

CertainNet: Sampling-free Uncertainty Estimation for Object Detection

Stefano Gasperini* Jan Haug* Mohammad-Ali Nikouei Mahani Alvaro Marcos-Ramiro
Nassir Navab Benjamin Busam Federico Tombari

Abstract—Estimating the uncertainty of a neural network plays a fundamental role in safety-critical settings. In perception for autonomous driving, measuring the uncertainty means providing additional calibrated information to downstream tasks, such as path planning, that can use it towards safe navigation. In this work, we propose a novel sampling-free uncertainty estimation method for object detection. We call it CertainNet, and it is the first to provide separate uncertainties for each output signal: objectness, class, location and size. To achieve this, we propose an uncertainty-aware heatmap, and exploit the neighboring bounding boxes provided by the detector at inference time. We evaluate the detection performance and the quality of the different uncertainty estimates separately, also with challenging out-of-domain samples: BDD100K and nuImages with models trained on KITTI. Additionally, we propose a new metric to evaluate location and size uncertainties. When transferring to unseen datasets, CertainNet generalizes substantially better than previous methods and an ensemble, while being real-time and providing high quality and comprehensive uncertainty estimates.

Index Terms—Computer Vision for Transportation, Deep Learning for Visual Perception, Uncertainty Estimation, Object Detection.

I. INTRODUCTION

WHILE neural networks have been widely employed in research and mobile applications, their usage is still limited in safety-critical real-world settings, such as

Manuscript received: September, 9, 2021; Accepted November, 7, 2021.

This paper was recommended for publication by Editor Cesar Cadena upon evaluation of the Associate Editor and Reviewers' comments.

This work was partly sponsored by the German Federal Ministry for Economic Affairs and Energy (grant number: 19A19005B) through the VDA KI-Absicherung project.

* The authors contributed equally.

Stefano Gasperini and Jan Haug are with the Faculty of Computer Science, Technische Universität München (TUM), 85748 Garching bei München, Germany, and also with BMW Group, 80788 München, Germany (stefano.gasperini@tum.de, jan.haug@tum.de).

Mohammad-Ali Nikouei Mahani was with BMW Group when this work was conducted, 80788 München, Germany (man.mahani@gmail.com).

Alvaro Marcos-Ramiro is with BMW Group, 80788 München, Germany (alvaro.marcos-ramiro@bmw.de).

Nassir Navab is with the Faculty of Computer Science, Technische Universität München (TUM), 85748 Garching bei München, Germany, and also with the Department of Computer Science, Johns Hopkins University, Baltimore, 21218 MD USA (nassir.navab@tum.de).

Benjamin Busam is with the Faculty of Computer Science, Technische Universität München (TUM), 85748 Garching bei München, Germany (b.busam@tum.de).

Federico Tombari is with the Faculty of Computer Science, Technische Universität München (TUM), 85748 Garching bei München, Germany, and also with Google, 8002 Zürich, Switzerland (tombari@in.tum.de).

This is a preprint of the article accepted at Robotics and Automation Letters (RA-L). © 2021 IEEE.



Fig. 1. Example predictions of the proposed CertainNet from the KITTI validation set [6]. Uncertainties for objectness, location and size are shown as the complement of the value at the top left of each box, the plus around the center, and dashed boxes respectively. The object class is color-coded.

autonomous driving [1]. The ability of a model to estimate the uncertainty of its predictions is a key enabler for its safe and reliable use in these contexts and unknown scenarios [2]. This makes uncertainty estimation a fundamental companion to object detection [3], semantic segmentation [4], visual odometry [5], and other relevant visual tasks, especially to handle out-of-domain data.

As one of the main perception tasks, object detection comprises a plurality of outputs (e.g., object location, size, and class), requiring both regression and classification, and rendering it particularly challenging for uncertainty estimation. Existing real-time approaches have focused on providing an uncertainty estimate only on a subset of these detection outputs (e.g., location and size [3, 7]). In traffic, there can be situations where an object's location and size are relatively certain, while its existence in the scene (i.e., objectness) is uncertain, e.g., a person in a dark area or entering a car with an unusual pose, as shown at the bottom of Figure 1. In these cases, it is particularly important to assess the objectness uncertainty, to avoid ignoring the detection solely based on a low confidence score. A method confidence is a score which does not correspond to the actual probability of being correct, while uncertainty reduces this gap [1].

However, prior works cover this only partially, as they leave behind the uncertainty over important parameters, such as the objectness, posing a potential safety issue. By extracting a dedicated uncertainty measure for each of the output signals, a system would provide deeper and better insights on the

scene, which could be exploited by downstream tasks (e.g., path planning). Furthermore, incorporating uncertainty can reduce the accuracy of the models [8], or significantly increase the latency [9, 10], hence leading to a challenging trade-off between accuracy, safety, and runtime.

In this work, we address these issues with a novel method that provides separate and dedicated uncertainties for all detection outputs, while preserving speed and accuracy. We name our approach CertainNet, and our contributions can be summarized as follows:

- We introduce a novel sampling-free method to estimate the model uncertainty of an object detector.
- We propose the first real-time approach to estimate the uncertainty of all signals of the detection output (i.e., objectness, location, size and class).
- We present new simple metrics to assess the quality of size and location uncertainties for object detection.

Additionally, we provide extensive evaluations and comparisons on three challenging public datasets, including out-of-domain data, reaching competitive accuracy, while improving significantly on uncertainty-aware metrics.

II. RELATED WORK

A. Uncertainty Estimation for Neural Networks

A learning-based system can express two kinds of uncertainty: the epistemic uncertainty caused by the model, and the aleatoric uncertainty due to the input data [1]. The former can arise from errors in the training procedure (e.g., ignoring data imbalance), a sub-optimal architecture or the knowledge gap with out-of-domain data, while the latter originates from the data itself, e.g., due to sensor noise, coarse approximations, ambiguous samples, or also inaccurate ground truth labels [1].

There are several types of uncertainty estimation approaches. In this Section we provide a brief overview following the categorization of Gawlikowski et al. [1]. **Single deterministic methods** predict with a deterministic model and provide an uncertainty estimate directly through that model [11], or via external computation [12, 13]. There are various ways for external methods to estimate the uncertainty, e.g., by computing the distance to the training data. DUQ, which is the work of Van Amersfoort et al. [12], belongs to this subcategory. The authors used Radial Basis Functions (RBF) to learn a transformation of features into a hyperspace. At inference time, they quantify the uncertainty as the distance in the hyperspace between the predicted features and a set of learned class centroids.

Another group of works is that of **Bayesian methods**, which infer the probability distribution of the model parameters [1]. In this category, Monte-Carlo Dropout (MC Dropout) [9] approximates this distribution by considering the set of models originating from keeping dropout active at test time. Thanks to the randomness of dropout, each forward pass produces different outputs, from which the uncertainty can be estimated according to their disagreement.

Then, **ensemble methods** also combine multiple outputs for both predictions and uncertainty estimates [14]. In this case the outputs are inferred by a set of different deterministic

models, trained independently. Subcategories of ensembles can be found depending on how these different models are obtained [1].

Lastly, **test-time augmentation methods** use a single deterministic model and augment its input during inference, collecting a variety of predictions, from which a final output and uncertainty estimate are computed [15].

A different categorization distinguishes sampling-based and sampling-free approaches, depending on how the uncertainty is estimated. The former group requires inferring multiple times on the same input with different models (for Bayesian and ensemble methods), or different inputs with the same model (with test-time augmentations), and then aggregating the results. This significantly increases the runtime [9, 10], which makes it impractical for most real-time scenarios, such as autonomous driving. Instead, the latter category does not require sampling, thanks to built-in quantification techniques, such as those typically used in single deterministic methods [12].

Our work builds on top of the findings of Van Amersfoort et al. [12]. In particular, we considerably extend and adapt their DUQ from image classification to a dense regression task for object detection, estimating the model objectness uncertainty at a given location, while significantly improving training stability and convergence. Additionally, we estimate the uncertainty on the variety of regressed outputs of object detection (e.g., object size), by aggregating the predictions provided by the detector.

B. Uncertainty Estimation in Object Detection

While many approaches quantify the uncertainty on basic tasks, posing the foundations of this domain [9, 12], fewer estimated it for safety-critical settings, such as autonomous driving [4, 16]. As we focus on estimating the uncertainty for object detection, in this Section we provide a brief summary of existing works in this area.

Choi et al. proposed Gaussian YOLOv3 (GYOLO) [3], extending the popular YOLOv3 detector [17] to estimate the object location and size uncertainties. They achieved this by predicting Gaussian parameters (i.e., mean and variance) for the position and dimensions of the boxes. Miller et al. [10] used MC Dropout [9] and thoroughly evaluated different clustering techniques, extracting variance and entropy of spatial and class predictions respectively. Harakeh et al. [16] also used MC Dropout [9] and replaced the usual non-maximum suppression (NMS) with clustering and Bayesian inference, obtaining an uncertainty quantification for the object size and class. Instead, Lee et al. [7] introduced with their Gaussian-FCOS a dedicated head on top of an anchor-less detector to estimate the localization uncertainty for each of the four box boundaries independently.

Likewise, our proposed CertainNet estimates the uncertainty for object detection, but in a different way. Unlike the works of Miller et al. [10] and Harakeh et al. [16], ours is sampling-free and real-time capable. Compared to GYOLO [3] and Gaussian-FCOS [7], ours does not regress the uncertainty as an explicit additional output, thereby reducing the potential

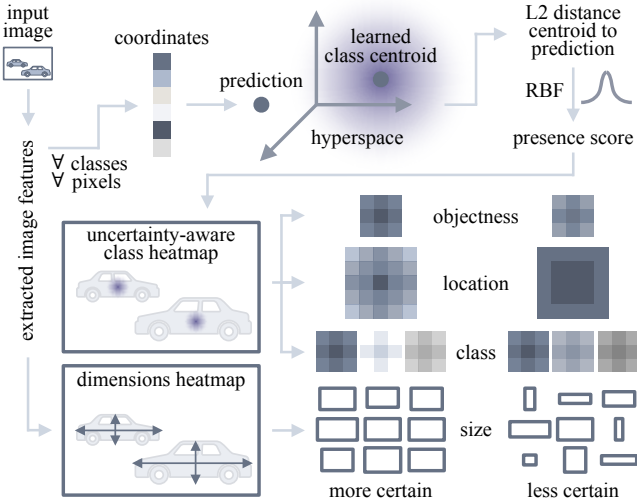


Fig. 2. The proposed CertainNet. The top shows how each uncertainty-aware objectness score is computed. At the bottom right, examples of more and less certain predictions for each output signal are shown.

impact of biases in the training data. Moreover, we are the first to quantify the uncertainty for each individual aspect of the detection output (i.e., objectness, location, size and class), increasing the explainability.

III. METHOD

In this work, we estimate the uncertainty for each aspect of object detection, improving the model explainability for safety-critical applications. Towards this end, we build on top of an anchor-less detection framework (Section III-A), which we render uncertainty-aware (Section III-B). We then compute the uncertainty for objectness, location, dimensions and class (Section III-C). Figure 2 shows an overview of the proposed method. This fragmentation of the uncertainty along the individual detection components introduces valuable information to improve the safety of downstream tasks.

A. Base Detection Framework

CenterNet [18] is a fast and accurate anchor-less detector. It predicts objects centers via a class heatmap, and regresses the bounding box size in a separate dimensions heatmap. CenterNet gives as intermediate output a bounding box for every pixel, hence providing a distribution of outputs to estimate our uncertainties. For these reasons, we use it as base detector. While we focus on 2D bounding boxes, 3D boxes can be predicted following [18], and the additional uncertainties (e.g., for object distance and length) can be computed analogously to the ones presented in this work.

B. Uncertainty-awareness

We make the class heatmap of the detector uncertainty-aware, by exploiting DUQ [12], extending it from recognizing out-of-domain samples in image classification tasks, for which it was proposed, to a dense regression task within our object detector. During training, we learn a set of class representatives (i.e., centroids), which are then compared with each prediction

at inference time. By doing so, similarly to DUQ [12], we assess the deviation from learned object prototypes, which is linked to the model uncertainty. We compute this for all regressed heatmap values, which represent the objectness scores. In particular, it is the comparison with the learned centroids, which makes every value of this heatmap uncertainty-aware.

Proposed for image classification, DUQ [12] learns to map each image to N_{cls} high dimensional spaces, one for each of N_{cls} classes. At training time, class centroids are updated as moving averages of positive samples. As Berger et al. [19] also pointed out, DUQ suffers from instability issues and deployment difficulties. Therefore, after performing a thorough analysis of DUQ and its feature space, we incorporate several modifications to better suit the object detection task, as well as to improve the overall convergence and stability of the model, as visualized in Figure 3.

Adaptation of DUQ to dense regression: Compared to the original implementation [12], due to the high dimensionality of the centroid space (e.g., 512 in [12]), and also the two additional dimensions required for our class heatmap output (i.e., width and height), speed and memory efficiency issues arise. We circumvent these by exploiting the shared operation across the pixels when transforming to the hyperspace. Towards this end, we use 1×1 convolutions, instead of individual multiplications as in [12].

Balanced centroid update: The centroids need to be a good representation of the data, as at inference time they determine both predictions and uncertainties, making the centroid update a critical step. In particular, the centroids are updated as trailing averages, of which the current estimates are computed to get closer to the current sets of predictions. In a regular setup [12], scaling by the amount of samples of a given class c is reasonable, as each minibatch t is representative of the whole dataset. However, with unbalanced classes and a varying amount of objects, that would lead to inconsistent update magnitudes, hindering convergence. We circumvent this and account for the minibatch variability by computing the centroid $\mathbf{e}_{c,t}$ as:

$$\mathbf{e}_{c,t} = \gamma \cdot \mathbf{e}_{c,t-1} + \frac{(1-\gamma)}{\sum_i \mathbf{y}_{c,t,i}^\lambda} \sum_i (\mathbf{y}_{c,t,i}^\lambda \circ \mathbf{W}_c f_\theta(\mathbf{x}_{t,i})) \quad (1)$$

This allows to set a weight on the center pixel higher than its surroundings, depending on λ . The right term is an average weighted and scaled by the ground truth class heatmap values $\mathbf{y}_{c,t}$, having a Gaussian around each object center. $\mathbf{W}_c f_\theta(\mathbf{x}_{t,i})$ are the predicted hyperspace coordinates, with f_θ being the feature extractor, and \mathbf{W}_c the hyperspace transformation for class c . γ is a hyperparameter, defined as the centroid momentum, as in [12].

Hyperspace regularization: To properly represent the training data, the centroids should lie among positive samples distributed along a Gaussian hypersphere. However, by visualizing the hyperspace via PCA and t-SNE, we noticed that this is not the case, with the centroid often falling around the boundary of an irregularly shaped training distribution (top of Figure 3). We circumvent this and aim at the ideal case, by

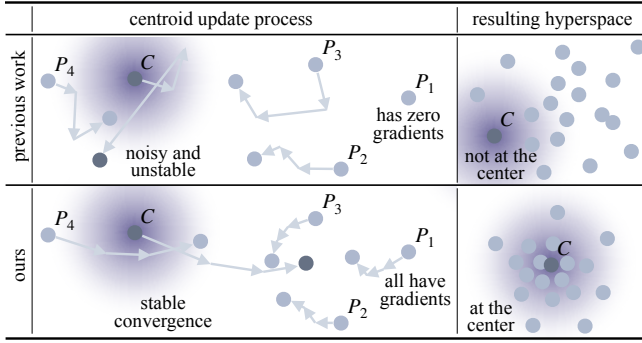


Fig. 3. Centroid update scheme, with the learned hyperspace. The bottom shows the impact of our modifications, compared to the unstable procedure of DUQ [12] at the top when applied on an unbalanced detection task.

regularizing this hyperspace with:

$$\mathcal{L}_{reg} = \frac{1}{\sum_{c,i} \mathbf{y}_{c,i}^\lambda} \sum_{c,i} (\mathbf{y}_{c,i}^\lambda \circ \|\mathbf{W}_{cf\theta}(\mathbf{x}_i) - \mathbf{e}_c\|_2^2) \quad (2)$$

which acts on the euclidean distance between the centroid and prediction for each pixel, and is weighted by the ground truth class heatmap values $\mathbf{y}_{c,i}$. Furthermore, we do not use the gradient penalty \mathcal{L}_{grad} of DUQ [12], as we find \mathcal{L}_{reg} to deliver more stable and consistent results.

Centroid momentum scheduling: The centroids are moved using steps with a magnitude proportional to their distance to the predictions, and to the centroid momentum γ . A large γ slows down training, while improving its stability, but reduces the robustness against the initialization. Instead, a low γ has the inverse effect. While a relatively low γ worked well for balanced classes and large batch sizes [12], it did not in our task. We apply a scheduling of γ similar to that of learning rates, to improve the stability in later training stages, and preserve robustness and speed at earlier ones. This effect is represented in Figure 3. Specifically, we reduce γ by a factor of 10 at fixed epochs.

Outliers protection: Additionally, we prevent outliers from impacting the centroids. In particular, we define an outlier as a prediction being further from its centroid than the triple of the length scale σ . This further improves the training stability and convergence, as instead of having the predictions following a highly moving target (i.e., centroid), they go towards a more stable target.

Length scale annealing: At early stages of training, the centroids position and the hyperspace transformation are not yet properly tuned, due to random initialization. This, together with the high dimensionality of the space, causes instability, as even the correct inputs can be mapped far from their centroid. We alleviate this problem and improve stability with an approach similar to that of simulated annealing [20]: we increase the length scale σ at the initial stages, to then slowly reduce it at each step towards the original value. This increases the gradients at larger distances, nearly zero otherwise, thereby allowing to improve convergence and predictions, as depicted in Figure 3.

C. Uncertainties Estimation

Objectness: Thanks to the class heatmap uncertainty-awareness (Section III-B), the objectness uncertainty \mathcal{U}_{obj} can be directly inferred as the complement of the score:

$$\mathcal{U}_{obj} = 1 - \hat{p} \quad \text{with: } \hat{p} = \max_c p(\mathbf{W}_{cf\theta}(\mathbf{x}) | \mathcal{N}(\mathbf{e}_c, \sigma)) \quad (3)$$

with $\mathcal{N}(\mathbf{e}_c, \sigma)$ being a Gaussian distribution centered at the centroid \mathbf{e}_c with variance σ , which is the length scale.

Location: An object location is computed according to the class heatmap. In our CertainNet, this output is uncertainty-aware (Section III-B). We make use of this by extracting the location uncertainty considering the distribution of scores around an object center. As can be seen in Figure 2, a location is more certain if the center peak is more defined and sharper, while it is less certain with a smoother and flatter peak. The uncertainty \mathcal{U}_x for the location in x direction depends on the horizontal variance and is normalized by the predicted width \hat{w} . We compute \mathcal{U}_x as follows:

$$\mathcal{U}_x = \frac{\sqrt{\text{var}_x}}{\hat{w}} \quad \text{with: } \text{var}_x = \frac{\sum_i o_x^2(i) \hat{p}_i \cos^\eta \alpha_i}{\sum_i \hat{p}_i \cos^\eta \alpha_i} \quad (4)$$

where i iterates over the pixels around the center, $o_x(i)$ is the x -offset from pixel i to the center, α_i is the angle between the x -axis and the pixel, and \hat{p}_i is the predicted objectness (Equation 3). η weights the influence of the scores that are not axis-aligned. The vertical location uncertainty \mathcal{U}_y is computed analogously, replacing \hat{w} with the predicted height \hat{h} and using var_y , obtained through the y -offset $o_y(i)$ and $\sin \alpha_i$.

Dimensions: As described in Section III-A, the model provides a bounding box for each pixel in the output heatmaps. We exploit this by computing the uncertainty on the object size $\mathcal{U}_{w,h}$, by comparing the boxes around an object center. Intuitively, the more these boxes are similar to one another, the more certain the model is on the object dimensions, and vice versa, as shown in Figure 2. Towards this end, we compute \mathcal{U}_w as an RMSE over the surrounding predicted dimensions for each object center, and weight this by the pixel objectness score:

$$\mathcal{U}_w = \frac{1}{\hat{w}} \sqrt{\frac{1}{\sum_i \hat{p}_i} \sum_i \hat{p}_i \cdot (\hat{w} - \hat{w}_i)^2} \quad (5)$$

where \hat{w} is the predicted object width and \hat{w}_i is the width at the heatmap pixel i around the center. The height uncertainty \mathcal{U}_h is computed analogously using the predicted heights \hat{h} .

Class: The object class is predicted via the class heatmap, which is uncertainty-aware in our model (Section III-B). As in Figure 2, a class prediction is more certain, if there is a peak only at such class with the other class probabilities being low, while it is less certain, the higher those other probabilities are. We formalize this, by comparing the objectness scores at each predicted center across the classes:

$$\mathcal{U}_{cls}(i) = \mathcal{U}_{cls}(i-1) + (1 - \mathcal{U}_{cls}(i-1)) \cdot \left(\frac{\hat{p}_{cls}(i)}{\hat{p}_{cls}(1)} \right)^i \quad (6)$$

with $\hat{p}_{cls}(i)$ being the i^{th} -highest class probability, so $i = 1$ for the predicted one. We compute \mathcal{U}_{cls} recursively, considering $\mathcal{U}_{cls}(0) = 0$ and reaching the final class uncertainty \mathcal{U}_{cls} as $\mathcal{U}_{cls}(N_{cls})$.

TABLE I
DETECTION PERFORMANCE COMPARISON WITH RELATED WORKS ON THE KITTI [6] VALIDATION SET ACROSS CLASSES AND DIFFICULTIES.

Method	mAP	Car			Pedestrian			Cyclist			all		fps
		Easy	Mod.	Hard	Easy	Mod.	Hard	Easy	Mod.	Hard	AUPR-In	AUPR-Out	
GaussianYOLO [3]	56.74	89.91	84.36	75.65	58.43	50.95	43.69	45.97	31.17	30.53	65.39	96.42	18.06
CenterNet [18]	71.49	92.30	89.15	82.17	76.53	67.53	59.37	73.48	52.63	50.25	68.33	99.82	20.21
5-Ensemble [10]	73.70	97.84	89.75	80.87	77.08	67.66	59.17	79.45	57.30	54.17	72.13	97.32	5.10
CertainNet [ours]	73.00	93.81	89.36	82.11	76.33	66.13	58.54	78.02	57.49	55.20	69.85	99.85	16.46

TABLE II

UNCERTAINTY QUALITY COMPARISON WITH RELATED WORKS ON THE VALIDATION SETS OF KITTI [6], BDD100K [21], AND nuIMAGES (nuIm.) [22]. ALL MODELS WERE TRAINED ON KITTI, AND TRANSFERRED (\rightarrow) TO THE OTHER DATASETS. ALL RESULTS ARE FOR *Car* (MODERATE FOR KITTI). ON CenterNet[‡] WE APPLIED OUR LOCATION AND DIMENSIONS UNCERTAINTY ESTIMATIONS AS POST-PROCESSING.

	Method	AP	objectness				location			dimensions			
			AUPR-In	AUPR-Out	AUROC	ECE \downarrow	UE \downarrow	CE \downarrow	UBQ	BR	CE \downarrow	UBQ	BR
KITTI	GaussianYOLO [3]	84.36	74.93	91.24	86.23	23.37	20.02	11.58	60.01	60.56	6.18	75.34	77.47
	CenterNet [‡] [18]	89.15	72.76	99.56	96.65	10.49	6.00	7.17	68.27	69.80	4.23	85.84	89.58
	5-Ensemble [10]	89.75	73.79	93.61	87.13	16.93	17.72	4.59	86.71	90.43	5.01	84.28	87.54
	CertainNet [ours]	89.36	78.16	99.81	98.00	4.60	4.98	4.54	75.92	77.81	4.47	86.76	91.05
\rightarrow BDD	GaussianYOLO [3]	31.03	88.73	79.72	85.84	22.80	20.68	9.48	57.32	58.49	7.67	57.94	59.19
	CenterNet [‡] [18]	34.51	81.35	98.06	91.74	7.78	14.03	25.72	35.90	36.76	10.53	73.43	79.55
	5-Ensemble [10]	26.75	96.42	81.84	91.99	25.64	15.26	3.31	84.83	93.54	9.40	84.81	93.50
	CertainNet [ours]	40.93	78.77	97.82	91.17	4.89	14.85	14.79	44.71	46.03	9.68	72.91	78.89
\rightarrow nuIm.	GaussianYOLO [3]	30.38	86.24	87.95	87.78	17.05	17.74	10.42	56.20	57.62	8.74	58.38	60.14
	CenterNet [‡] [18]	43.93	75.94	98.93	91.58	5.51	15.70	27.07	35.69	36.21	8.39	74.79	79.21
	5-Ensemble [10]	35.31	93.93	90.90	93.04	17.74	13.98	3.37	89.48	98.51	8.36	89.31	98.21
	CertainNet [ours]	53.14	79.97	99.28	94.24	7.80	11.90	16.21	43.57	44.46	7.73	75.36	80.25

IV. EXPERIMENTS AND RESULTS

A. Experimental setup

Datasets: We conducted our experiments on three public autonomous driving datasets, namely KITTI [6], the recent nuImages extension of nuScenes [22], and the dashcam-based BDD100K [21]. **KITTI** is a popular benchmark for object detection, and it has been recorded in Germany. We followed the standard 3DOP split [18], which comprises 3712 training and 3769 validation images. Moreover, we used the standard *car*, *pedestrian* and *cyclist* classes. **NuImages** is a highly diverse large-scale dataset, collected in Boston and Singapore. It includes rain, snow and night conditions. **BDD100K** is another diverse large-scale dataset. It was crowd-sourced from a variety of dashboard cameras and different vehicles, around cities of the USA. It also includes various challenging weather and lighting conditions, such as snow, fog and night. In particular, we used nuImages and BDD100K to assess the generalization capability to out-of-domain data, which is critical for autonomous driving and uncertainty estimation. Towards this end, we evaluated KITTI models (without any fine-tuning) on the validation set of nuImages, which includes 3249 images, sized 1600x900, and the validation set of BDD100K, which contains 10K samples, sized 1280x720. Due to the sub-optimal class overlap between KITTI and the two transfer datasets, for the evaluation we focused on the *car* class.

Evaluation metrics: We evaluated the object detection performance on the standard AP, with the standard IoU threshold of 0.7 for KITTI *cars* and 0.5 for the other classes and datasets. We split the evaluation according to the specific

uncertainty estimates, which we thoroughly evaluated. For the objectness uncertainty we computed various metrics, such as the area under the precision-recall (AUPR) curve, as AUPR-In and AUPR-Out [10], which should be considered jointly. The former assesses the ability to accept all correct detections and minimise wrong acceptances. Conversely, the latter looks at the rejection of negative samples and the rate of wrong rejections. Additionally, we computed the area under the receiver operating characteristic curve (AUROC), comparing true and false positives rates, as used by Miller et al. [10]. Following [23], we calculated the expected calibration error (ECE) as the average difference between confidence and detection accuracy, weighted by the occurrence. Moreover, we evaluated the minimum uncertainty error (UE) as the ability to accept correct detections and reject incorrect ones based on the uncertainty estimate [10]. In the top of Figure 4, we show how location and size uncertainties have a similar effect on the box uncertainty boundaries, despite leading to different box distributions. This allows to evaluate them in the same way. We adapted the expected calibration error (ECE) [23] for our detection case, as the calibration error (CE): average deviation between predicted uncertainty and measured error for matched detections. We independently computed the CE on x and y for the location, and on *width* and *height* for the size, reporting the average of each pair. All metrics in this work are percentage-based.

Uncertainty boundary quality: We introduce a new metric to evaluate the uncertainties on the object location and size. As shown at the bottom of Figure 4, the goal is assessing how

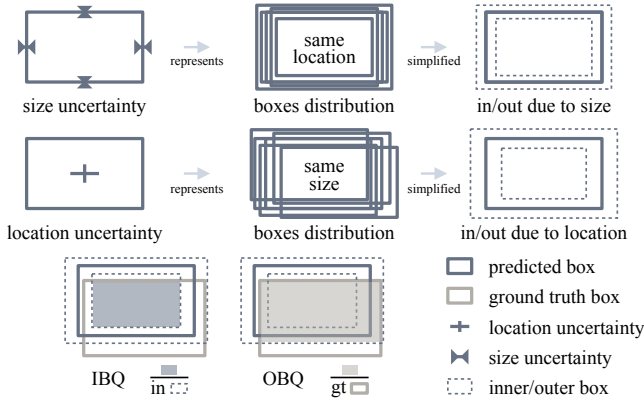


Fig. 4. The top shows the effect of size and location uncertainties on the box uncertainty boundaries, allowing to evaluate them in a similar fashion. The bottom shows the IBQ and OBQ terms of our novel UBQ metric.

much of an object is in the uncertainty boundaries. We name this metric uncertainty boundary quality (UBQ), and compute it as the average across matched boxes of:

$$UBQ_i = \frac{1}{2} (IBQ_i + OBQ_i) BR_i \quad (7)$$

with BR being the boundary ratio between inner and outer:

$$BR_i = \frac{1}{2} \left(\frac{\text{width}(in_i)}{\text{width}(out_i)} + \frac{\text{height}(in_i)}{\text{height}(out_i)} \right) \quad (8)$$

UBQ is composed of an inner (IBQ) and an outer (OBQ) term, which are averaged and multiplied by the ratio (BR), to penalize extremely broad and trivial confidence intervals. Figure 4 shows how IBQ_i and OBQ_i are computed:

$$IBQ_i = \frac{A(gt_i \cap in_i)}{A(in_i)} \quad OBQ_i = \frac{A(gt_i \cap out_i)}{A(gt_i)} \quad (9)$$

where A is the area of a bounding box i , gt is the ground truth box, in and out are the inner and outer predictions. As for the CE, our UBQ metric can be used for both location and size. We report UBQ as a summary, and BR as the ratio.

Network architecture: All our models were based on a DLA architecture [24]. We selected it as compared to others it offers a good speed-accuracy trade-off [18]. Thus, we followed the DLA-34 configuration of CenterNet, with deformable convolutions [18]. We then made the class heatmap head uncertainty-aware, with the hyperspace transformation and the RBF kernel, as described in Section III-B. We used a single model for all three classes of KITTI.

Implementation details: We trained our models with an Adam optimizer for 70 epochs, with 10 additional epochs for the objectness head. We used an initial learning rate of $1.25e-4$, reduced by a factor of 10 after epochs 45 and 60. We chose a batch size of 16 and a resolution of 1280x384 for KITTI, resizing BDD and nuImages samples to 896x512. We selected 512 as centroid dimensionality, and initialized their momentum γ to 0.9, increasing it to (0.99, 0.999, 0.9999) at epochs (5, 20, 60) respectively. The length scale σ was initialized at $2.5e-1$, with a decay rate of 0.999 per step for the sigma annealing, until reaching $5e-2$. We weighted the hyperspace regularization loss \mathcal{L}_{reg} by $1e-2$, and set λ to 20

TABLE III
ABLATION STUDY ON THE KITTI [6] VALIDATION SET ACCORDING TO *Car* MODERATE DETECTION AND UNCERTAINTY QUALITY.

ID	Description	AP	AUPR-In	AUPR-Out
A0	adapted DUQ [12]	76.41	81.64	99.62
A1	A0 - $\mathcal{L}_{grad} + \mathcal{L}_{reg}$	87.33	76.26	99.15
A2	A1 + balanced update	87.98	74.12	99.71
A3	A2 + outlier protection	88.14	77.34	99.56
A4	A3 + momentum schedul.	88.76	75.27	99.72
A5	A4 + length scale anneal.	88.85	81.08	99.80
A6	A5 + freeze last 10 epochs	89.36	78.16	99.81

for balancing the score. To compensate the added loss on the objectness, the weight on the dimensions loss was increased to 0.2. We set $\eta = 4$, and kept the remaining hyperparameters the same as proposed for CenterNet [18]. We trained our models using PyTorch on a single NVIDIA Tesla V100 32GB GPU and evaluated the runtime (fps) on a single NVIDIA Quadro RTX 4000 8GB GPU.

Details for prior works: For a fair comparison, we re-trained all methods on the same standard dataset split [18], using the official implementations and hyperparameters, starting from weights pretrained on ImageNet [25], until convergence. For GYOLO [3] we used the best settings with input size 704 on Darknet. When evaluating location and size uncertainties, we extended CenterNet [18] with our location and dimensions uncertainty estimations, on top of its outputs. For the ensemble, we trained 5 DLA-34-based CenterNets (randomly initialized, with data shuffle and augmentations), and combined their predictions with the merging strategy of Miller et al. [10], via DBSCAN. For CenterNet, GYOLO and the ensemble, the objectness uncertainty is evaluated using the standard confidence scores provided.

B. Quantitative Results

Detection: Table I shows the detection performance of our approach compared to related works on KITTI [6]. Our uncertainty-aware CertainNet outperformed the CenterNet [18] baseline on the overall mAP, as well as most classes and difficulties, while sharing the same architecture. Our additional safety-related outputs reduced the runtime by 18% to 16.46 fps, thereby remaining real-time. Moreover, the sampling-based 5-Ensemble [10] models achieved the overall highest mAP, and the best score in most classes and difficulties, with ours securing a close second place. However, the superior performance of the ensemble came at the cost of a significantly higher runtime. GYOLO [3] was slightly faster than ours, but produced the worst detections among tested methods. As shown in Table I, its AP scores are significantly lower than those reported by Choi et al. in [3], due to the different dataset splits used for KITTI. In particular, Choi et al. split training and validation randomly across single frames. We used instead a standard split [18], reducing the similarity between training and validation sets.

Uncertainty estimation: In Table II we report a variety of uncertainty metrics for different datasets, for the main class *Car*. On KITTI, our CertainNet achieved higher scores and

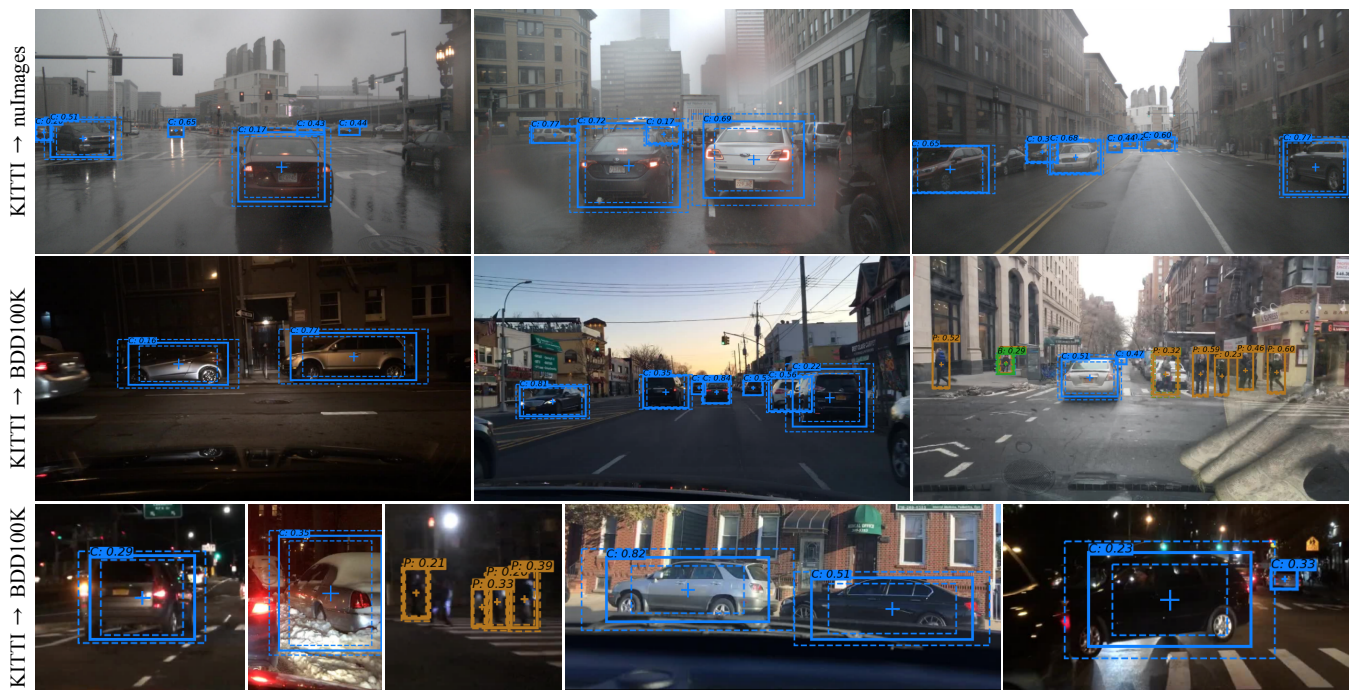


Fig. 5. Example predictions of the proposed CertainNet on out-of-domain data. The model was trained only on KITTI [6] and applied on nuImages [22] and BDD100K [21] without any fine-tuning. The bottom row shows crops. The boxes follow the format introduced in Figure 1, with KITTI \rightarrow KITTI.

lower errors than the other methods across most metrics. Remarkably, our approach was able to achieve the best objectness results across the board, often with a substantial margin (e.g., on the calibration ECE), thanks to its effective uncertainty-aware heatmap. Interestingly, in terms of uncertainties, the ensemble could not achieve satisfactory results, except for the location. This can be attributed to the disagreement of the 5 models, and to the difficulty of merging contrasting predictions in the case of object detection. Table II shows also the suboptimal performance of GYOLO [3], which had a tendency to overestimate the uncertainty boundaries, especially for the location (BR around 60%), despite its training time modifications for location and size, and having been trained on the same dataset. Furthermore, the flexibility of our uncertainty quantification techniques allowed to produce estimates also for the location and size of the uncertainty-unaware CenterNet [18]. However, this configuration could not match the uncertainty quality of our CertainNet, especially for the location, computed on the class heatmap. This shows how our training time modifications contributed to the uncertainty-awareness of the whole model.

Out-of-domain data: As all models were trained on KITTI, the results for BDD100K and nuImages in Table II show the ability of each method to generalize to unseen scenarios in out-of-domain data. As also pointed out in [26, 27], multiple differences define a substantial domain gap between these two datasets and KITTI: weather and lighting conditions, image resolution and aspect ratio, camera mounting position and type (e.g., dashcam often with reflections on BDD), different continents and street layouts, amount of dynamic scenes, as well as different vehicles on sale in those regions. All these render the transfer task particularly difficult. Most notably,

our CertainNet reached the highest AP on both datasets, by a significant margin on both BDD and nuImages. This shows the robustness of our approach, as well as the strong generalization provided by our uncertainty estimations, compared to CenterNet [18]. Interestingly, the ensemble, which performed well on KITTI (i.e., the training dataset), underperformed on both out-of-domain datasets. We attribute this to the rule-based clustering step required to merge the predictions [10], which was tuned on KITTI. This reiterates that the ensemble’s aggregation step is critical on object detection, and is prone to deliver suboptimal results on out-of-domain data, which is fundamental in autonomous driving. GYOLO [3] followed a similar trend, but suffered the transfer to out-of-distribution less than the ensemble, thanks to its learning-based techniques that generalized better. Nevertheless, the uncertainty estimations of the ensemble seemed high quality for the most part, especially on the object location and size. However, this is due to the rather low number of detections that all related works provided, as discussed below. Overall, our CertainNet achieved the highest AP, while maintaining high quality uncertainty estimates, which correctly increased on out-of-domain data.

Amount of detections: The uncertainty evaluation metrics used in this work are highly influenced by the number of predictions or that of matched detections. Therefore, high scores and low error rates can be achieved by a model that outputs a low amount of high confidence boxes. We found this to be the case for the methods in Table II: in nuImages [22] our method detected 4905 *cars*, CenterNet 3714, 5-Ensemble 3229, and GYOLO 2248; similarly in BDD [21], our CertainNet 26802, CenterNet 20210, ensemble 16440, and GYOLO 15395. These values also show the limitations of the ensemble, confirming that contrasting predictions within its models,

especially in highly uncertain scenarios, such as out-of-domain data (e.g., BDD100K and nuImages in Table II), can happen and thereby significantly degrade the overall detection and uncertainty performance of the system. In general, each of the uncertainty metrics evaluates only a specific aspect and should be considered together with others, and combined with the detection performance (i.e., AP and mAP).

Ablation study: The importance of our modifications over DUQ [12] is testified by the results shown in Table III, where A6 represents our full approach. In particular, the baseline A0, which is the adaptation of DUQ [12] from image classification to object detection, could not converge properly, due to the severe limitations described in Section III-B, worsened by the high imbalance of KITTI [6]. The biggest improvement was brought by A1 replacing the gradient penalty \mathcal{L}_{grad} of DUQ with our hyperspace regularization \mathcal{L}_{reg} , which also improved the consistency of training convergence. Moreover, the AP increased monotonously with the introduction of each of our modifications (Section III-B), while maintaining good uncertainty estimates throughout, thereby showing their significant impact over DUQ (i.e., A0).

C. Qualitative Results

In Figures 1 and 5 we show predictions of our method on challenging scenes of KITTI and out-of-domain samples respectively. The images confirm the large domain gap between KITTI and the two transfer datasets, as described in Section IV-B. The night and rainy settings in Figure 5 from nuImages [22] and BDD100K [21] are completely different from the sunny and overcast scenarios of KITTI [6] on which the model was trained, making these very difficult samples. Nevertheless, the proposed CertainNet correctly detected most objects. Compared to the KITTI predictions in Figure 1, the objectness certainties (values at the top left of each box) are lower for the unusual and difficult to see out-of-domain vehicles in Figure 5. This shows that the model recognized them as out-of-distribution samples. Interestingly, the uncertainty estimates for location and size improved the overlap between sub-optimal predictions and the object, e.g., bottom left of Figure 5, as tested by our UHQ metric.

V. CONCLUSION

In this work we introduced CertainNet, a novel method to estimate the uncertainty of every aspect of 2D object detection. Extensive evaluations showed the benefit of quantifying the uncertainty, especially when transferring to out-of-distribution data for safety-critical applications. Our method substantially improved the generalization ability over previous works, with more accurate predictions and high quality uncertainty estimates, while running in real-time. Therefore, the proposed CertainNet constitutes a valuable contribution towards robust object detection for autonomous driving.

REFERENCES

[1] J. Gawlikowski, C. R. N. Tassi, M. Ali, J. Lee, M. Humt, J. Feng, *et al.*, "A survey of uncertainty in deep neural networks," *arXiv*, 2107.03342, 2021.

[2] D. Feng, L. Rosenbaum, and K. Dietmayer, "Towards safe autonomous driving: Capture uncertainty in the deep neural network for LiDAR 3D vehicle detection," in *IEEE ITSC*, 2018, pp. 3266–3273.

[3] J. Choi, D. Chun, H. Kim, and H.-J. Lee, "Gaussian YOLOv3: An accurate and fast object detector using localization uncertainty for autonomous driving," in *IEEE/CVF ICCV*, 2019, pp. 502–511.

[4] J. Postels, F. Ferroni, H. Coskun, N. Navab, and F. Tombari, "Sampling-free epistemic uncertainty estimation using approximated variance propagation," in *IEEE/CVF ICCV*, 2019, pp. 2931–2940.

[5] N. Yang, L. v. Stumberg, R. Wang, and D. Cremers, "D3VO: Deep depth, deep pose and deep uncertainty for monocular visual odometry," in *IEEE/CVF CVPR*, 2020, pp. 1281–1292.

[6] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *IEEE/CVF CVPR*, 2012, pp. 3354–3361.

[7] Y. Lee, J.-w. Hwang, H.-I. Kim, K. Yun, and Y. Kwon, "Localization uncertainty estimation for anchor-free object detection," *arXiv*, 2006.15607, 2020.

[8] T. Cortinhal, G. Tzelepis, and E. E. Aksoy, "SalsaNext: Fast, uncertainty-aware semantic segmentation of LiDAR point clouds," in *ISVC*, 2020, pp. 207–222.

[9] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *ICML*, 2016, pp. 1050–1059.

[10] D. Miller, F. Dayoub, M. Milford, and N. Sünderhauf, "Evaluating merging strategies for sampling-based uncertainty techniques in object detection," in *IEEE ICRA*, 2019, pp. 2348–2354.

[11] M. Sensoy, L. Kaplan, and M. Kandemir, "Evidential deep learning to quantify classification uncertainty," in *NeurIPS*, 2018, pp. 3183–3193.

[12] J. Van Amersfoort, L. Smith, Y. W. Teh, and Y. Gal, "Uncertainty estimation using a single deep deterministic neural network," in *ICML*, 2020, pp. 9690–9700.

[13] M. Raghu, K. Blumer, R. Sayres, Z. Obermeyer, B. Kleinberg, S. Mullainathan, and J. Kleinberg, "Direct uncertainty prediction for medical second opinions," in *ICML*, 2019, pp. 5281–5290.

[14] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," in *NeurIPS*, 2017, pp. 6405–6416.

[15] G. Wang, W. Li, M. Aertsen, J. Deprest, S. Ourselin, and T. Vercauteren, "Aleatoric uncertainty estimation with test-time augmentation for medical image segmentation with convolutional neural networks," *Neurocomputing*, vol. 338, pp. 34–45, 2019.

[16] A. Harakeh, M. Smart, and S. L. Waslander, "BayesOD: A Bayesian approach for uncertainty estimation in deep object detectors," in *IEEE ICRA*, 2020, pp. 87–93.

[17] J. Redmon and A. Farhadi, "YOLOv3: An incremental improvement," *arXiv*, 1804.02767, 2018.

[18] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *arXiv*, 1904.07850, 2019.

[19] C. Berger, M. Paschali, B. Glocker, and K. Kamnitsas, "Confidence-based out-of-distribution detection: A comparative study and analysis," in *UNSURE MICCAI Workshop*, pp. 122–132, 2021.

[20] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, "Optimization by simulated annealing," *Science*, vol. 220, no. 4598, pp. 671–680, 1983.

[21] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell, "BDD100k: A diverse driving dataset for heterogeneous multitask learning," in *IEEE/CVF CVPR*, 2020, pp. 2636–2645.

[22] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liang, Q. Xu, *et al.*, "nusenes: A multimodal dataset for autonomous driving," in *IEEE/CVF CVPR*, 2020, pp. 11 621–11 631.

[23] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, "On calibration of modern neural networks," in *ICML*, 2017, pp. 1321–1330.

[24] F. Yu, D. Wang, E. Shelhamer, and T. Darrell, "Deep layer aggregation," in *IEEE/CVF CVPR*, 2018, pp. 2403–2412.

[25] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, *et al.*, "ImageNet Large Scale Visual Recognition Challenge," *IJCV*, vol. 115, no. 3, pp. 211–252, 2015.

[26] S. Gasperini, P. Koch, V. Dalletta, N. Navab, B. Busam, and F. Tombari, "R4Dyn: Exploring radar for self-supervised monocular depth estimation of dynamic scenes," in *IEEE 3DV*, 2021, pp. 751–760.

[27] V. Guizilini, R. Ambrus, S. Pillai, A. Raventos, and A. Gaidon, "3D packing for self-supervised monocular depth estimation," in *IEEE/CVF CVPR*, 2020, pp. 2485–2494.