Proceedings of the IEEE international conference on computer vision. pp. 2980–2988 (2017)
18. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014)
19. Lyu, Z., Gutierrez, N., Rajguru, A., Beksi, W.J.: Probabilistic object detection via deep ensembles. In: European Conference on Computer Vision. pp. 67–75. Springer (2020)
20. Mahler, R.P.: Advances in statistical multisource-multitarget information fusion. Artech House (2014)
21. Miller, D., Nicholson, L., Dayoub, F., Sünderhauf, N.: Dropout sampling for robust object detection in open-set conditions. In: 2018 IEEE International Conference on Robotics and Automation (ICRA). pp. 3243–3249. IEEE (2018)
22. Motro, M., Ghosh, J.: Scaling data association for hypothesis-oriented MHT. In: 2019 22th International Conference on Information Fusion (FUSION). pp. 1–8. IEEE (2019)
23. Murty, K.G.: An algorithm for ranking all the assignments in order of increasing cost. Operations research **16**(3), 682–687 (1968)
24. Pinto, J., Xia, Y., Svensson, L., Wymeersch, H.: An uncertainty-aware performance measure for multi-object tracking. IEEE Signal Processing Letters **28**, 1689–1693 (2021)
25. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. In: Advances in Neural Information Processing Systems. vol. 28 (2015)
26. Skaftø, N., Jørgensen, M., Hauberg, S.: Reliable training and estimation of variance networks. In: Advances in Neural Information Processing Systems. vol. 32 (2019)
27. Sun, Z., Cao, S., Yang, Y., Kitani, K.M.: Rethinking transformer-based set prediction for object detection. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3611–3620 (2021)
28. Valdenegro-Toro, M.: I find your lack of uncertainty in computer vision disturbing. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1263–1272 (2021)
29. Vo, B., Vo, B., Hoang, H.G.: An efficient implementation of the generalized labeled multi-Bernoulli filter. IEEE Transactions on Signal Processing **65**(8), 1975–1987 (2017)
30. Vo, B.T., Vo, B.N.: Labeled random finite sets and multi-object conjugate priors. IEEE Transactions on Signal Processing **61**(13), 3460–3475 (2013)
31. Williams, J.L.: Hybrid Poisson and multi-Bernoulli filters. In: 2012 15th International Conference on Information Fusion. pp. 1103–1110. IEEE (2012)
32. Williams, J.L.: Marginal multi-Bernoulli filters: RFS derivation of MHT, JIPDA, and association-based MeMBer. IEEE Transactions on Aerospace and Electronic Systems **51**(3), 1664–1687 (2015)
33. Wirges, S., Reith-Braun, M., Lauer, M., Stiller, C.: Capturing object detection uncertainty in multi-layer grid maps. In: 2019 IEEE Intelligent Vehicles Symposium (IV). pp. 1520–1526. IEEE (2019)
34. Wu, Y., Kirillov, A., Massa, F., Lo, W.Y., Girshick, R.: Detectron2. <https://github.com/facebookresearch/detectron2> (2019)
35. Zhang, Y., Hare, J., Prugel-Bennett, A.: Deep set prediction networks. In: Advances in Neural Information Processing Systems. vol. 32 (2019)

36. Zhou, X., Koltun, V., Krähenbühl, P.: Probabilistic two-stage detection. arXiv preprint arXiv:2103.07461 (2021)
37. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable DETR: Deformable transformers for end-to-end object detection. In: International Conference on Learning Representations (2021)

A Cost Matrix

Suppose that we have a PMB density with Poisson intensity $\lambda(\cdot)$ and with m Bernoulli components where the i -th Bernoulli component has probability of existence r_i and existence-conditioned object density $p_i(\cdot)$. Note that we can model an MB as a PMB with Poisson intensity $\lambda(\cdot) = 0$, which means that it is enough to describe how to handle PMB densities. To evaluate the multi-object density of a PMB, we have to calculate all possible assignments, which can be computationally intractable when working with many elements. However, we can approximate the PMB likelihood by only considering the assignments with the highest likelihood. This section shows how to find these assignments by solving an optimal assignment problem and how to select the corresponding cost matrix.

Before formulating the optimal assignment problem, we remind ourselves of the problem setting. For an object with state $y_j \in \mathbb{Y} = \{y_1, \dots, y_n\}$, the single object likelihood is proportional to $\lambda(y_j)$ if it is associated to the PPP and proportional to $r_i p_i(y_j)$ if it is associated to the i -th Bernoulli component. If this Bernoulli component is not associated to any object states, then the likelihood is $1 - r_i$.

Next, to formulate the optimization problem of finding the assignment with highest likelihood we introduce an association variable. Define a surjective association $\theta : \{1, \dots, n\} \rightarrow \{0, 1, \dots, m\}$ such that $\theta(i) = \theta(j) \in \{1, \dots, m\}$ if and only if $i = j$. If $\theta(j) = 0$, object y_j is associated to the PPP, while $\theta(j) = i > 0$ means that object state y_j is associated to the i -th Bernoulli component. Further, let Θ be the set of all such θ . Then, the PMB likelihood for the set of objects \mathbb{Y} can be expressed as

$$\begin{aligned} f_{\text{PMB}}(\mathbb{Y}) &= \sum_{\theta \in \Theta} \prod_{j: \theta(j) > 0} r_{\theta(j)} p_{\theta(j)}(y_j) \prod_{i: \nexists \theta(j)=i \forall j} 1 - r_i \prod_{j: \theta(j)=0} \lambda(y_j) \exp(-\bar{\lambda}), \\ &\propto \sum_{\theta \in \Theta} \prod_{j: \theta(j) > 0} r_{\theta(j)} p_{\theta(j)}(y_j) \prod_{i: \nexists \theta(j)=i \forall j} 1 - r_i \prod_{j: \theta(j)=0} \lambda(y_j). \end{aligned} \quad (15)$$

When searching for the most likely associations θ , we disregard $\exp(-\bar{\lambda}) = \exp(-\int \lambda(y') dy')$ since it is independent of θ .

We note that (15) captures the association of every Bernoulli component. If the i -th Bernoulli component does not appear in the second product in (15), then it must appear in the first product in (15). We also note that the factor $\prod_{i=1}^m (1 - r_i)$ is independent of the association θ . This means that dividing (15) by $\prod_{i=1}^m (1 - r_i)$ yields

$$f_{\text{PMB}}(\mathbb{Y}) \propto \sum_{\theta \in \Theta} \prod_{j: \theta(j) > 0} \frac{r_{\theta(j)} p_{\theta(j)}(y_j)}{1 - r_{\theta(j)}} \prod_{j: \theta(j)=0} \lambda(y_j), \quad (16)$$

and the association that maximizes the PMB likelihood can be found as

$$\theta^* = \arg \max_{\theta} \prod_{j: \theta(j) > 0} \frac{r_{\theta(j)} p_{\theta(j)}(y_j)}{1 - r_{\theta(j)}} \prod_{j: \theta(j)=0} \lambda(y_j). \quad (17)$$

Equivalently, we can search for the association that minimises the negative logarithm of (17), which gives us

$$\theta^* = \arg \min_{\theta} - \sum_{j: \theta(j) > 0} \log \frac{r_{\theta(j)} p_{\theta(j)}(y_j)}{1 - r_{\theta(j)}} - \sum_{j: \theta(j) = 0} \log \lambda(y_j). \quad (18)$$

In order to formulate this minimization as an optimal assignment problem, we make a slight change in notation. Each association map θ can be represented by a $(m+n) \times n$ assignment matrix A consisting of 0s or 1s. There is a bijective mapping between θ and A : for $i = 1, \dots, m$, $j = 1, \dots, n$, $A_{i,j} = 1$ if and only if $\theta(j) = i$. For $i = m+j$, $j = 1, \dots, n$, $A_{i,j} = 1$ if and only if $\theta(j) = 0$, whereas $A_{i,j}$ is always 0 when $i > m$ and $i \neq m+j$. The assignment matrix hence must satisfy the constraints $\sum_i A_{i,j} = 1$, $\forall j$, and $\sum_j A_{i,j} \leq 1$, $\forall i$.

The cost matrix of the corresponding assignment matrix A is the $(m+n) \times n$ matrix C where

$$C_{i,j} = -\log \frac{r_i p_i(y_j)}{1 - r_i}, \quad i = 1, \dots, m, \quad j = 1, \dots, n, \quad (19)$$

and where the lower part of C is a “diagonal” matrix such that $C_{m+j,j} = -\log \lambda(y_j)$, $j = 1, \dots, n$, and entries not on the diagonal are ∞ . The cost of assignment matrix A is then given by the Frobenius inner product

$$\text{trace}(A^T C) = \sum_{i=1}^{m+n} \sum_{j=1}^n C_{i,j} A_{i,j}, \quad (20)$$

and the problem of finding the optimal assignment A^* becomes

$$A^* = \arg \min_A \sum_{i=1}^{m+n} \sum_{j=1}^n C_{i,j} A_{i,j} \quad (21a)$$

$$\text{s.t.} \quad \sum_{i=1}^{m+n} A_{i,j} = 1, \quad \sum_{j=1}^n A_{i,j} \leq 1, \quad (21b)$$

$$C_{i,j} = \begin{cases} -\log \left(\frac{p_i(y_j)}{1 - r_i} r_i \right) & \text{if } i \leq m, \\ -\log \lambda(y_j) & \text{if } i = m+j, \\ \infty & \text{otherwise.} \end{cases} \quad (21c)$$

B Experimental Details

B.1 Model Implementation

For implementing the DETR, RetinaNet and Faster-RCNN models, we used the probabilistic extension [13] of the Detectron2 [34] object detection framework. In that framework, models are trained to predict the covariance matrix Σ_b for the corresponding bounding box b . Specifically, models output the parameters of a

lower triangular matrix \mathbf{L} of the Cholesky decomposition $\Sigma_b = \mathbf{L}\mathbf{L}^T$. While originally trained with a Gaussian distribution, we found that using an independent Laplace distribution for each parameter in b yielded better results. Using the diagonal elements of \mathbf{L} as $[\sigma_1, \sigma_2, \sigma_3, \sigma_4]$, we find the scale of each Laplace distribution as $s_i = \sigma_i/\sqrt{2}$. The choice of a diagonal matrix is partially motivated by the evaluations in [13]. While their study was limited to Gaussian distributions, they found that diagonal covariance matrices perform on par with, or better than, their full equivalent.

B.2 Training Details

Fine-tuning toward MB-NLL was done given the pre-trained weights in [13]⁴. For DETR and Faster-RCNN, models trained with ES were used as a starting point, while the model trained with NLL was used for RetinaNet. Faster-RCNN and RetinaNet were fine-tuned for 135,000 gradient steps, where the learning rate was dropped by a factor of 10 at 105,000 iterations, and again at 125,000 iterations. The initial learning rate was set to 0.001 for RetinaNet and 0.0025 for Faster-RCNN. DETR was also fine-tuned for 135,000 iterations, but with an initial learning rate of $5 \cdot 10^{-5}$ and with learning rate drops at 60,000 and 100,000 iterations. Otherwise, no changes to hyperparameters from the standard Detectron2 framework were done.

As both RetinaNet and Faster-RCNN rely on non-maximum suppression to remove duplicate detections, we applied NMS when training with the MB-NLL loss. Here, we used the standard IoU threshold of 0.5 and used the top 100 detections. For DETR, this was not necessary since it predicts the set of objects directly.

B.3 Inference Details

Following the COCO standard, detectors are limited to 100 predictions per image. We do not apply any confidence thresholding, but for RetinaNet and Faster-RCNN the 100 predictions with highest existence probability after NMS are used. For DETR, no selection is needed as the model only produces 100 predictions per image.

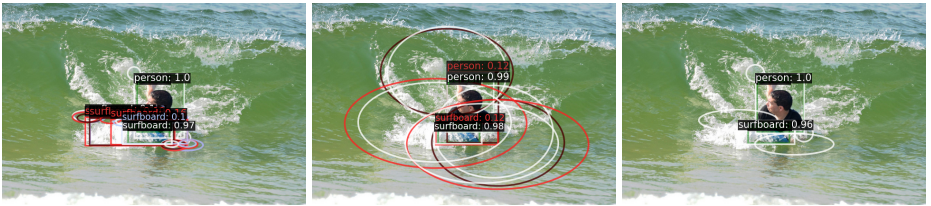
C Additional Results

C.1 Qualitative Results

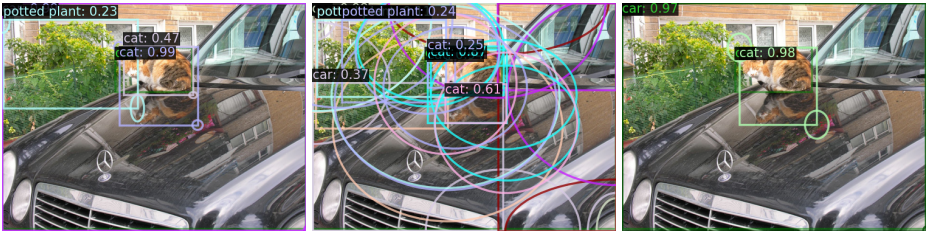
In addition to the example detections shown in Section 4, we provide further examples for DETR, RetinaNet and Faster-RCNN in figures 5, 6 and 7. We can identify similar trends in these examples as in the ones described in our results. First, Fig. 7b shows additional examples where the assignment is ambiguous. There, the cat predictions for Faster-RCNN trained with MB-NLL either have

⁴ <https://github.com/asharakeh/probdet>

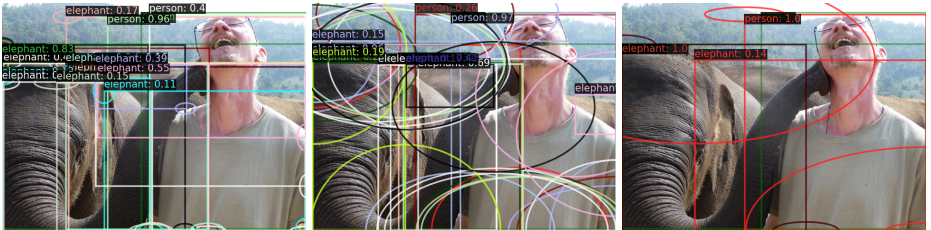
large regression or classification errors, depending on which one of them is assigned to the true object. Second, we see that MB-NLL reduces the number of confident false detections across all models. When comparing the models trained with ES and MB-NLL, the reduction of false detections can also be interpreted as a different representation of spatial uncertainty. In both Fig. 5a and Fig. 6a, there is uncertainty in where the surfboard ends on the left side. The MB-NLL models have a single prediction with larger regression uncertainty, while the ES models have many detections, each with relatively small spatial uncertainty. Further, the MB-NLL loss is normalized with the number of predictions during training.



(a) Example 1.



(b) Example 2.

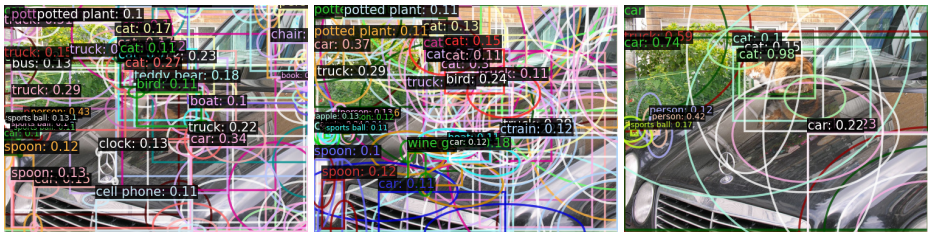


(c) Example 3.

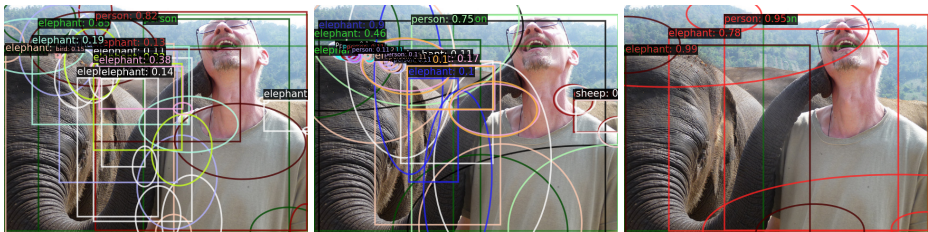
Fig. 5: Examples from COCO validation data with predictions made by DETR detectors trained with ES (left), NLL (middle), and MB-NLL (right). True objects are shown in green and without confidence values. Predictions with $r < 0.1$ are not shown.



(a) Example 1.

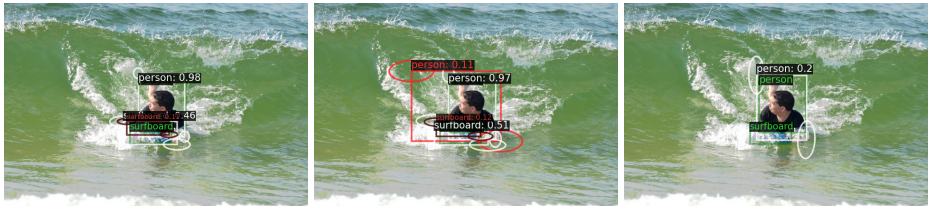


(b) Example 2.

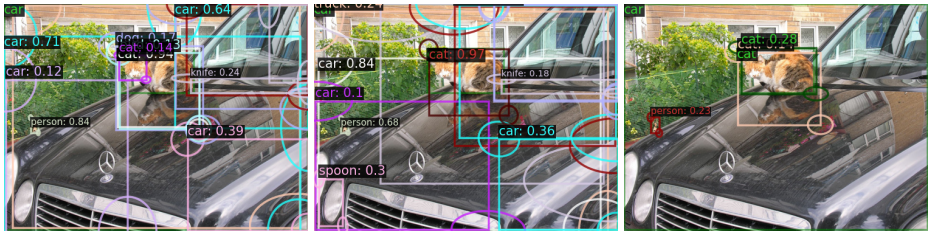


(c) Example 3.

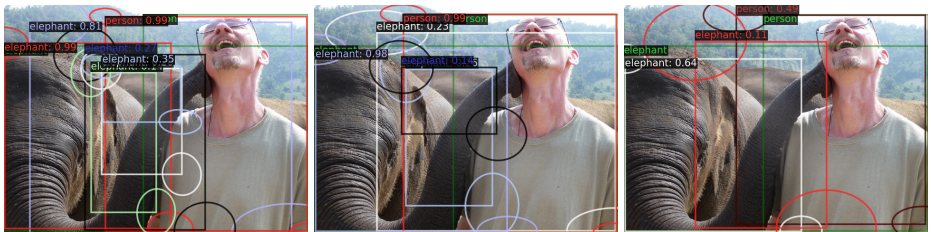
Fig. 6: Examples from COCO validation data with predictions made by RetinaNet detectors trained with ES (left), NLL (middle), and MB-NLL (right). True objects are shown in green and without confidence values. Predictions with $r < 0.1$ are not shown.



(a) Example 1.



(b) Example 2.



(c) Example 3.

Fig. 7: Examples from COCO validation data with predictions made by Faster-RCNN detectors trained with ES (left), NLL (middle), and MB-NLL (right). True objects are shown in green and without confidence values. Predictions with $r < 0.1$ are not shown.

C.2 PMB-NLL Decomposition

To complement the PMB-NLL decomposition in Table 2, we also provide corresponding histograms for DETR, RetinaNet and Faster-RCNN in figures 8, 9 and 10. The histograms show how predictions contribute to the overall PMB-NLL in the assignment with the highest likelihood. Further, for matched predictions, histograms are also decomposed based on the size of the true object. Here, we follow the COCO standard for defining small, medium and large objects. For the regression of matched Bernoullis and PPP, the values have been limited to 40 for enhanced visualizations. For the classification of matched Bernoullis, the upper limit is set to 3.

Generally, the models are worse at detecting small objects, which is shown by them being assigned to the PPP more often than medium or large objects. Further, the detectors are more confident when predicting larger objects, as can be seen from the classification histograms over matched Bernoullis. This is of course expected as large objects are inherently easier to detect. Lastly, we can observe that the histograms for models trained with ES or NLL tend to have more outliers, i.e., predictions whose values have been clipped to the visualization limits.

D Comparison to DETR loss

The DETR object detector [2] popularized the concept of treating object detection as a direct set prediction task. Similar to our method, they rely on a one-to-one matching between predictions and ground truth objects. We aim to compare their formulation to ours, highlighting similarities and differences.

D.1 DETR Loss Revisited

We start by reviewing the methods used in DETR. To find the matching between predictions and true objects, they rely on the Hungarian algorithm to minimize

$$\hat{\sigma} = \arg \min_{\sigma \in \mathfrak{S}_N} \sum_i^N \mathcal{L}_{\text{match}}(y_i, \hat{y}_{\sigma(i)}), \quad (22)$$

where $\hat{y} = \{\hat{y}_i\}_{i=1}^N$ is the set of predictions and y is the ground truth set of objects, where we assume y to be padded to size N with \emptyset (no object). An object $y_i = (c_i, b_i)$ consists of a class label (which can be \emptyset) and a bounding box $b_i \in \mathbb{R}^4$, while a prediction $\hat{y}_i = (\hat{p}_{i,\text{cls}}(c_i), \hat{b}_i)$ consists of a class distribution $\hat{p}_{i,\text{cls}}(c_i)$ which assigns probabilities to all classes including \emptyset and a predicted bounding box \hat{b}_i . Further, \mathfrak{S}_N is the set of all permutations of N elements and $\mathcal{L}_{\text{match}}(y_i, \hat{y}_{\sigma(i)})$ is a pair-wise matching cost defined as

$$\mathcal{L}_{\text{match}}(y_i, \hat{y}_{\sigma(i)}) = -\mathbb{1}_{c_i \neq \emptyset} \hat{p}_{\sigma(i), \text{cls}}(c_i) + \mathbb{1}_{c_i \neq \emptyset} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\sigma(i)}), \quad (23)$$

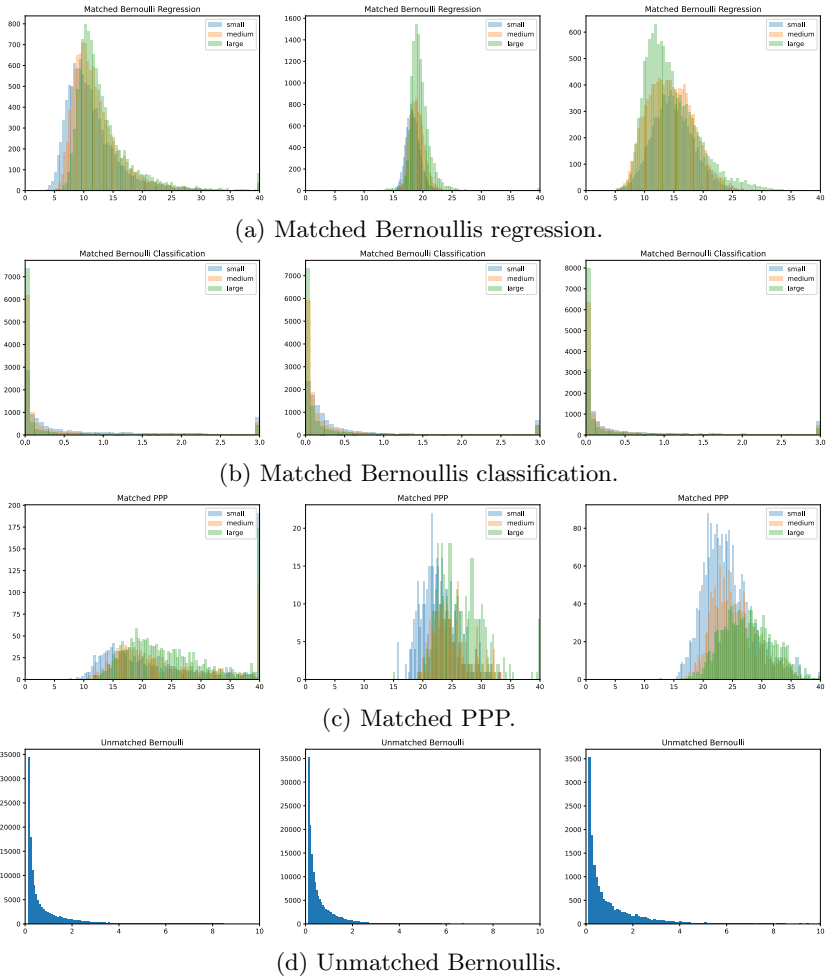


Fig. 8: Histograms over PMB-NLL decomposition for DETR trained with different loss functions: ES (left), NLL (middle), and MB-NLL (right). Note varying y-axes across models.

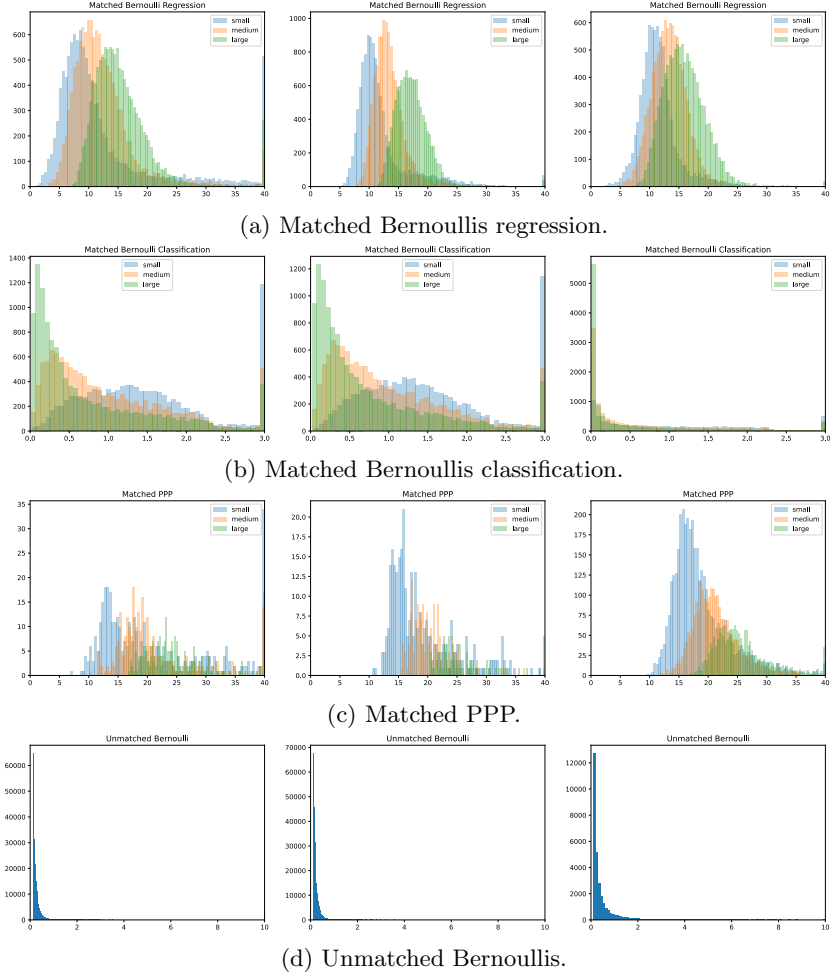
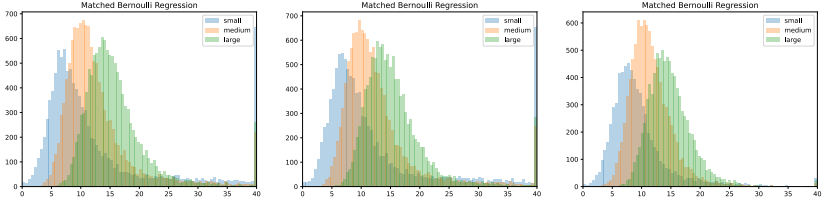
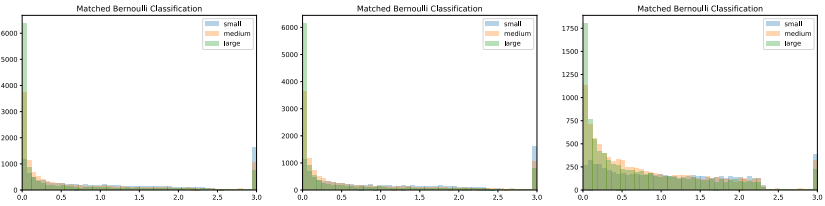


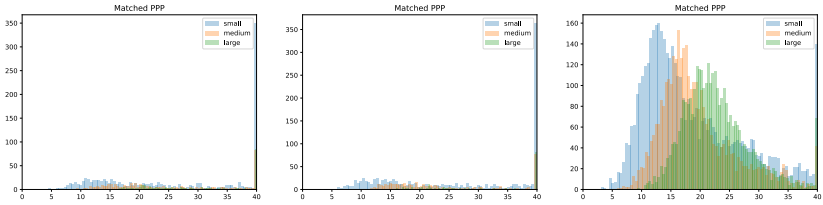
Fig. 9: Histograms over PMB-NLL decomposition for Retinanet trained with different loss functions: ES (left), NLL (middle), and MB-NLL (right). Note varying y-axes across models.



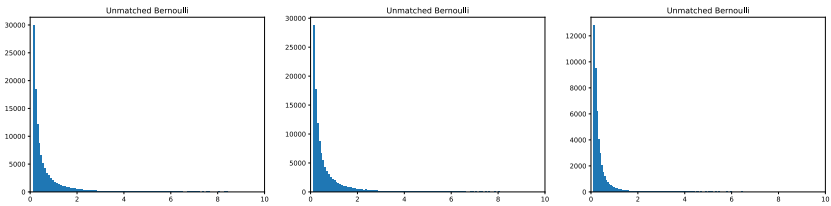
(a) Matched Bernoullis regression.



(b) Matched Bernoullis classification.



(c) Matched PPP.



(d) Unmatched Bernoullis.

Fig. 10: Histograms over PMB-NLL decomposition for Faster-RCNN trained with different loss functions: ES (left), NLL (middle), and MB-NLL (right). Note varying y-axes across models.

with

$$\mathcal{L}_{\text{box}}(b_i, \hat{b}_{\sigma(i)}) = \lambda_{\text{iou}} \mathcal{L}_{\text{iou}}(b_i, \hat{b}_{\sigma(i)}) + \lambda_{\text{L1}} \|b_i - \hat{b}_{\sigma(i)}\|_1, \quad (24)$$

where hyperparameters are typically set to $\lambda_{\text{iou}} = 2$ and $\lambda_{\text{L1}} = 5$. We can note that the matching cost for $c_i = \emptyset$ is zero.

Given the optimal matching $\hat{\sigma}$, the final loss is calculated as

$$\mathcal{L}_{\text{Hungarian}}(y, \hat{y}) = \sum_{i=1}^N \left[-\log(\hat{p}_{\hat{\sigma}(i), \text{cls}}(c_i)) + \mathbf{1}_{c_i \neq \emptyset} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)}) \right]. \quad (25)$$

However, rather than using the negative log-likelihood for class predictions, the log-probability term is down-weighted by a factor 10 when $c_i = \emptyset$. This is motivated in [2] to handle class imbalance.

D.2 MB-NLL Relation to DETR

This section aims at describing the MB-NLL using the same notation as DETR; to simplify the comparison we focus on MB-NLL rather than PMB-NLL. We start by comparing the matching costs before moving on to the loss formulation.

Matching Cost. We can compare (23) with the cost of matching objects to Bernoulli predictions used in this work. As described in Section 3.1, we use $r = 1 - \hat{p}_{\text{cls}}(\emptyset)$. Further, for the Bernoulli predictions, the class distribution is assumed to be conditioned on existence, i.e., it has non-zero probability only for foreground classes. Hence, to clarify the relation to DETR we define

$$p_{\text{cls}}(c) = \begin{cases} 0 & \text{if } c = \emptyset, \\ \hat{p}_{\text{cls}}(c)/r & \text{otherwise,} \end{cases} \quad (26)$$

as the class distribution over foreground classes. We scale the predicted distribution $\hat{p}_{\text{cls}}(c)$ by $\frac{1}{r}$ such that $p_{\text{cls}}(c)$ becomes a proper distribution and fulfills $\sum_c p_{\text{cls}}(c) = 1$.

In Appendix A, we showed how to find the assignment that maximizes the likelihood (15) and consequently minimizes the negative log-likelihood. To enable easier comparison, we want to use the same notation as DETR and formulate a minimization over permutations

$$\hat{\sigma} = \arg \min_{\sigma \in \mathfrak{S}_N} \sum_i^N \mathcal{L}_{\text{match}, \text{MB}}(y_i, p_{\sigma(i)}), \quad (27)$$

where $\mathcal{L}_{\text{match}, \text{MB}}$ is a pair-wise matching cost. Using (15), we can express $\mathcal{L}_{\text{match}, \text{MB}}$ as

$$\begin{aligned} \mathcal{L}_{\text{match}, \text{MB}}(y_i, p_{\sigma(i)}) &= -\mathbf{1}_{c_i \neq \emptyset} (\log(r_{\sigma(i)} p_{\sigma(i)}(y_i))) - \mathbf{1}_{c_i = \emptyset} \log(1 - r_{\sigma(i)}) \\ &= -\log(\hat{p}_{\sigma(i), \text{cls}}(c_i)) - \mathbf{1}_{c_i \neq \emptyset} \log(p_{\sigma(i), \text{reg}}(b_i)), \end{aligned} \quad (28)$$

since the cost of assigning a prediction to a true object is $-\log(r_{\sigma(i)}p_{\sigma(i)}(y_i))$, while the cost of assigning it to background is $-\log(1 - r_{\sigma(i)})$. In (28), we have also used the fact that $p_i(y_i) = p_{i,\text{cls}}(c_i)p_{i,\text{reg}}(b_i)$, and the relation

$$\hat{p}_{i,\text{cls}}(c_i) = \begin{cases} r p_{i,\text{cls}}(c_i) & \text{if } c_i \neq \emptyset \\ 1 - r & \text{if } c_i = \emptyset, \end{cases} \quad (29)$$

to obtain an expression that resembles (23) and (25).

Further, if we assume $p_{\sigma(i),\text{reg}}(b_i)$ to be a Laplace distribution with independent box parameters $b = [b^1, b^2, b^3, b^4]$, with means $\hat{b} = [\hat{b}^1, \hat{b}^2, \hat{b}^3, \hat{b}^4]$ and scales $\hat{s} = [\hat{s}^1, \hat{s}^2, \hat{s}^3, \hat{s}^4]$, we can rewrite

$$\begin{aligned} -\log(p_{\sigma(i),\text{reg}}(b_i)) &= -\log\left(\prod_{k=1}^4 \frac{1}{2\hat{s}_{\sigma(i)}^k} \exp\left(-\frac{|b_i^k - \hat{b}_{\sigma(i)}^k|}{\hat{s}_{\sigma(i)}^k}\right)\right), \\ &= -\sum_{k=1}^4 \log\left(\frac{1}{2\hat{s}_{\sigma(i)}^k} \exp\left(-\frac{|b_i^k - \hat{b}_{\sigma(i)}^k|}{\hat{s}_{\sigma(i)}^k}\right)\right), \\ &= \sum_{k=1}^4 \frac{|b_i^k - \hat{b}_{\sigma(i)}^k|}{\hat{s}_{\sigma(i)}^k} + \log(\hat{s}_{\sigma(i)}^k) + \log(2). \end{aligned} \quad (30)$$

As $\log(2)$ is present in all matching costs between pairs of predictions and true objects, it does not affect the optimal assignment and can be disregarded. Further, if we let $\hat{s}_i^k = s, \forall i, k$, we can use the same argument to disregard $\log(\hat{s}_{\sigma(i)}^k)$ and replace $-\log(p_{j,\text{reg}}(b_i))$ in (28) with $\|b_i - \hat{b}_{\sigma(i)}\|_1/s$, obtaining

$$\mathcal{L}_{\text{match,MB}}(y_i, p_{\sigma(i)}) = -\log(\hat{p}_{\sigma(i),\text{cls}}(c_i)) + \mathbf{1}_{c_i \neq \emptyset} \|b_i - \hat{b}_{\sigma(i)}\|_1/s. \quad (31)$$

We can now compare the expressions for the matching losses in (31) and (23), used by MB-NLL and DETR, respectively, and analyze their similarities and differences. Rather than using the log-probabilities, the original DETR matching loss uses the class probabilities directly. In [2] this is motivated as making the classification part of the cost comparable to the regression part. Interestingly, we find that the classification log-probability is comparable to the L1 regression under the constant scale assumption. Further, the classification cost in (23) only evaluates the probability of the true class when the object is not \emptyset . For background, the cost is set to zero. In contrast, the MB matching cost also considers the log-probability of the prediction being background and favours predictions for which the probability of background $\hat{p}_{\text{cls}}(\emptyset)$ is small. The reason for also considering the cost of assigning predictions to background in MB-NLL is that predictions with large existence probabilities infer large penalties in the final loss function if they are not assigned to a true object.

We further highlight the difference in how $\hat{p}_{\text{cls}}(\emptyset)$ is handled by the two matching costs with an example. Imagine a scenario that contains two predictions and one true object, a car. Both predictions have the same regression error,

but differ in their classification. Suppose that both predictions have the same $\hat{p}_{\text{cls}}(\text{car}) = rp_{\text{cls}}(\text{car})$ but that the first prediction has a small r and a large $p_{\text{cls}}(\text{car})$ whereas the second prediction has a large r and a small $p_{\text{cls}}(\text{car})$. In (23), both these predictions would be treated as equally good. For MB-NLL however, it is better to assign the second prediction to the car, simply because that implies that the first prediction, which has a small r , is assigned to the background.

For the regression part, (23) contains both an additional IoU-loss, and the hyperparameters $\lambda_{\text{iou}}, \lambda_{\text{L1}}$ when compared to (31). While the IoU-part has no related term in (31), we find an inverse relationship for λ_{L1} and the assumed constant Laplace scale s , i.e., $\lambda_{\text{L1}} = \frac{1}{s}$. Thus, the choice $\lambda_{\text{L1}} = 5$ is equivalent of assuming a constant Laplace scale $s = 0.2$, and increasing λ_{L1} translates to assuming smaller spatial uncertainties.

Loss Function. The MB-NLL training loss

$$\mathcal{L}_{\text{MB}}(y, \hat{y}) = \sum_i^N \left[-\log(\hat{p}_{\hat{\sigma}(i)}(c_i)) + \mathbf{1}_{c_i \neq \emptyset} \sum_{k=1}^4 \frac{|b_i^k - \hat{b}_{\sigma(i)}^k|}{\hat{s}_{\sigma(i)}^k} + \log(\hat{s}_{\sigma(i)}^k) \right], \quad (32)$$

is very similar to the matching loss, but makes use of the scaling parameters $\hat{s}_{\sigma(i)}^k$ predicted by our networks, whereas the matching loss uses a fixed scaling parameter s .

We see that both (25) and (32) contain two terms, one for classification and one for regression. However, in contrast to the original DETR loss (25), we do not down-weight the log-probability when predictions are assigned to \emptyset . We believe this is one of the contributing factors to the reduced number of false detections when training with the MB-NLL loss. By down-weighting the penalty for predictions assigned to \emptyset , the detector is encouraged to produce artificially high classification confidence. In mAP, the measure that DETR most likely has been optimized toward, the high classification confidence is generally not penalized as it only relies on the ranked confidences among predictions and not the absolute confidence values.

Similar to the matching cost, (32) lacks the IoU cost found in (25), but has an identical L1-loss when $\lambda_{\text{L1}} = \frac{1}{\hat{s}_{\sigma(i)}^1} = \dots = \frac{1}{\hat{s}_{\sigma(i)}^4}$. Naturally, L1-losses increase with larger λ_{L1} , but using (32) this can also be explained as assuming smaller regression uncertainties. Finally, when training with MB-NLL, the relation between the Laplace scale and λ_{L1} also shows how the network can self-regulate the regression penalties. For the network to choose suitable scales, it has to balance the two terms $\log(\hat{s}_{\sigma(i)}^k)$ and $\frac{|b_i^k - \hat{b}_{\sigma(i)}^k|}{\hat{s}_{\sigma(i)}^k}$. While the first term encourages smaller scales, choosing too small scales may yield a large cost if the regression performance $|b_i^k - \hat{b}_{\sigma(i)}^k|$ is poor.