

Probabilistic Regression for Visual Tracking

Martin Danelljan Luc Van Gool Radu Timofte

Computer Vision Lab, D-ITET, ETH Zürich, Switzerland

Abstract

Visual tracking is fundamentally the problem of regressing the state of the target in each video frame. While significant progress has been achieved, trackers are still prone to failures and inaccuracies. It is therefore crucial to represent the uncertainty in the target estimation. Although current prominent paradigms rely on estimating a state-dependent confidence score, this value lacks a clear probabilistic interpretation, complicating its use.

In this work, we therefore propose a probabilistic regression formulation and apply it to tracking. Our network predicts the conditional probability density of the target state given an input image. Crucially, our formulation is capable of modeling label noise stemming from inaccurate annotations and ambiguities in the task. The regression network is trained by minimizing the Kullback-Leibler divergence. When applied for tracking, our formulation not only allows a probabilistic representation of the output, but also substantially improves the performance. Our tracker sets a new state-of-the-art on six datasets, achieving 59.8% AUC on LaSOT and 75.8% Success on TrackingNet. The code and models are available at <https://github.com/visionml/pytracking>.

1. Introduction

Visual object tracking is the task of estimating the state of a target object in each frame of a video sequence. Most commonly, the state is represented as a bounding box encapsulating the target. Different flavors of the problem arise from the type of given prior information about the scenario, such as object class [1] or static camera [38]. In its most general form, however, virtually no prior knowledge is assumed and the initial state of the target is given during inference. This imposes severe challenges, since the method must learn a model of the target during tracking itself.

Along with a myriad of other computer vision tasks, such as object detection [24, 26, 35], pose estimation [6, 42, 48], and keypoint detection [36, 46], visual tracking can fundamentally be formulated as a regression problem. In this general view, the goal is thus to learn a model, typically a deep

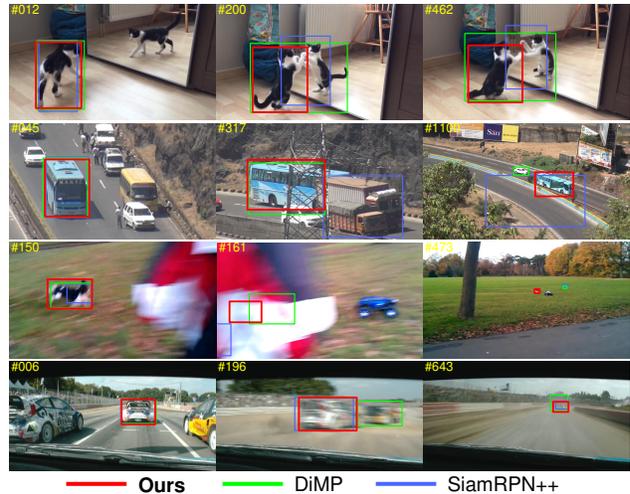


Figure 1. A comparison of our approach with state-of-the-art trackers DiMP [3] and SiamRPN++ [28]. In tracking, estimating the uncertainty of the target state is important in the presence of similar objects (first row), occlusions (second row), when determining failures (third row), and in cases of blur or other obstructions (bottom row). Unlike current state-of-the-art, our approach predicts a probability distribution $p(y|x)$ of the target state y conditioned on the input image x , providing a clear interpretation of the output. The proposed probabilistic formulation further improves the overall performance of the tracker, including the cases shown above.

neural network, capable of predicting the target state in each frame. While current and past approaches apply a variety of techniques to address this problem, most of the successful methods have a crucial aspect in common. Namely, that the task of regressing the target state y^* in a frame x is achieved by learning to predict a confidence value $s(y, x)$ for any given state y . The target state is then estimated by maximizing the predicted confidence $y^* = \arg \max_y s(y, x)$.

The aforementioned confidence-based regression strategy is shared by the previously dominant Discriminative Correlation Filter (DCF) paradigm [5, 10, 12, 15, 22, 31, 40] and the more recent Siamese trackers [2, 17, 28, 29, 45, 51]. Both employ a convolutional operation to predict a target confidence $s(y, x)$ at each spatial position y , in order to localize the target. Recent work [3, 9] also demonstrated the effectiveness of training a network branch to predict

the confidence $s(y, x)$ of the entire target box y to achieve highly accurate bounding box regression. Due to the vast success of these confidence-based regression techniques, we first set out to unify much of the recent progress in visual tracking under this general view.

One distinctive advantage of confidence-based regression is its ability to flexibly represent *uncertainties*, encoded in the predicted confidence values $s(y, x)$. In contrast, a direct regression strategy $y = f(x)$ forces the network to commit to a single prediction y , providing no other information. However, the confidence value $s(y, x)$ itself has no clear interpretation, since it simply acts as a quantity to be maximized. The range of values and characteristic of the predicted confidence largely depends on the choice of loss and strategy for generating the corresponding pseudo labels for training. This provides severe challenges when designing strategies for estimating and reasoning with the uncertainty in the prediction. Such measures are highly relevant in tracking, *e.g.* to decide whether to update, if the target is lost, or how uncertain the output is (see Figure 1). We aim to address these limitations by taking a probabilistic view.

Contributions: We propose a formulation for learning to predict the conditional probability density $p(y|x)$ of the target state y given an input image x . Unlike the confidence value $s(y, x)$, the density $p(y|x)$ has a clear and direct interpretation, allowing the computation of absolute probabilities. We assume no particular family of distributions, such as Gaussian, instead letting $p(y|x)$ be directly parametrized by the network architecture itself. Specifically, the density $p(y|x)$ is represented by the continuous generalization of the SoftMax operation, previously employed in energy-based models [27] and recently in [18]. In contrast to these previous works, we also model the uncertainty in the annotations themselves. This is shown to be crucial in visual tracking to counter noise in the annotations and ambiguities in the regression task itself. The network is trained by minimizing the Kullback-Leibler divergence between the predicted density and the label distribution.

We demonstrate the effectiveness of our general approach by integrating it into the recent state-of-the-art tracker DiMP [3]. Our resulting tracker does not only allow a fully probabilistic representation $p(y|x)$ of the predicted target state. Comprehensive experiments on *seven* benchmark datasets show that our probabilistic representation and training significantly *improves* the performance of the tracker. Our Probabilistic DiMP (PrDiMP) outperforms previous state-of-the-art by a large margin, particularly on available large-scale datasets, including LaSOT (+2.9% AUC) and TrackingNet (+1.8% Success).

2. Regression by Confidence Prediction

In machine learning, regression is fundamentally the problem of learning a mapping $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ from an

input space \mathcal{X} to a *continuous* output space \mathcal{Y} , given a set of example pairs $\{(x_i, y_i)\}_i \subset \mathcal{X} \times \mathcal{Y}$. For our purposes, \mathcal{X} constitutes the space of images. The most straightforward take on regression is to directly learn the function f_θ , parametrized as *e.g.* a deep neural network with weights θ , by minimizing a loss $L(\theta) = \sum_i \ell(f_\theta(x_i), y_i)$. Here, the function ℓ measures the discrepancy between the prediction $f_\theta(x_i)$ and corresponding ground-truth value y_i . While the choice of loss ℓ is highly problem dependent, popular alternatives include the L^p family, $\ell(y, y') = \|y - y'\|_p^p$.

2.1. General Formulation

While direct regression is successfully applied for many computer vision problems, including optical flow [41] and depth estimation [16], it has proven less suitable for other vision tasks. Examples of the latter include visual tracking [2, 11, 22], object detection [24, 26, 35] and human pose estimation [6, 42, 48]. In these problems, networks are often trained to instead predict a confidence score, which is then maximized in order to achieve the final estimate. Confidence prediction has prevailed over standard direct regression in these circumstances thanks to two key advantages. First, confidence prediction can capture the presence of uncertainties, multiple hypotheses and ambiguities in the output space \mathcal{Y} . The network does not have to commit to a single estimate $f_\theta(x) = y$. Second, the network can more easily exploit symmetries shared by \mathcal{X} and \mathcal{Y} , such as translational invariance in the case of image 2D-coordinate regression tasks, which are particularly suitable for CNNs.

Formally, we define confidence-based regression as learning a function $s_\theta : \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{R}$ that predicts a scalar confidence score $s_\theta(y, x) \in \mathbb{R}$ given an output-input pair (y, x) . The final estimate $f(x) = y^*$ is obtained by maximizing the confidence w.r.t. to y ,

$$f(x) = \arg \max_{y \in \mathcal{Y}} s_\theta(y, x). \quad (1)$$

The regression problem is thus transformed to learning the function s_θ from the data $\{(x_i, y_i)\}_i$. This is generally performed by defining a function $a : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ for generating a pseudo label $a(y, y_i)$, acting as the ground-truth confidence value for the prediction $s_\theta(y, x_i)$. The confidence prediction network can then be trained by minimizing the loss $L = \sum_i L(\theta; x_i, y_i)$, where

$$L(\theta; x_i, y_i) = \int_{\mathcal{Y}} \ell(s_\theta(y, x_i), a(y, y_i)) dy. \quad (2)$$

The function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ now instead measures the discrepancy between the predicted confidence value $s_\theta(y, x_i)$ and the corresponding label value $a(y, y_i)$. In practice, a variety of losses ℓ and pseudo label functions a are employed, depending on the task at hand. In the next section, some more popular examples are studied, where our discussion focuses on the visual tracking problem in particular.

2.2. In Visual Tracking

In visual tracking, the task is to regress the state of the target object in each frame of the video, given its initial location. The state is most often represented as an axis-aligned bounding box $y \in \mathbb{R}^4$. Compared to other vision tasks, this problem is particularly challenging since an example appearance of the target object is only provided at test-time. The tracker must thus learn a model based on this first example in order to locate the object in each frame.

Due to this challenging nature of the problem, the majority of approaches until very recently, focused on regressing the center 2D image coordinate $y \in \mathbb{R}^2$ of the target object, and then optionally using this model to estimate the one-parameter scale factor by a multi-scale search. This class of methods include the widely popular Discriminative Correlation Filter (DCF) approaches [5, 10, 12, 22], most of the more recent Siamese networks [2, 17, 45, 51] and other earlier approaches [50]. The formulation (1), (2) is explicitly utilized in the theory of Structural SVMs [43], employed in the well-known Struck tracker [19]. In DCF-based methods, a convolutional layer is trained *online*, *i.e.* during tracking, to predict a target confidence score

$$s_\theta(y, x) = (w_\theta * \phi(x))(y). \quad (3)$$

Here, w_θ is the convolution kernel and $\phi(x)$ are the features extracted from the image x , typically by a CNN with frozen weights [12, 31]. The result of the convolution (3) is evaluated at the spatial location y to obtain the confidence $s_\theta(y, x)$. The DCF paradigm adopts a squared loss $\ell(s, a) = (s - a)^2$ on the confidence predictions, which enables efficient optimization of (2) w.r.t. w_θ in the Fourier domain [5, 12]. Nearly all DCF methods employ a Gaussian confidence pseudo label $a(y, y_i) = e^{-\frac{\|y - y_i\|^2}{2\sigma^2}}$ centered at the target position y_i in frame x_i .

In contrast to DCF, Siamese trackers [2, 17, 28, 29, 45, 51] aim to fully learn the parameters θ of the network in an offline training stage. This is performed by learning an embedding space ϕ_θ in which similarities between a target template z and frame x can be computed as a correlation,

$$s_\theta(y, x) = (\phi_\theta(z) * \phi_\theta(x))(y). \quad (4)$$

Siamese methods often employ a binary cross entropy loss

$$\ell(s, a) = a \log(1 + e^{-s}) + (1 - a) \log(1 + e^s) \quad (5)$$

in (2) to train the network parameters θ . That is, target localization is treated as a dense binary classification problem, where the pseudo label $a(y, y_i) \in [0, 1]$ represents the target/background class, or more generally, a Bernoulli distribution. It is commonly set to $a(y, y_i) = 1$ in the target vicinity $\|y - y_i\| < r$ and $a(y, y_i) = 0$ otherwise [2].

To achieve an accurate prediction of the full target bounding box, a few recent trackers [3, 9, 28, 29] have

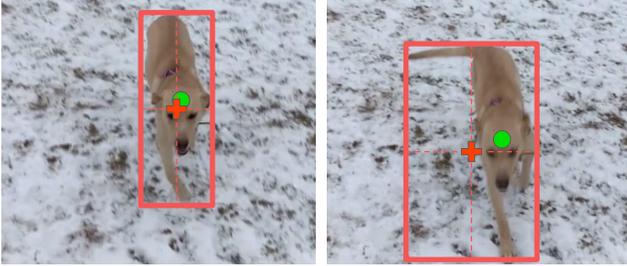
achieved remarkable performance by separating the tracking problem into two parts. First, the object is coarsely localized using techniques reminiscent of the aforementioned approaches, that are robust to similar background objects, clutter and occlusions. In the second stage, a separate network branch is employed for regressing the target bounding box. For this purpose, the ATOM tracker [9] employs an IoU-Net [24] based network head $s_\theta(y, x)$, which scores any input bounding-box $y \in \mathbb{R}^4$. This head is trained in an offline learning stage to predict the Intersection-over-Union (IoU) overlap $a(y, y_i) = \text{IoU}(y, y_i)$ using the squared error $\ell(s, a) = (s - a)^2$ in (2). In this case, the integral (2) is approximated by sampling bounding boxes during training. During tracking, the optimal box (1) is achieved by gradient-based maximization of the predicted confidence.

More recently, Bhat *et al.* [3] proposed the DiMP tracker by designing a meta-learning based network architecture which predicts discriminative target model weights $w_\theta = \psi_\theta(\{(\phi_\theta(z_j), y_j)\}_j)$ in (3) from a set of sample pairs $\{(z_j, y_j)\}_j$. The predicted weights are then employed for the first-stage robust target localization, and updated during tracking through a learned recurrent optimization process. The target model predictor ψ_θ is learned end-to-end using a robust version of the squared error ℓ and Gaussian confidence labels $a(y, y_i)$. For the second stage, it adopts the bounding-box regression technique proposed in ATOM.

3. Method

We propose a probabilistic regression model that integrates all advantages of confidence-based regression. However, unlike the aforementioned confidence-based models, our approach generates a predictive probability distribution $p(y|x, \theta)$ as output. The network is trained by minimizing the KL divergence between the predictive density $p(y|x, \theta)$ and the conditional ground-truth distribution $p(y|y_i)$, which models label noise and ambiguities in the task itself. During inference, a point estimate of the regressed value is obtained by maximizing the predicted density.

Our approach possesses a few important advantages compared to the confidence-based regression methods. In the latter, the prediction $s_\theta(y, x)$ can be difficult to interpret, and its value largely depend on the pseudo label function a and employed loss ℓ . In contrast, the probabilistic nature of our method allows reasoning about the uncertainty in the output. Moreover, in our approach, the pseudo label function a is replaced by the label-conditional distribution $p(y|y_i)$, which models noise and uncertainty in the annotation y_i . Lastly, in contrast to confidence-based regression, our approach does not require a choice of loss ℓ . Instead, we directly minimize the Kullback-Leibler (KL) divergence between the predictive distribution and the ground-truth. Next, we provide a general formulation of the proposed regression model, and apply it to tracking in Section 4.



⊕ Ground-truth center coordinate ● Target center prediction

Figure 2. Trackers are most often trained to predict the center coordinate of the ground-truth bounding box (red). This is a natural choice for the left frame and aligns well with the tracker prediction (green). Only two frames later (right), the motion of the tail has led to a radical shift in the ground-truth center location, which now lies in the background. This is not necessarily a natural definition of the target center coordinate, due to the minor change in the object appearance. Target center regression is thus an ambiguous task, where it is unclear how to define the “correct” value y_i . Our formulation models such ambiguity and uncertainty in the regression task by a distribution $p(y|y_i)$ of “correct” values.

3.1. Representation

In this section, we formulate an approach for effectively training the network to predict a probability distribution $p(y|x, \theta)$ of the output y given the input x . The density itself is represented using the formulation previously employed in probabilistic energy-based deep learning [27] and the recent deep conditional target densities [18],

$$p(y|x, \theta) = \frac{1}{Z_\theta(x)} e^{s_\theta(y, x)}, \quad Z_\theta(x) = \int e^{s_\theta(y, x)} dy. \quad (6)$$

As for the confidence-based methods described in Section 2, $s_\theta : \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{R}$ is a deep neural network mapping the output-input pair (y, x) to a scalar value. The expression (6) converts this value to a probability density by exponentiation and division by the normalizing constant $Z_\theta(x)$. In fact, it should be noted that (6) is a direct generalization of the SoftMax operation to an arbitrary output space \mathcal{Y} .

Since the output of the network represents a probability density over \mathcal{Y} , we can learn the network parameters θ by applying techniques for fitting a probability distribution to data. Given training sample pairs $\{(x_i, y_i)\}_i$, the simplest approach is to minimize the negative log-likelihood,

$$-\log p(y_i|x_i, \theta) = \log \left(\int e^{s_\theta(y, x_i)} dy \right) - s_\theta(y_i, x_i). \quad (7)$$

This strategy was recently successfully applied for a number of computer vision tasks [18], including bounding-box regression in visual tracking. One advantage of the negative log-likelihood loss (7) is that it only employs the training sample (x_i, y_i) itself, without further assumptions. However, this brings an important limitation, discussed next.

3.2. Label Uncertainty and Learning Objective

Compared to the negative log-likelihood loss (7), the confidence-based paradigm described in Section 2 enjoys a certain flexibility stemming from the pseudo label function $a(y, y_i)$. In practice, the design of $a(y, y_i)$ has been shown to be critical for tracking performance [4, 39]. We believe that this is mostly due to the inherent ambiguity of the task and the uncertainty in the label y_i itself. Most tracking approaches focus on regressing the center coordinate $y \in \mathbb{R}^2$ of the target in the image. However, for most objects, this is an inherently ambiguous and ill-defined task. While the center coordinate can be defined as the center of mass of the target bounding box, this is hardly a visually intuitive definition for a human, or a tracking algorithm for that matter.

Consider the example in Figure 2. When the target dog in the video raises its tail, the center of mass changes radically and ends up at a background pixel. On the other hand, the appearance and location of the object is almost unchanged. The tracker would thus naturally predict a similar target center location as before. This demonstrates that the definition of the target center is largely ambiguous and that the center of mass is often confusing for the tracker. The pseudo label function $a(y, y_i)$ can encapsulate this ambiguity by having a wider high-confidence peak, which has been shown beneficial for training tracking models [4]. Another source of uncertainty is label noise. Accurate bounding box annotation is a difficult task, especially in the presence of occlusions, motion blur, and for small objects, as shown in Figure 3. In other words, multiple annotators would naturally disagree for a given object, with some level of variation. This variation, or noise, in the annotation is most often ignored when training the network.

We propose to probabilistically model label noise and task ambiguities for the regression problem as a conditional ground-truth distribution $p(y|y_i)$. It characterizes the probability density of the ground-truth output value y , given the annotation y_i . Instead of the negative log-likelihood (7), we train the network to minimize the KL divergence to $p(y|y_i)$,

$$\begin{aligned} \text{KL}(p(\cdot|y_i), p(\cdot|x_i, \theta)) &= \int p(y|y_i) \log \frac{p(y|y_i)}{p(y|x_i, \theta)} dy \\ &\sim \log \left(\int e^{s_\theta(y, x_i)} dy \right) - \int s_\theta(y, x_i) p(y|y_i) dy. \end{aligned} \quad (8)$$

Here, \sim denotes equality up to a constant term. The sec-



Figure 3. Examples of noisy, incorrect or ambiguous ground truth bounding box annotations y_i from different datasets. These aspects are modeled by our label distribution $p(y|y_i)$.

ond line in (8) corresponds to the cross entropy between the two distributions and the discarded constant term is the negative entropy $\int p(y|y_i) \log p(y|y_i) dy$ of the label distribution. See Appendix A for a detailed derivation.

The loss (8) naturally integrates information about the uncertainty $p(y|y_i)$ in the annotated sample (x_i, y_i) . Unlike the pseudo label function $a(y|y_i)$ employed in confidence-based regression, $p(y|y_i)$ has a clear interpretation as a probability distribution. In fact, $p(y|y_i)$ could be empirically estimated by obtaining multiple annotations for a small subset of the data. In the case of a Gaussian model $p(y|y_i) = \mathcal{N}(y|y_i, \sigma^2)$, the variance σ^2 can be estimated as the empirical variance of these annotations. In this work, we simply consider σ^2 a hyper-parameter.

3.3. Training

In this section, we consider strategies for training the network parameters θ based on the loss (8). In practice, this requires approximating the two integrals in (8). We consider two techniques for this purpose, namely grid sampling and Monte Carlo integration with importance sampling.

Grid Sampling: For 2D image coordinate regression problems, *e.g.* the case of regressing the center of the tracked target, $y \in \mathcal{Y} \subset \mathbb{R}^2$ represents a location in the image. In this case, translational invariance is efficiently exploited by parametrizing $s_\theta(y, x) = f_\theta(x)(y)$, where f_θ is a Convolutional Neural Network (CNN). $s_\theta(y, x)$ is thus obtained by evaluating the output of the CNN at image coordinate y . Let $\{y^{(k)}\}_{k=1}^K \subset \mathcal{Y}$ be the set of uniform grid locations evaluated by the CNN $f_\theta(x)$ when applied to an image sample x . Further, let A be the area of a single grid cell. The uniform grid sampling automatically provided by the CNN generates the following approximation of the loss (8),

$$L_i = \log \left(A \sum_{k=1}^K e^{s_\theta(y^{(k)}, x_i)} \right) - A \sum_{k=1}^K s_\theta(y^{(k)}, x_i) p(y^{(k)} | y_i). \quad (9)$$

The final loss is then obtained by averaging L_i over all samples i in the mini-batch.

Monte Carlo Integration: For the more general regression problems, grid sampling does not necessarily provide any computational benefits. On the contrary, it scales poorly to higher dimensions and may induce a sampling bias due to the rigid grid. In the more general case, we therefore adopt the Monte Carlo (MC) based sampling strategy proposed in [18]. Specifically, we draw samples $y_i^{(k)} \sim q(y|y_i)$ from a proposal distribution $q(y|y_i)$ during training. The same samples are employed to approximate both integrals in (8),

$$L_i = \log \left(\frac{1}{K} \sum_{k=1}^K \frac{e^{s_\theta(y_i^{(k)}, x_i)}}{q(y_i^{(k)} | y_i)} \right) - \frac{1}{K} \sum_{k=1}^K s_\theta(y_i^{(k)}, x_i) \frac{p(y_i^{(k)} | y_i)}{q(y_i^{(k)} | y_i)}. \quad (10)$$

To accurately approximate the original loss (8), the proposal distribution $q(y|y_i)$ should ideally cover the label distribution $p(y|y_i)$ as well as regions with high predicted density $p(y|x_i, \theta)$. In [18] it was shown that a simple Gaussian mixture centered at the annotation y_i sufficed for a variety of tasks, including bounding box regression.

The loss (10) requires multiple evaluations of the network $s_\theta(y_i^{(k)}, x_i)$. In practice, however, computer vision architectures popularly employ deep backbone feature extractors $\phi_\theta(x)$, such as ResNet [20], generating a powerful representation of the image. The output value y can be fused at a late stage, such that $s_\theta(y, x) = f_\theta(y, \phi_\theta(x))$. This allows the computationally demanding feature extraction $\phi_\theta(x_i)$ to be shared among all samples $y_i^{(k)}$. Specifically for our purpose, such architectures have been successfully employed for bounding box regression in object detection and visual tracking problems [3, 9, 18, 24].

4. Tracking Approach

We apply the general probabilistic regression formulation introduced in Section 3 for the challenging and diverse task of visual target tracking.

4.1. Baseline Tracker: DiMP

We employ the recent state-of-the-art tracker DiMP [3] as our baseline. As briefly discussed in Section 2.2, the DiMP model contains two output branches.

Target Center Regression (TCR): The center regression branch aims to coarsely localize the target in the image by only regressing its center coordinate. This branch emphasizes robustness over accuracy. It consists of a linear convolutional output layer, whose weights w_θ are predicted by the network as an unrolled optimization process that minimizes an L^2 -based discriminative learning loss. This allows the tracker to robustly differentiate the target object from similar objects in the background. The target center confidence at location $y^{tc} \in \mathbb{R}^2$ in frame x is predicted similarly to (3), *i.e.* $s_\theta^{tc}(y^{tc}, x) = (w_\theta * \phi_\theta(x))(y^{tc})$, where ϕ_θ is the backbone feature extractor. This branch is trained in a meta-learning setting, with a confidence-based objective (2) using Gaussian pseudo labels a^{tc} and a robust L^2 loss,

$$\ell(s, a) = \begin{cases} (s - a)^2, & a > T \\ \max(0, s)^2, & a \leq T \end{cases}. \quad (11)$$

During tracking, the target center is regressed by densely computing confidence scores $s_\theta^{tc}(y^{tc}, x)$ within a wide search region in the frame x . We refer to [3] for details.

Bounding Box Regression (BBR): The BBR branch adopts the target conditional IoU-Net-based [24] architecture proposed in [9]. As discussed in Section 2.2, this branch predicts a confidence score $s_\theta^{bb}(y^{bb}, x)$ for a given

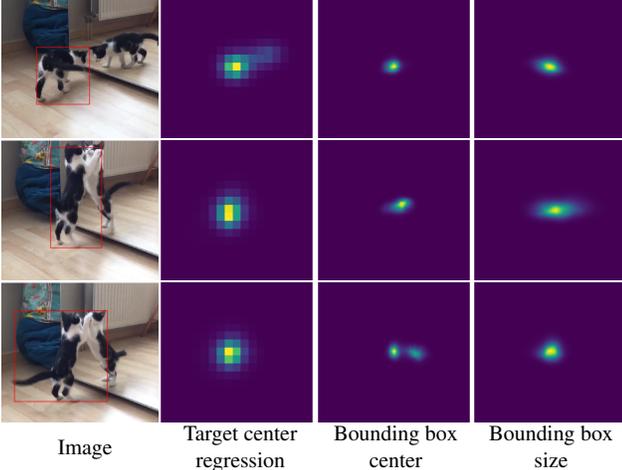


Figure 4. Visualization of the probability densities $p(y^{tc}|x, \theta)$ and $p(y^{bb}|x, \theta)$ predicted by the target center and bounding box regression branch respectively. The densities are centered at the predicted state (red box). The network captures uncertainties in the state, e.g. larger variance or multiple modes, for challenging cases. More examples and discussion are provided in Appendix D.

box $y^{bb} \in \mathbb{R}^4$. It is trained using the bounding box IoU as pseudo label $a^{bb}(y^{bb}, y_i^{bb})$ and the standard L^2 loss ℓ in (2). During tracking, the BBR branch is applied to fit an accurate bounding box to the target using gradient-based maximization of $s_\theta^{bb}(y^{bb}, x)$ w.r.t. y^{bb} . We refer to [9] for details.

4.2. Our Tracker: Probabilistic DiMP

We introduce a tracking approach with fully probabilistic output representations, obtained by integrating our regression formulation into both branches of the baseline DiMP. Example predicted densities are visualized in Figure 4.

Target Center Regression: We represent the predicted distribution of the target center coordinate $p(y^{tc}|x, \theta)$ by applying (6) to the network output $s_\theta^{tc}(y^{tc}, x)$. Since this branch is fully convolutional, we approximate the KL-divergence loss (8) for training using the grid sampling strategy (9). The conditional ground-truth density is set to a Gaussian $p(y^{tc}|y_i^{tc}) = \mathcal{N}(y^{tc}; y_i^{tc}, \sigma_{tc}^2)$ with the *same* variance parameter σ_{tc}^2 used for the corresponding pseudo label function a^{tc} in the baseline DiMP ($\frac{1}{4}$ th of the target size). For the optimization module, which predicts the convolution weights w_θ for the center regression branch, we use the KL-divergence loss (9) with an added L^2 regularization term. We modify the steepest descent based architecture [3] to employ a second order Taylor expansion, since the original Gauss-Newton approximation is limited to least-squares objectives. Our approach benefits from the fact that the resulting objective (9) is convex in w_θ for the linear predictor $s_\theta^{tc}(y^{tc}, x) = (w_\theta * \phi_\theta(x))(y^{tc})$, and thanks to efficient analytic expressions of the gradient and Hessian. See Appendix B for a detailed description of the optimizer module.

Bounding Box Regression: We use the same architec-

ture $s_\theta^{bb}(y^{bb}, x)$ as in [9, 3] and apply it in our probabilistic formulation (6). We follow the work of [18], which extended the same ATOM BBR module [9] to the probabilistic setting using the negative log-likelihood loss (7) and an MC-based approximation. In this work, we further integrate the label distribution $p(y^{bb}|y_i^{bb})$ to model the noise and uncertainty in the bounding box annotations, and minimize the KL-divergence (8) using MC sampling (10). Specifically, we use an isotropic Gaussian distribution $p(y^{bb}|y_i^{bb}) = \mathcal{N}(y^{bb}; y_i^{bb}, \sigma_{bb}^2)$ and set $\sigma_{bb} = 0.05$. For a fair comparison, we use the same proposal distribution $q(y^{bb}|y_i^{bb}) = \frac{1}{2}\mathcal{N}(y^{bb}; y_i^{bb}, 0.05^2) + \frac{1}{2}\mathcal{N}(y^{bb}; y_i^{bb}, 0.5^2)$ and bounding box parametrization as in [18].

Details: Our entire network is trained jointly end-to-end using the *same* strategy and settings as for the original DiMP [3], by integrating it into the publicly available PyTracking framework [7]. The training splits of the LaSOT [13], GOT10k [23], TrackingNet [33] and COCO [30] are used, running 50 epochs with 1000 iterations each. We also preserve the tracking procedure and settings in DiMP, only performing minimal changes, which are forced by the probabilistic output representation provided by our model. Due to different scaling of the network output, we accordingly change the threshold for which the target is reported missing and the gradient step length used for bounding box regression. We refer to [3, 7] for detailed description of training and inference settings. Our code is available at [7].

5. Experiments

We perform a detailed analysis of our approach along with comprehensive comparisons with state-of-the-art on seven datasets. The ResNet-18 and 50 versions of our tracker operates at about 40 and 30 FPS respectively.

5.1. Comparison of Regression Models

We first analyze the impact of different regression formulations for tracking. In each case, we train a separate version of the DiMP baseline, with the ResNet-18 backbone. We compare four different approaches. **L2:** Standard squared loss, used in the baseline DiMP for Bounding Box Regression (BBR). **R-L2:** The Robust L2 loss (11), employed in the baseline DiMP for the Target Center Regression (TCR) branch. **NLL:** The probabilistic Negative Log-Likelihood formulation (7) proposed in [18]. **Ours:** Trained using the KL-divergence (8) as described in Section 4.2.

Following [3], we perform this analysis of our approach on the combined OTB-100 [47], UAV123 [32] and NFS [14] dataset, totaling 323 diverse videos, and report the average over 5 runs. We report the Overlap Precision OP_T , i.e. percentage of frames with an IoU overlap larger than T , along with the main metric $AUC = \int_0^1 OP_T$ (see [47]).

Complete Tracker: We first analyze the performance when each regression model is applied to the entire tracker,

Model	Complete tracker Both BBR and TCR			Bounding Box reg. TCR formulation: R-L2			Target Center reg. BBR formulation: Ours		
	AUC	OP _{0.50}	OP _{0.75}	AUC	OP _{0.50}	OP _{0.75}	AUC	OP _{0.50}	OP _{0.75}
L2	63.1	78.0	50.2	63.8	79.2	50.6	64.8	80.9	54.1
R-L2	63.8	79.2	50.6	63.8	79.2	50.6	65.8	82.0	54.1
NLL	63.0	78.5	51.5	65.0	81.0	52.8	63.2	79.0	52.6
Ours	65.5	81.6	54.1	65.8	82.0	54.1	65.5	81.6	54.1

Table 1. Analysis of four different regression models for tracking on the combined OTB100-NFS-UAV123 dataset. In the left section, the model is applied to both branches (BBR and TCR) of the network. In the center and right sections, we exclusively analyze their impact on BBR and TCR respectively. In the former case, R-L2 is always employed for TCR, while the right section uses Our approach for BBR in all cases. See text for details.

i.e. for both the BBR and TCR branch. The results are reported in the left section of Table 1. The R-L2, which corresponds to standard DiMP, improves 0.8% AUC over the standard L2 loss. Our model outperforms the R-L2 baseline with a gain of 1.7% in AUC.

Bounding Box Regression (BBR): We exclusively analyze the impact of each model for BBR by using the baseline DiMP R-L2 formulation for TCR in all cases. Thus, only the BBR branch is affected. Results are provided in the center section of Table 1. The baseline [3, 9], employing L2 loss to predict IoU, achieves 63.8 AUC. The NLL formulation [18] gives a substantial 1.2% gain. By modeling uncertainty in the annotation and using the KL loss (8), our approach achieves an additional improvement of 0.8%.

Target Center Regression (TCR): Similarly, we compare the models for TCR by employing our approach for the BBR branch in all cases. The results, reported in the right section of Table 1, show that probabilistic NLL is not suitable for TCR. This is likely due to the inherent ambiguity of the problem, which is not accounted for. By explicitly modeling label uncertainty $p(y|y_i)$, our formulation achieves a 2.3% gain in AUC, further outperforming the standard L2 model. The R-L2 model, which is specifically designed for TCR [3], achieves a marginally better performance. However, note that this result is achieved when combined with *our* formulation for BBR.

In conclusion, our probabilistic formulation outperforms the baseline DiMP, employing L2 and Robust L2 for TCR

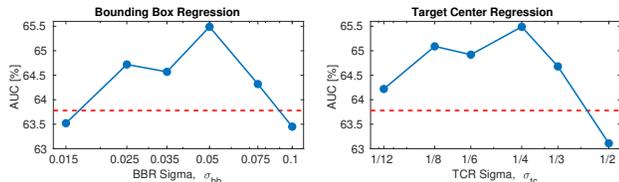
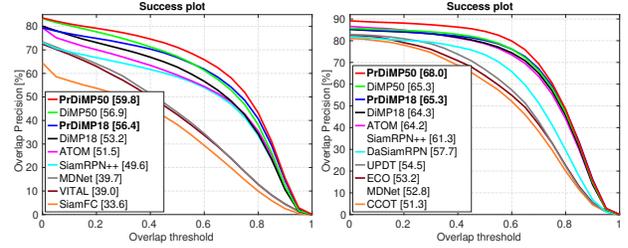


Figure 5. The impact of modeling the label uncertainty $p(y|y_i) = \mathcal{N}(y; y_i, \sigma^2)$ by varying the standard deviation σ for bounding box regression (left) and target center regression (right). We show AUC on the combined OTB100-NFS-UAV123 datasets and display the baseline DiMP-18 result as a red dashed line.



(a) LaSOT [13]

(b) UAV123 [32]

Figure 6. Success plots, showing OP_T , on LaSOT and UAV123, showing average over 5 runs for our method. Our approach outperforms DiMP by a large margin in AUC, shown in the legend.

and BBR respectively, by a large 1.7% margin in AUC, while achieving the best overall performance for all regression models. In the subsequent experiments, we employ our formulation for both the BBR and TCR branch.

5.2. Analysis of Label Uncertainty

In this section, we further analyze the importance of our ability to model the uncertainty $p(y|y_i)$ in the annotation y_i . This is probed by investigating the impact of the standard deviation parameter σ for our Gaussian model of the label noise $p(y|y_i) = \mathcal{N}(y; y_i, \sigma^2)$ (see Section 4.2). This is individually performed for the bounding box (σ_{bb}) and target center (σ_{tc}) regression branch. We use the same experiment setup as described in the previous section, reporting the average of five runs over the combined OTB100-NFS-UAV123 dataset, and using the ResNet-18 backbone.

Figure 5 shows the tracking performance in AUC when varying the standard deviations σ_{bb} and σ_{tc} over a wide range of values. The baseline DiMP-18 performance is shown in red for reference. A similar overall trend is observed in both cases. A too large standard deviation σ leads to poor results, since the over-estimation in the label uncertainty forces the tracker to output too uncertain predictions. However, a small σ causes overfitting and over-confident predictions, leading to sub-optimal performance. Properly modeling the label uncertainty is thus critical for visual tracking, due to the ambiguities in the regression tasks and annotation noise. We also note that by outperforming the baseline DiMP over a wide range of σ_{bb} and σ_{tc} values, our approach is not sensitive to specific settings.

5.3. State-of-the-Art

We compare our approach, termed **PrDiMP**, on seven tracking benchmarks. We evaluate two versions, PrDiMP18 and PrDiMP50, employing ResNet-18 and ResNet-50 respectively as backbone. The same settings and parameters are used for one-pass evaluation on all datasets. To ensure the significance in the results, we report the average over 5 runs for all datasets, unless the specific protocol requires otherwise. Additional results are provided in Appendix C.

LaSOT [13]: We first compare on the large-scale La-

	SiamFC	MDNet	UPDT	DaSiam-	ATOM	SiamRPN++	DiMP18	DiMP50	PrDiMP18	PrDiMP50
	[2]	[34]	[4]	RPN [52]	[9]	[28]	[3]	[3]		
Precision	53.3	56.5	55.7	59.1	64.8	69.4	66.6	68.7	69.1	70.4
Norm. Prec.	66.6	70.5	70.2	73.3	77.1	80.0	78.5	80.1	80.3	81.6
Success (AUC)	57.1	60.6	61.1	63.8	70.3	73.3	72.3	74.0	75.0	75.8

Table 2. Results on the TrackingNet [33] test set in terms of precision, normalized precision, and success (AUC). Both our PrDiMP versions outperform previous methods by a significant margin.

SOT dataset. The test set contains 280 long videos (2500 frames in average), thus emphasizing the robustness of the tracker, along with its accuracy. The success plot in Figure 6a, showing the overlap precision OP_T as a function of the threshold T , is computed as the average over 5 runs. Trackers are ranked w.r.t. their AUC score, shown in the legend. Our approach outperforms the previous best tracker, DiMP, by a large margin of 3.2% and 2.9% in AUC when using ResNet-18 and ResNet-50 respectively. The improvement in OP_T is most prominent for $T > 0.3$, demonstrating superior accuracy achieved by our bounding box regression.

TrackingNet [33]: This is a large-scale tracking dataset with high variety in terms of classes and scenarios. The test set contains over 500 videos without publicly available ground-truth. The results, shown in Table 2, are obtained through an online evaluation server. Both our versions outperform all previous approaches in the main Success (AUC) metric by a significant margin. With the same ResNet-50 backbone, our approach achieves a gain of 2.5% and 1.8% in AUC over SiamRPN++ [28] and DiMP [3] respectively.

VOT2018 [25]: Next, we evaluate on the 2018 edition of the Visual Object Tracking challenge. We compare with the top methods in the challenge [25], as well as more recent methods. The dataset contains 60 videos. Trackers are restarted at failure by the evaluation system. The performance is then decomposed into accuracy and robustness, defined using IoU overlap and failure rate respectively. The main EAO metric takes both these aspects into account. The results, computed over 15 repetitions as specified in the protocol, are shown in Table 3. Our PrDiMP50 achieves the best overall performance, with the highest accuracy and competitive robustness compared to previous methods.

GOT10k [23]: This dataset contains 10,000 sequences for training and 180 for testing. We follow the defined protocol [23] and *only* train on the specified GOT10k training set for this experiments, while keeping all other settings the same. By having no overlap in object classes between training and testing, GOT10k also benchmarks the generalizability of the tracker to novel objects. The results in Table 4, obtained through a evaluation server, are reported in terms of SR_T and AO, which are equivalent to OP_T and

	RCO	UPDT	DaSiam-	MFT	LADCF	ATOM	SiamRPN++	DiMP18	DiMP50	PrDiMP18	PrDiMP50
	[25]	[4]	RPN [52]	[25]	[49]	[9]	[28]	[3]	[3]		
EAO	0.376	0.378	0.383	0.385	0.389	0.401	0.414	0.402	0.440	0.385	0.442
Robustness	0.155	0.184	0.276	0.140	0.159	0.204	0.234	0.182	0.153	0.217	0.165
Accuracy	0.507	0.536	0.586	0.505	0.503	0.590	0.600	0.594	0.597	0.607	0.618

Table 3. Results on the VOT2018 challenge dataset [25] in terms of expected average overlap (EAO), robustness and accuracy.

	CF2	ECO	CCOT	GOTURN	SiamFC	SiamFCv2	ATOM	DiMP18	DiMP50	PrDiMP18	PrDiMP50
	[31]	[8]	[12]	[21]	[2]	[44]	[9]	[3]	[3]		
$SR_{0.50}$	29.7	30.9	32.8	37.5	35.3	40.4	63.4	67.2	71.7	71.5	73.8
$SR_{0.75}$	8.8	11.1	10.7	12.4	9.8	14.4	40.2	44.6	49.2	50.4	54.3
AO	31.5	31.6	32.5	34.7	34.8	37.4	55.6	57.9	61.1	61.2	63.4

Table 4. State-of-the-art comparison on the GOT10k test set [23]. The metrics average overlap (AO) and success rate (SR_T) are equivalent to AUC and OP_T respectively (described in Section 5.1). The evaluated methods are trained only on the GOT10k training set, and must therefore generalize to novel classes.

	DaSiam-	CCOT	MDNet	ECO	ATOM	SiamRPN++	UPDT	DiMP18	DiMP50	PrDiMP18	PrDiMP50
	RPN [52]	[9]	[12]	[34]	[8]	[28]	[4]	[3]	[3]		
OTB-100	65.8	68.2	67.8	69.1	66.9	69.6	70.2	66.0	68.4	68.0	69.6
NFS	-	48.8	42.2	46.6	58.4	-	53.7	61.0	62.0	63.3	63.5

Table 5. Comparison with state-of-the-art on the OTB-100 [47] and NFS [14] datasets in terms of overall AUC score. The average value over 5 runs is reported for our approach.

AUC [23]. Compared to DiMP, our method achieve large gains of 3.3% and 2.3% in terms of the main AO metric, when using ResNet-18 and -50 respectively.

UAV123 [32]: This challenging dataset, containing 123 videos, is designed to benchmark trackers for UAV applications. It features small objects, fast motions, and distractor objects. The results are shown in Figure 6b, where OP_T is plotted over IoU thresholds T and AUC is shown in the legend. ATOM [9] and DiMP50 obtain 64.2% and 65.3% respectively. Our PrDiMP50 achieves a remarkable 68.0%, surpassing previous methods by a significant margin.

OTB-100 [47]: For reference, we report results on the OTB-100 dataset. While this dataset has served an important role in the development of trackers since its release, it is known to have become highly saturated over recent years, as shown in Table 5. Still, our approach performs similarly to the top correlation filter methods, such as UPDT [4], while on par with the end-to-end trained SiamRPN++.

NFS [14]: Lastly, we report results on the 30 FPS version of the Need for Speed (NFS) dataset, containing fast motions and challenging distractors. As shown in Table 5, our approach achieves a substantial improvement over the previous state-of-the-art on this dataset.

6. Conclusions

We propose a probabilistic regression formulation, where the network is trained to predict the conditional density $p(y|x, \theta)$ of the output y given the input x . The density is parametrized by the architecture itself, allowing the representation of highly flexible distributions. The network is trained by minimizing the KL-divergence to the label distribution $p(y|y_i)$, which is introduced to model annotation noise and task ambiguities. When applied for the tracking task, our approach outperforms the baseline DiMP [3] and sets a new state-of-the-art on six datasets.

Acknowledgments: This work was partly supported by the ETH Zürich Fund (OK), a Huawei Technologies Oy (Finland) project, an Amazon AWS grant, and Nvidia.

References

- [1] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe. Tracking without bells and whistles. In *ICCV*, 2019. [1](#)
- [2] Luca Bertinetto, Jack Valmadre, João F Henriques, Andrea Vedaldi, and Philip HS Torr. Fully-convolutional siamese networks for object tracking. In *ECCV workshop*, 2016. [1](#), [2](#), [3](#), [8](#), [13](#)
- [3] Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte. Learning discriminative model prediction for tracking. In *IEEE/CVF International Conference on Computer Vision, ICCV, Seoul, South Korea*, pages 6181–6190, 2019. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [8](#), [11](#), [12](#), [13](#)
- [4] Goutam Bhat, Joakim Johlander, Martin Danelljan, Fahad Shahbaz Khan, and Michael Felsberg. Unveiling the power of deep tracking. In *15th European Conference on Computer Vision ECCV, Munich, Germany*, pages 493–509, 2018. [4](#), [8](#), [13](#)
- [5] David S. Bolme, J. Ross Beveridge, Bruce A. Draper, and Yui Man Lui. Visual object tracking using adaptive correlation filters. In *CVPR*, 2010. [1](#), [3](#)
- [6] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7291–7299, 2017. [1](#), [2](#)
- [7] Martin Danelljan and Goutam Bhat. PyTracking: Visual tracking library based on PyTorch. <https://github.com/visionml/pytracking>, 2019. Accessed: 16/09/2019. [6](#)
- [8] Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. ECO: efficient convolution operators for tracking. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR, Honolulu, HI, USA*, pages 6931–6939, 2017. [8](#), [13](#)
- [9] Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. ATOM: accurate tracking by overlap maximization. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR, Long Beach, CA, USA*, pages 4660–4669, 2019. [1](#), [3](#), [5](#), [6](#), [7](#), [8](#), [12](#), [13](#)
- [10] Martin Danelljan, Gustav Häger, Fahad Shahbaz Khan, and Michael Felsberg. Learning spatially regularized correlation filters for visual tracking. In *IEEE International Conference on Computer Vision, ICCV, Santiago, Chile*, pages 4310–4318, 2015. [1](#), [3](#)
- [11] Martin Danelljan, Gustav Häger, Fahad Shahbaz Khan, and Michael Felsberg. Discriminative scale space tracking. *IEEE Transactions on Pattern Analysis Machine Intelligence, TPAMI*, 39(8):1561–1575, 2017. [2](#)
- [12] Martin Danelljan, Andreas Robinson, Fahad Shahbaz Khan, and Michael Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. In *14th European Conference on Computer Vision ECCV, Amsterdam, The Netherlands*, pages 472–488, 2016. [1](#), [3](#), [8](#), [13](#)
- [13] Heng Fan, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Hexin Bai, Yong Xu, Chunyuan Liao, and Haibin Ling. Lasot: A high-quality benchmark for large-scale single object tracking. *CoRR*, abs/1809.07845, 2018. [6](#), [7](#), [12](#)
- [14] Hamed Kiani Galoogahi, Ashton Fagg, Chen Huang, Deva Ramanan, and Simon Lucey. Need for speed: A benchmark for higher frame rate object tracking. In *ICCV*, 2017. [6](#), [8](#), [13](#)
- [15] Hamed Kiani Galoogahi, Terence Sim, and Simon Lucey. Correlation filters with limited boundaries. In *CVPR*, 2015. [1](#)
- [16] Clément Godard, Oisín Mac Aodha, and Gabriel J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, pages 6602–6611, 2017. [2](#)
- [17] Qing Guo, Wei Feng, Ce Zhou, Rui Huang, Liang Wan, and Song Wang. Learning dynamic siamese network for visual object tracking. In *ICCV*, 2017. [1](#), [3](#)
- [18] Fredrik K. Gustafsson, Martin Danelljan, Goutam Bhat, and Thomas B. Schön. DCTD: deep conditional target densities for accurate regression. *CoRR*, abs/1909.12297, 2019. [2](#), [4](#), [5](#), [6](#), [7](#)
- [19] Sam Hare, Stuart Golodetz, Amir Saffari, Vibhav Vineet, Ming-Ming Cheng, Stephen L. Hicks, and Philip H. S. Torr. Struck: Structured output tracking with kernels. *IEEE Trans. Pattern Anal. Mach. Intell.*, 38(10):2096–2109, 2016. [3](#)
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [5](#)
- [21] David Held, Sebastian Thrun, and Silvio Savarese. Learning to track at 100 fps with deep regression networks. In *ECCV*, 2016. [8](#), [13](#)
- [22] João F. Henriques, Rui Caseiro, Pedro Martins, and Jorge Batista. High-speed tracking with kernelized correlation filters. *TPAMI*, 37(3):583–596, 2015. [1](#), [2](#), [3](#)
- [23] Lianghua Huang, Xin Zhao, and Kaiqi Huang. GOT-10k: A large high-diversity benchmark for generic object tracking in the wild. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. [6](#), [8](#), [13](#)
- [24] Boru Jiang, Ruixua Luo, Jiayuan Mao, Tete Xiao, and Yunying Jiang. Acquisition of localization confidence for accurate object detection. In *ECCV*, 2018. [1](#), [2](#), [3](#), [5](#)
- [25] Matej Kristan, Ales Leonardis, Jiri Matas, Michael Felsberg, Roman Pfugfelder, Luka Cehovin Zajc, Tomas Vojir, Goutam Bhat, Alan Lukezic, Abdelrahman Eldesokey, Gustavo Fernandez, and et al. The sixth visual object tracking vot2018 challenge results. In *ECCV workshop*, 2018. [8](#)
- [26] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 734–750, 2018. [1](#), [2](#)
- [27] Yann LeCun, Sumit Chopra, Raia Hadsell, Fu Jie Huang, and et al. A tutorial on energy-based learning. In *PREDICTING STRUCTURED DATA*. MIT Press, 2006. [2](#), [4](#)
- [28] Bo Li, Wei Wu, Qiang Wang, Fangyi Zhang, Junliang Xing, and Junjie Yan. Siamrpn++: Evolution of siamese visual tracking with very deep networks. In *CVPR*, 2019. [1](#), [3](#), [8](#), [12](#), [13](#)
- [29] Bo Li, Junjie Yan, Wei Wu, Zheng Zhu, and Xiaolin Hu. High performance visual tracking with siamese region proposal network. In *CVPR*, 2018. [1](#), [3](#)

- [30] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, 2014. 6
- [31] Chao Ma, Jia-Bin Huang, Xiaokang Yang, and Ming-Hsuan Yang. Hierarchical convolutional features for visual tracking. In *ICCV*, 2015. 1, 3, 8, 13
- [32] Matthias Mueller, Neil Smith, and Bernard Ghanem. A benchmark and simulator for uav tracking. In *ECCV*, 2016. 6, 7, 8
- [33] Matthias Müller, Adel Bibi, Silvio Giancola, Salman Al-Subaihi, and Bernard Ghanem. Trackingnet: A large-scale dataset and benchmark for object tracking in the wild. In *ECCV*, 2018. 6, 8
- [34] Hyeonseob Nam and Bohyung Han. Learning multi-domain convolutional neural networks for visual tracking. In *CVPR*, 2016. 8, 12, 13
- [35] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39:1137–1149, 2015. 1, 2
- [36] Tomas Simon, Hanbyul Joo, Iain A. Matthews, and Yaser Sheikh. Hand keypoint detection in single images using multiview bootstrapping. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, pages 4645–4653, 2017. 1
- [37] Yibing Song, Chao Ma, Xiaohe Wu, Lijun Gong, Linchao Bao, Wangmeng Zuo, Chunhua Shen, Rynson W. H. Lau, and Ming-Hsuan Yang. VITAL: visual tracking via adversarial learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR, Salt Lake City, UT, USA*, pages 8990–8999, 2018. 12
- [38] Chris Stauffer and W. Eric L. Grimson. Adaptive background mixture models for real-time tracking. In *CVPR*, pages 2246–2252, 1999. 1
- [39] Yao Sui, Ziming Zhang, Guanghui Wang, Yafei Tang, and Li Zhang. Real-time visual tracking: Promoting the robustness of correlation filter learning. In *European Conference on Computer Vision ECCV*, pages 662–678, 2016. 4
- [40] Chong Sun, Dong Wang, Huchuan Lu, and Ming-Hsuan Yang. Correlation tracking via joint discrimination and reliability learning. In *CVPR*, 2018. 1
- [41] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 8934–8943, 2018. 2
- [42] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, 2019. 1, 2
- [43] Ioannis Tsochantaridis, Thorsten Joachims, Thomas Hofmann, and Yasemin Altun. Large margin methods for structured and interdependent output variables. *J. Mach. Learn. Res.*, 6:1453–1484, 2005. 3
- [44] Jack Valmadre, Luca Bertinetto, João F Henriques, Andrea Vedaldi, and Philip H. S. Torr. End-to-end representation learning for correlation filter based tracking. In *CVPR*, 2017. 8, 13
- [45] Qiang Wang, Zhu Teng, Junliang Xing, Jin Gao, Weiming Hu, and Stephen J. Maybank. Learning attentions: Residual attentional siamese network for high performance online visual tracking. In *CVPR*, 2018. 1, 3
- [46] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 4724–4732, 2016. 1
- [47] Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. Object tracking benchmark. *TPAMI*, 37(9):1834–1848, 2015. 6, 8, 13
- [48] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 466–481, 2018. 1, 2
- [49] Tianyang Xu, Zhen-Hua Feng, Xiao-Jun Wu, and Josef Kittler. Learning adaptive discriminative correlation filters via temporal consistency preserving spatial feature selection for robust visual tracking. *CoRR*, abs/1807.11348, 2018. 8
- [50] Jianming Zhang, Shugao Ma, and Stan Sclaroff. MEEM: robust tracking via multiple experts using entropy minimization. In *ECCV*, 2014. 3
- [51] Lichao Zhang, Abel Gonzalez-Garcia, Joost van de Weijer, Martin Danelljan, and Fahad Shahbaz Khan. Learning the model update for siamese trackers. In *IEEE/CVF International Conference on Computer Vision, ICCV, Seoul, South Korea*, pages 4009–4018, 2019. 1, 3
- [52] Zheng Zhu, Qiang Wang, Li Bo, Wei Wu, Junjie Yan, and Weiming Hu. Distractor-aware siamese networks for visual object tracking. In *ECCV*, 2018. 8, 13

Appendix

This Appendix provides detailed derivations and results. Appendix A derive the employed loss in (8) from the KL divergence. We then provide details about the target center regression module in Appendix B. Lastly, we report more detailed results in Appendix C and provide additional visualizations of the predicted distributions in Appendix D.

A. Derivation of KL Divergence Loss

Here, we derive the loss (8) from the KL divergence between the predicted distribution $p(y|x_i, \theta)$ and the ground-truth density $p(y|y_i)$. Starting from the definition of the KL divergence and inserting our formulation (6), we achieve,

$$\begin{aligned} \text{KL}(p(\cdot|y_i), p(\cdot|x_i, \theta)) &= \int p(y|y_i) \log \frac{p(y|y_i)}{p(y|x_i, \theta)} dy \\ &= \int p(y|y_i) \log \frac{p(y|y_i)}{e^{s_\theta(y, x_i)} / Z_\theta(x_i)} dy \\ &= \int p(y|y_i) \left(\log p(y|y_i) - \log e^{s_\theta(y, x_i)} + \log Z_\theta(x_i) \right) dy \\ &= \int p(y|y_i) \log p(y|y_i) dy \\ &\quad + \log Z_\theta(x_i) - \int p(y|y_i) s_\theta(y, x) dy \\ &\sim \log \left(\int e^{s_\theta(y, x_i)} dy \right) - \int s_\theta(y, x_i) p(y|y_i) dy. \end{aligned} \quad (12)$$

In the last row, we have discarded the first term (the negative entropy of $p(y|y_i)$) and substituted the definition of the partition function $Z_\theta(x_i)$ from (6).

B. Target Center Regression Module

In this section, we give a detailed description and derivation of the optimization module employed for target center regression in our PrDiMP tracker. The goal of the optimizer module is to predict the weights w_θ of the target center regression component,

$$s_\theta^{\text{tc}}(y^{\text{tc}}, x) = (w_\theta * \phi_\theta(x))(y^{\text{tc}}). \quad (13)$$

Here, x is an input image, $y^{\text{tc}} \in \mathbb{R}^2$ is an image coordinate and ϕ_θ is the backbone feature extractor (see Section 4.2). The weights w_θ are learned from a set of training (support) images $\{z_j\}_{j=1}^n$ and corresponding target bounding box annotations $\{y_j^{\text{bb}}\}_{j=1}^n$. As in the baseline DiMP [3], these images are sampled from an interval within each sequence during training. During tracking, the images z_j are obtain by performing augmentations on the first frame, and by gradual update of the memory.

In DiMP, the optimizer module is derived by applying the steepest descent algorithm to a least-squares objective. In our case however, the employed loss function

does not admit a least-squares formulation. In this work, we therefore replace the Gauss-Newton approximation with the quadratic Newton approximation in order to compute the step-length necessary for the steepest descent algorithm. We derive closed form solutions of all operations, ensuring simple integration of the optimizer module as a series of deep neural network layers.

Similarly to our offline training objective, we let the optimizer module minimize the KL-divergence based learning loss (8). We also add an L^2 regularization term to benefit generalization to unseen frames. The loss is thus formulated as follows,

$$L(w_\theta) = \sum_{j=1}^n \gamma_j L_{\text{CE}}(\tilde{z}_j * w_\theta; p_j) + \frac{\lambda}{2} \|w_\theta\|^2. \quad (14)$$

The non-negative scalars λ and γ_j control the impact of the regularization term and sample z_j respectively. We also make the following definitions for convenience,

$$\tilde{z}_j = \phi_\theta(z_j) \quad \text{Extracted image features.} \quad (15a)$$

$$p_j^{(k)} = p(y^{\text{tc}, (k)} | y_j^{\text{tc}}) \quad \text{Label density at location } k. \quad (15b)$$

$$s_j^{(k)} = (\tilde{z}_j * w_\theta)(y^{\text{tc}, (k)}) \quad \text{Target scores at location } k. \quad (15c)$$

$$\hat{p}_j^{(k)} = \frac{\exp(s_j^{(k)})}{\sum_{l=1}^K \exp(s_j^{(l)})} \quad \text{Spatial SoftMax of } s_j. \quad (15d)$$

Note that we use superscript $k \in \{1, \dots, K\}$ to denote spatial grid location $y^{\text{tc}, (k)} \in \mathbb{R}^2$. In the following, the quantities in (15b)-(15d) are either seen as vectors in \mathbb{R}^K or 2D-maps $\mathbb{R}^{H \times W}$ (with $K = HW$), as made clear from the context.

In (14), the per-sample loss L_{CE} is the grid approximation (9) of the original KL-divergence objective. Without loss of generality, we may assume $A = 1$, obtaining

$$\begin{aligned} L_{\text{CE}}(s; p) &= \log \left(\sum_{k=1}^K e^{s^{(k)}} \right) - \sum_{k=1}^K s^{(k)} p^{(k)} \\ &= \log(\mathbf{1}^T e^s) - p^T s. \end{aligned} \quad (16)$$

Here, $\mathbf{1}^T = [1, \dots, 1]$ denotes a vector of ones. Note that the grid approximation thus corresponds to the SoftMax-Cross Entropy loss, commonly employed for classification.

To derive the optimization module, we adopt the steepest descent formulation [3], but employ the Newton approximation discussed above. This results in the following optimization strategy,

$$w_\theta^{(i+1)} = w_\theta^{(i)} - \alpha^{(i)} \nabla L(w_\theta^{(i)}) \quad (17a)$$

$$\alpha^{(i)} = \frac{\nabla L(w_\theta^{(i)})^T \nabla L(w_\theta^{(i)})}{\nabla L(w_\theta^{(i)})^T H(w_\theta^{(i)}) \nabla L(w_\theta^{(i)})}. \quad (17b)$$

Algorithm 1 Target Model Prediction ψ_θ for Center Reg.

Require: Training (support) images $\{z_j\}_{j=1}^n$
Require: Corresponding box annotations $\{\tilde{y}_j^{\text{bb}}\}_{j=1}^n$

- 1: $\tilde{z}_j = \phi_\theta(z_j)$, $j = 1, \dots, n$ *Extract features*
- 2: $w_\theta^{(0)} \leftarrow \psi_\theta^{\text{init}}(\{(z_j, \tilde{y}_j^{\text{bb}})\}_{j=1}^n)$ *Initialize weights [3]*
- 3: **for** $i = 0, \dots, N_{\text{iter}} - 1$ **do** *Optimizer module loop*
- 4: $s_j \leftarrow \tilde{z}_j * w_\theta^{(i)}$ *Using (15a)*
- 5: $\hat{p}_j \leftarrow \text{SoftMax}(s_j)$ *Using (15d)*
- 6: $g \leftarrow \nabla L(w_\theta^{(i)})$ *Using (19)*
- 7: $\alpha^{(i)} = \frac{g^\top g}{g^\top H(w_\theta^{(i)}) g}$ *Using (21) for the denom.*
- 8: $w_\theta^{(i+1)} \leftarrow w_\theta^{(i)} - \alpha^{(i)} g$ *Update weights*
- 9: **end for**

Here, $\nabla L(w_\theta^{(i)})$ and $H(w_\theta^{(i)})$ is the gradient and Hessian of L (14), evaluated at the current estimate $w_\theta^{(i)}$ of the weights. To implement (17), we efficiently compute these quantities by deriving closed-form expressions of both.

First, we can first easily compute the gradient and hessian of (16) w.r.t. s as,

$$\nabla_s L_{\text{CE}}(s; p) = \hat{p} - p \quad (18a)$$

$$\frac{\partial^2}{\partial s^2} L_{\text{CE}}(s; p) = \text{diag}(\hat{p}) - \hat{p}\hat{p}^\top \quad (18b)$$

Here, \hat{p} is the SoftMax of s as defined in (15d). By applying (18a) together with the chain rule, we can compute the gradient of (14) as,

$$\begin{aligned} \nabla L(w_\theta) &= \sum_{j=1}^n \gamma_j \left[\frac{\partial s_j}{\partial w_\theta} \right]^\top \nabla_s L_{\text{CE}}(s_j; p_j) + \lambda w_\theta \\ &= \sum_{j=1}^n \gamma_j [\tilde{z}_j *]^\top (\hat{p}_j - p_j) + \lambda w_\theta. \end{aligned} \quad (19)$$

We denote the transpose of the linear convolution operator $w \mapsto \tilde{z} * w$ as $[\tilde{z} *]^\top$, which corresponds to the transpose of the Jacobian $\frac{\partial s_j}{\partial w_\theta} = [\tilde{z} *]^\top$. By another differentiation w.r.t. w_θ , using (18b), the chain rule and the linearity of s_j in w_θ , we obtain the Hessian of (14) as

$$\begin{aligned} H(w_\theta) &= \frac{\partial^2}{\partial w_\theta^2} L(w_\theta) \\ &= \sum_{j=1}^n \gamma_j \left[\frac{\partial s_j}{\partial w_\theta} \right]^\top \left[\frac{\partial}{\partial s} L_{\text{CE}}(s_j; p_j) \right] \left[\frac{\partial s_j}{\partial w_\theta} \right] + \lambda I \\ &= \sum_{j=1}^n \gamma_j [\tilde{z}_j *]^\top (\text{diag}(\hat{p}_j) - \hat{p}_j \hat{p}_j^\top) [\tilde{z}_j *] + \lambda I. \end{aligned} \quad (20)$$

Here, I denotes the identity matrix. We further obtain a simple expression of the denominator of the step-length in

(17) by evaluating the product,

$$\begin{aligned} g^\top H(w_\theta) g &= \\ &= \sum_{j=1}^n \gamma_j (\tilde{z}_j * g)^\top (\text{diag}(\hat{p}_j) - \hat{p}_j \hat{p}_j^\top) (\tilde{z}_j * g) + \lambda g^\top g \\ &= \sum_{j=1}^n \gamma_j v_j^\top (\hat{p}_j \cdot (v_j - \hat{p}_j^\top v_j)) + \lambda g^\top g, \quad v_j = \tilde{z}_j * g. \end{aligned} \quad (21)$$

Here, \cdot denotes element-wise multiplication. We summarize the full optimization module in Algorithm 1.

C. Detailed Results

We provide more detailed results from the state-of-the-art comparison performed in Section 5.3.

C.1. LaSOT

In addition to the success plot shown in Section 5.3, we here provide the normalized precision plot over the LaSOT [13] test set, containing 280 videos. The normalized precision score NPr_D is computed as the percentage of frames where the normalized distance (relative to the target size) between the predicted and ground-truth target center location is less than a threshold D . NPr_D is plotted over a range of thresholds $D \in [0, 0.5]$. The trackers are ranked using the area under this curve, which is shown in the legend (see [13] for details). The normalized precision plot is shown in Figure 7. For the ResNet-18 and 50 versions of our PrDiMP, we report the average result over 5 runs. We compare with state-of-the-art trackers DiMP [3], ATOM [9], SiamRPN++ [28], MDNet [34], VITAL [37],

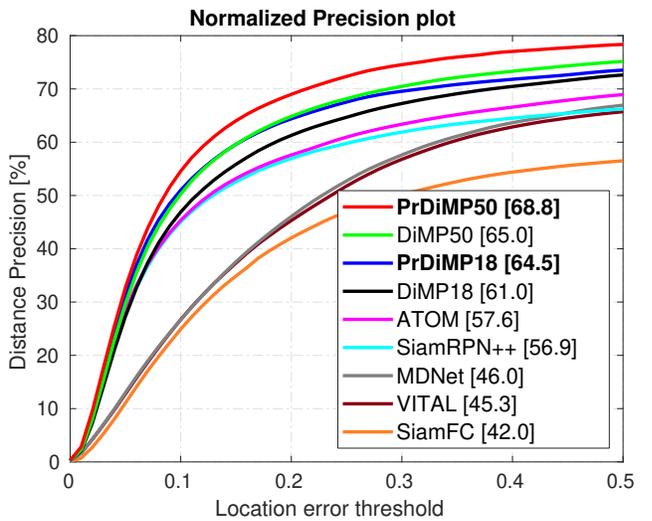


Figure 7. Normalized precision plot on the LaSOT dataset [13]. The average normalized precision is shown in the legend. Our approach outperforms previous trackers by a large margin.

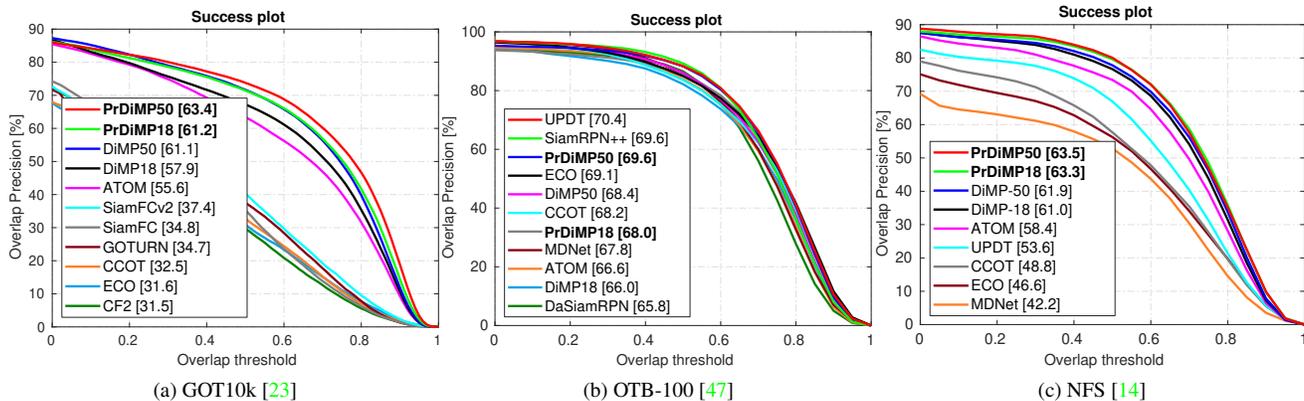


Figure 8. Success plots on the GOT10k [23] (a), OTB-100 [47] (b), and NFS [14] (c) datasets, showing the percentage of frames with a ground-truth IoU overlap larger than a threshold. The area-under-the-curve (AUC) metric for each tracker is shown in the legend. Our approach outperforms previous approaches by a significant margin on the challenging GOT10k and NFS datasets, while performing similarly to the top trackers on the largely saturated OTB-100 dataset.

and SiamFC [2]. Our approach outperforms previous state-of-the-art by a large margin. Compared to the ResNet-50 based SiamRPN++ [28] and DiMP-50 [3], our PrDiIMP50 version achieves absolute gains of 11.9% and 3.8% respectively. As demonstrated by the plot, this gain is obtained by an improvement in both accuracy (small thresholds) and robustness (large thresholds).

C.2. GOT10k

The success plot over the 180 videos in the GOT10k [23] test set is shown in Figure 8a. It displays the overlap precision OP_T as a function of the IoU threshold T . Overlap precision OP_T itself is defined as the percentage of frames with a ground-truth IoU overlap larger than a threshold T . The final area-under-the-curve (AUC) score is shown in the legend. We compare our approach with trackers with available results: DiMP [3], ATOM [9], SiamFCv2 (CFNet) [44], SiamFC [2], GOTURN [21], CCOT [12], ECO [8], and CF2 (HCF) [31]. Our PrDiIMP versions outperform previous methods, in particular for large overlap thresholds. This demonstrates the superior accuracy of our probabilistic bounding box regression formulation.

C.3. OTB-100

We provide the success plot over the 100 videos in the OTB dataset [47] in Figure 8b. We compare with state-of-the-art trackers UPDT [4], SiamRPN++ [28], ECO [8], DiMP [3], CCOT [12], MDNet [34], ATOM [9], and DaSiamRPN [52]. Despite the highly saturated nature of this dataset, our tracker performs among the top methods, providing a significant gain over the baseline DiMP.

C.4. NFS

Figure 8c show the success plot over the 30 FPS version of the challenging NFS dataset [14]. We compare with top trackers with available results: DiMP [3], ATOM [9],

UPDT [4], CCOT [12], ECO [8], and MDNet [34]. In this case, both our ResNet-18 and 50 versions achieve similar results, significantly outperforming the previous state-of-the-art DiMP-50.

D. Visualization of Predicted Distributions

In this section, we provide additional visualizations of the predicted probability distributions for target center and bounding box regression branches in our tracker. We visualize the output as in Figure 4. Several example frames for two challenging sequences are shown in Figure 9. As during standard tracking, the predicted distribution $p(y^c|x, \theta)$ for the target center regression (second column) is computed by applying the fully convolutional center regression branch on the search region centered at the previous target location.

To visualize the probability distribution $p(y^{bb}|x, \theta)$ predicted by the bounding box regression branch, we evaluate the density in a grid. Note that the bounding box $y^{bb} \in \mathbb{R}^4$ is 4-dimensional, and we therefore cannot visualize the full distribution as a 2-dimensional heatmap. We therefore plot two *slices* of this density, showing the variation in the bounding box location and size as follows. Our bounding box is parametrized as $y^{bb} = (c_x/w_0, c_y/h_0, \log w, \log h)$, where (c_x, c_y) is the center position, (w, h) is the size (width and height), and (w_0, h_0) is a constant reference size. The latter is set to the current estimate of the target size. The distribution of the bounding box center (third column in Figure 9) is obtained by predicting the density value $p(y^{bb}|x, \theta) \sim \exp(s_\theta^{bb}(y^{bb}, x))$ in a dense grid of center coordinates (c_x, c_y) , while keeping the size (w, h) constant at the current target estimate. To visualize the distribution over the bounding box size (fourth column), we conversely evaluate $p(y^{bb}|x, \theta) \sim \exp(s_\theta^{bb}(y^{bb}, x))$ in a dense grid of log-size coordinates $(\log w, \log h)$, while keeping the box centered at the current target estimate (c_x, c_y) .

In Figure 9, the target center density is visualized in

the range $\pm 2.5\sqrt{wh}$ relative to the previous target location. The bounding box center density is plotted within the range $c_x \pm w$ and $c_y \pm h$. The distribution over the bounding box size is plotted from $1/3$ to 3 times the estimated target size (w, h) , *i.e.* $\pm \log 3$. Outputs for two challenging sequences are visualized in Figure 9. The left part shows a cat and its mirror image. Due to their similarity and proximity, it is in many frames difficult to predict the exact bounding box of the cat. In these cases, the predicted distribution captures this uncertainty. For example, in the fourth row, two clear modes are predicted, which correspond to aligning the box with the real cat or with the right edge of the reflection. Moreover, the last row shows a failure case, where the box briefly expands. However, note that the size probability distribution (right column) is highly uncertain. This information could thus be used to indicate low reliability of the

estimated bounding box size.

The right part of Figure 9 depicts a challenging sequence with multiple distractors. The target center regression, which has a wider view of the scene, captures the presence of distractors in uncertain cases. In the last row, the tracker briefly fails by jumping to a distractor object. However, there is a strong secondary mode in the target center regression distribution that indicates the true target. Thus, the estimated distribution accurately captures the uncertainty in this ambiguous case. Moreover, in row 4-7, the target object is small with many nearby similar-looking objects. This makes bounding box regression extremely hard, even for a human. Our network can predict flexible distributions reflecting meaningful uncertainties for both bounding box position and size, when encountered with these difficulties.

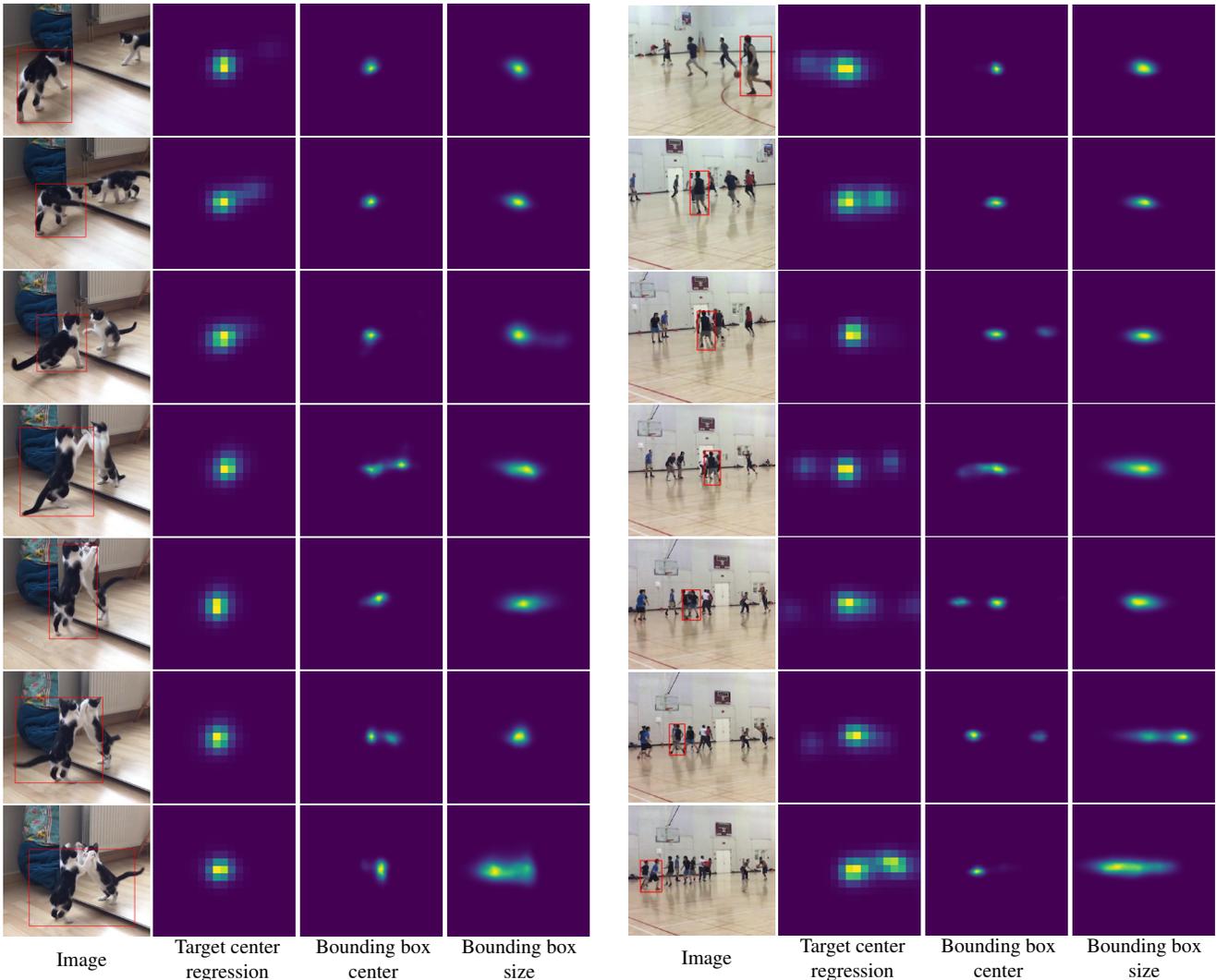


Figure 9. Visualization of the probability densities $p(y^{tc}|x, \theta)$ and $p(y^{bb}|x, \theta)$ predicted by the target center and bounding box regression branch respectively. We illustrate the output for example frames in two highly challenging sequences, where capturing uncertainty is important. Refer to the text for details.