

# Siamese Box Adaptive Network for Visual Tracking

Zedu Chen<sup>1</sup>, Bineng Zhong<sup>1,6\*</sup>, Guorong Li<sup>2</sup>, Shengping Zhang<sup>3,4</sup>, Rongrong Ji<sup>5,4</sup>

<sup>1</sup>Department of Computer Science and Technology, Huaqiao University

<sup>2</sup>School of Computer Science and Technology, University of Chinese Academy of Sciences

<sup>3</sup>Harbin Institute of Technology, <sup>4</sup>Peng Cheng Laboratory

<sup>5</sup>Department of Artificial Intelligence, School of Informatics, Xiamen University

<sup>6</sup>Key Laboratory of Intelligent Perception and Systems for High-Dimensional Information of Ministry of Education, Nanjing University of Science and Technology

zeduchen@stu.hqu.edu.cn, bnzhong@hqu.edu.cn, liguorong@ucas.ac.cn  
s.zhang@hit.edu.cn, rrji@xmu.edu.cn

## Abstract

Most of the existing trackers usually rely on either a multi-scale searching scheme or pre-defined anchor boxes to accurately estimate the scale and aspect ratio of a target. Unfortunately, they typically call for tedious and heuristic configurations. To address this issue, we propose a simple yet effective visual tracking framework (named Siamese Box Adaptive Network, SiamBAN) by exploiting the expressive power of the fully convolutional network (FCN). SiamBAN views the visual tracking problem as a parallel classification and regression problem, and thus directly classifies objects and regresses their bounding boxes in a unified FCN. The no-prior box design avoids hyper-parameters associated with the candidate boxes, making SiamBAN more flexible and general. Extensive experiments on visual tracking benchmarks including VOT2018, VOT2019, OTB100, NFS, UAV123, and LaSOT demonstrate that SiamBAN achieves state-of-the-art performance and runs at 40 FPS, confirming its effectiveness and efficiency. The code will be available at <https://github.com/hqucv/siamban>.

## 1. Introduction

Visual tracking is a fundamental but challenging task in computer vision. Given the target state in the initial frame of a sequence, the tracker needs to predict the target state in each subsequent frame. Despite great progress in recent years, visual tracking still faces challenges due to occlusion, scale variation, background clutters, fast motion, illumination variation, and appearance variations.

In a real-world video, the target scale and aspect ratio are

\*Corresponding author.

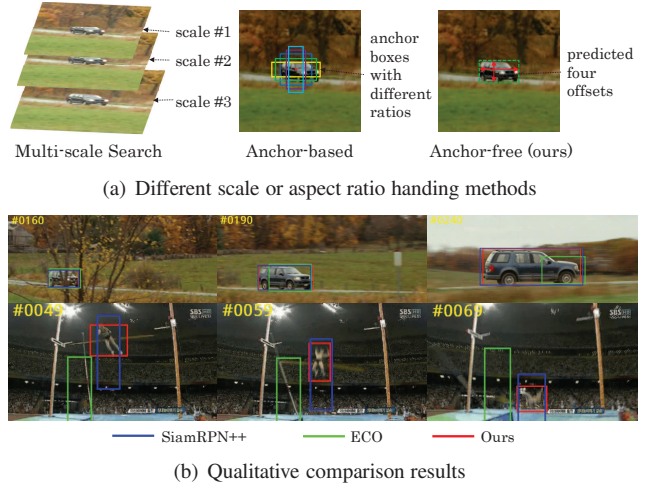


Figure 1. (a) Methods used to estimate the target scale or aspect ratio: multi-scale search (such as SiamFC, ECO), anchor-based (such as SiamRPN, SiamRPN++), and anchor-free (such as ours). (b) Some representative experiment results from our SiamBAN tracker and two state-of-the-art trackers. Observed from the visualization results, our tracker is better than the other two trackers in terms of scale and aspect ratio.

also changing due to target or camera movement and target appearance changes. Accurately estimating the scale and aspect ratio of the target becomes a challenge in the field of visual tracking. However, many existing trackers ignore this problem and rely on a multi-scale search to estimate the target size. For example, the current state-of-the-art correlation filter based trackers [6, 3] rely on their classification components, and the target scale is simply estimated by multi-scale search. Recently, Siamese network based visual trackers [21, 52, 20] introduce a region proposal network

(RPN) to obtain accurate target bounding boxes. However, in order to handle different scales and aspect ratios, they need to carefully design anchor boxes based on heuristic knowledge, which introduces many hyper-parameters and computational complexity.

In contrast, neuroscientists have shown that the bio-visual primary visual cortex can quickly and effectively extract the contours or boundaries of the observed objects from complex environments [29]. That is to say, humans can identify the object position and boundary without candidate boxes. So can we design an accurate and robust visual tracking framework without relying on candidate boxes? Inspired by the anchor-free detectors [14, 47, 31, 51, 37], the answer is yes. By exploiting the expressive power of the fully convolutional network (FCN), we propose a simple yet effective visual tracking framework named Siamese box adaptive network (SiamBAN) to address the challenge of accurately estimating the scale and aspect ratio of the target. The framework consists of a Siamese network backbone and multiple box adaptive heads, which does not require pre-defined candidate boxes and can be optimized end-to-end during training. SiamBAN classifies the target and regresses bounding boxes directly in a unified FCN, transforming the tracking problem into a classification-regression problem. Specifically, it directly predicts the foreground-background category score and a 4D vector of each spatial position on the correlation feature maps. The 4D vector depicts the relative offset from the four sides of the bounding box to the center point of the feature location corresponding to the search region. During inference, we use a search image centered on the previous position of the target. Through the bounding box corresponding to the position of the best score, we can get the displacement and size change of the target between frames.

The main contributions of this work are threefold.

- We design a Siamese box adaptive network, which can perform end-to-end offline training with deep convolutional neural networks [12] on well-annotated datasets [34, 30, 25, 15, 9].
- The no-prior box design in SiamBAN avoids hyper-parameters associated with the candidate boxes, making our tracker more flexible and general.
- The proposed SiamBAN not only achieves state-of-the-art results, but also runs at 40 FPS on tracking benchmarks including VOT2018 [17], VOT2019 [18], OTB100 [43], NFS [16], UAV123 [27] and LaSOT [9].

## 2. Related Works

Visual tracking is one of the most active research topics in computer vision in recent decades. A comprehensive survey of the related trackers is beyond the scope of this paper,

so we only briefly review two aspects that are most relevant to our work: Siamese network based visual trackers and anchor-free object detectors.

### 2.1. Siamese Network Based Visual Trackers

Recently, Siamese network based trackers have attracted great attention from the visual tracking community due to their end-to-end training capabilities and high efficiency [1, 11, 41, 21, 20, 49]. SiamFC [1] adopts the Siamese network as a feature extractor and first introduces the correlation layer to combine feature maps. Owing to its light structure and no need to model update, SiamFC runs efficiently at 86 FPS. DSiam [11] learns a feature transformation to handle the target appearance variation and to suppress background. RASNet [41] embeds diverse attention mechanisms in the Siamese network to adapt the tracking model to the current target. However, these methods need a multi-scale test to cope with scale variation and cannot handle aspect ratio changes due to target appearance variations. In order to get a more accurate target bounding box, SiamRPN [21] introduces the RPN [32] into the SiamFC. SPM-Tracker [39] proposes a series-parallel matching framework to enhance the robustness and discrimination power of SiamRPN. SiamRPN++ [20], SiamMask [42] and SiamDW [49] remove the influence factors such as padding in different ways, and introduce modern deep neural networks such as ResNet [12], ResNeXt [44] and MobileNet [13] into the Siamese network based visual trackers. Although anchor-based trackers [21, 39, 20] can handle changes in scale and aspect ratio, it is necessary to carefully design and fix the parameters of the anchor boxes. Design parameters often requires heuristic adjustments and involves many tricks to achieve good performance. In contrast to anchor-based trackers, our tracker avoids hyper-parameters associated with the anchor boxes and is more flexible and general.

### 2.2. Anchor-free Object Detectors

Recently, anchor-free object detection has attracted the attention of the object detection community. However, anchor-free detection is not a new concept. DenseBox [14] first introduced an FCN framework to jointly perform face detection and landmark localization. UnitBox [47] offered another option for performance improvement by carefully designing optimization losses. YOLOv1 [31] proposed to divide the input image into a grid and then predicted bounding boxes and class probabilities on each grid cell.

Recently, many new anchor-free detectors have emerged. These detection methods can be roughly classified into key-point based object detection [19, 50, 46] and dense detection [51, 37]. Specifically, CornerNet [19] proposed to detect an object bounding box as a pair of keypoints. ExtremeNet [50] presented to detect four extreme points and one center

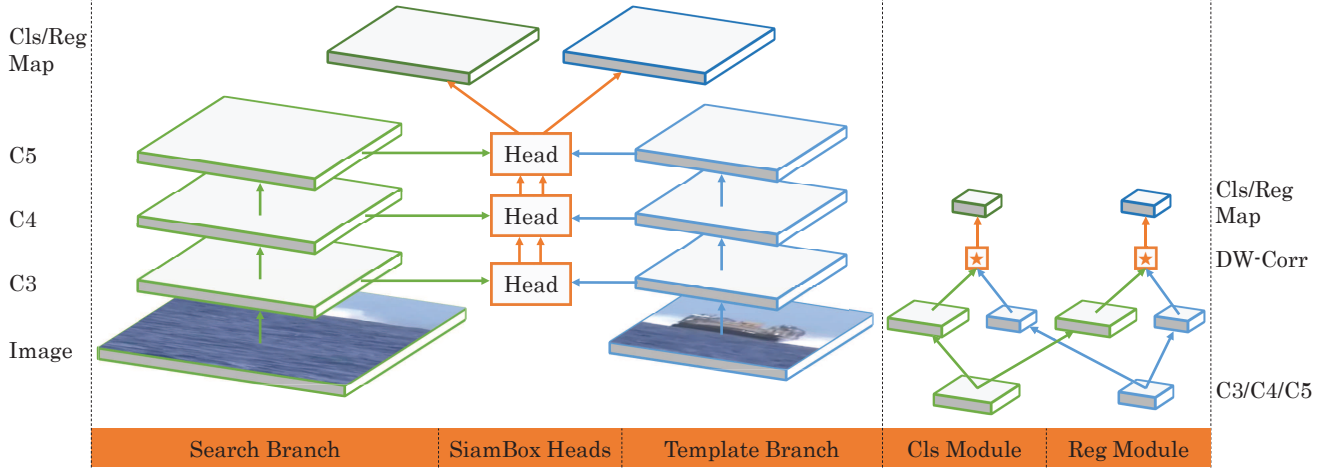


Figure 2. The framework of the proposed Siamese box adaptive network. The left sub-figure shows its main structure, where C3, C4, and C5 denote the feature maps of the backbone network, Cls Map and Reg Map denote the feature maps of the SiamBAN heads output. The right sub-figure shows each SiamBAN head, where DW-Corr means depth-wise cross-correlation operation.

point of objects using a standard keypoint estimation network. RepPoints [46] introduced the representative points, a new representation of objects to model fine-grained localization information and identify local areas significant for object classification. FSAF [51] proposed feature selective anchor-free module to address the limitations imposed by heuristic feature selection for anchor-based single-shot detectors with feature pyramids. FCOS [37] proposed to directly predict the possibility of object existence and the bounding box coordinates without anchor reference.

Compared to object detection, there are two key challenges in the visual tracking task, i.e. unknown categories and discrimination between different objects. The anchor-free detectors usually assume the categories of the objects to be detected are pre-defined. However, the categories of the targets are unknown before tracking. Meanwhile, anchor-free detectors typically focus on detecting the objects from different categories, while in tracking, it is necessary to determine whether the two objects are the same one. Therefore, a template branch that can encode the appearance information is needed in our framework to identify the target and background.

### 3. SiamBAN Framework

In this section, we describe the proposed SiamBAN framework. As shown in Figure 2, SiamBAN consists of a Siamese network backbone and multiple box adaptive heads. The Siamese network backbone is responsible for computing the convolutional feature maps of the template patch and the search patch, which uses an off-the-shelf convolutional network. The box adaptive head includes a classification module and a regression module. Specifically,

the classification module performs foreground-background classification on each point of the correlation layer, and the regression module performs bounding box prediction on the corresponding position.

#### 3.1. Siamese Network Backbone

Modern deep neural networks [12, 44, 13] have proven to be effective in Siamese network based trackers [20, 42, 49], and now we can use them such as ResNet, ResNeXt, and MobileNet in Siamese network based trackers. In our tracker, we adopt ResNet-50 [12] as the backbone network. Although ResNet-50 with continuous convolution striding can learn more and more abstract feature representations, it reduces feature resolution. However, Siamese network based trackers need detailed spatial information to perform dense predictions. To deal with this problem, we remove the downsampling operations from the last two convolution blocks. At the same time, in order to improve the receptive field, we use atrous convolution [4], which is proven to be effective for visual tracking [21, 42]. In addition, inspired by multi-grid methods [40], we adopt different atrous rates in our model. Specifically, we set the stride to 1 in the *conv4* and *conv5* blocks, the atrous rate to 2 in the *conv4* block, and the atrous rate to 4 in the *conv5* block.

The Siamese network backbone consists of two identical branches. One is called the template branch, which receives the template patch as input (denoted as  $z$ ). The other is called the search branch, which receives the search patch as input (denoted as  $x$ ). The two branches share parameters in a convolutional neural network to ensure that the same transformation is applied to both patches. In order to reduce the computational burden, we add a  $1 \times 1$  convolution to reduce the output features channel to 256, and use only the

features of the template branch center  $7 \times 7$  regions [38, 20], which can still capture the entire target region. For convenience, the output features of the Siamese network are represented as  $\varphi(z)$  and  $\varphi(x)$ .

### 3.2. Box Adaptive Head

As shown in Figure 2 (right), box adaptive head consists of a classification module and a regression module. Both modules receive features from the template branch and the search branch. So we adjust and copy  $\varphi(z)$  and  $\varphi(x)$  to  $[\varphi(z)]_{cls}$ ,  $[\varphi(z)]_{reg}$  and  $[\varphi(x)]_{cls}$ ,  $[\varphi(x)]_{reg}$  to the corresponding module. According to our design, each point of the correlation layer of the classification module needs to output two channels for foreground-background classification, and each point of the correlation layer of the regression module needs to output four channels for prediction of the bounding box. Each module combines the feature maps using a depth-wise cross-correlation layer [20]:

$$\begin{aligned} P_{w \times h \times 2}^{cls} &= [\varphi(x)]_{cls} \star [\varphi(z)]_{cls}, \\ P_{w \times h \times 4}^{reg} &= [\varphi(x)]_{reg} \star [\varphi(z)]_{reg}, \end{aligned} \quad (1)$$

where  $\star$  denotes the convolution operation with  $[\varphi(z)]_{cls}$  or  $[\varphi(z)]_{reg}$  as the convolution kernel,  $P_{w \times h \times 2}^{cls}$  denotes classification map,  $P_{w \times h \times 4}^{reg}$  indicates regression map. It is worth noting that our tracker outputs 5 times fewer variables than anchor-based trackers [21, 20] with 5 anchor boxes.

For each location on the classification map  $P_{w \times h \times 2}^{cls}$  or the regression map  $P_{w \times h \times 4}^{reg}$ , we can map it to the input search patch. For example, the location  $(i, j)$  corresponding to the location on the search patch is  $[w_{im} - (\lfloor \frac{w}{2} \rfloor - i) \times s, h_{im} - (\lfloor \frac{h}{2} \rfloor - j) \times s]$  (denoted as  $(p_i, p_j)$ ).  $w_{im}$  and  $h_{im}$  represent the width and height of the input search patch and  $s$  represents the total stride of the network, which is the center of the receptive field of the position  $(i, j)$ . For the regression, the anchor-based trackers [21, 52, 20] treat the location  $(p_i, p_j)$  as the center of the anchor box and regress the location  $(p_i, p_j)$ , width  $a_w$  and height  $a_h$ . That is, for the position  $(i, j)$ , the regression can adjust all of its offset values, but the classification is still performed in the original position, which may result in inconsistencies in classification and regression. So we do not adjust the location  $(p_i, p_j)$  and only calculate its offset value to the bounding box. In addition, since our regression targets are positive real numbers, we apply  $\exp(x)$  at the last level of the regression module to map any real number to  $(0, +\infty)$ .

### 3.3. Multi-level Prediction

After utilizing ResNet-50 with atrous convolution, we can use multi-level features for prediction. Although the spatial resolutions of the *conv3*, *conv4* and *conv5* blocks of our backbone network are the same, they have atrous convolutions with different expansion rates, so the difference

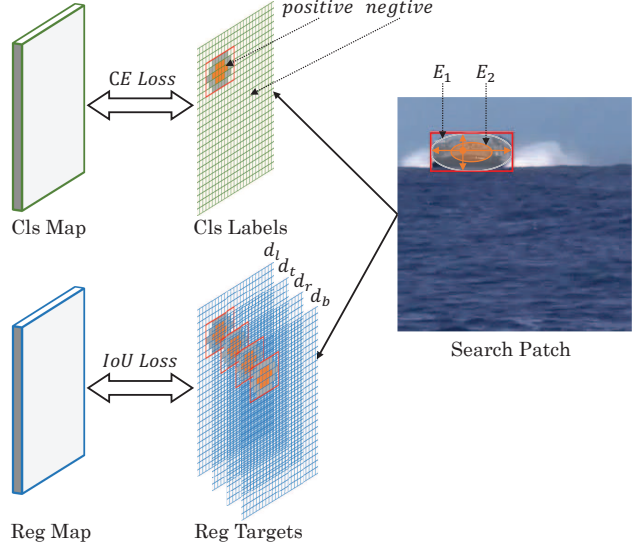


Figure 3. Illustrations of classification labels and regression targets. Prediction values and supervision signals are as shown, where  $E_1$  represents ellipse  $E_1$  and  $E_2$  represents ellipse  $E_2$ . We use a cross entropy and an IoU loss for classification and box regression, respectively.

between their receptive fields is large, and the captured information is naturally different. As pointed by CF2 [26], features from earlier layers can capture fine-grained information, which is useful for precise localization; while features from latter layers can encode abstract semantic information, which is robust to target appearance changes. In order to take full advantage of different characteristics of multi-level features, we use multiple box adaptive heads for prediction. The classification maps and the regression maps obtained by each detection head are adaptively fused:

$$\begin{aligned} P_{w \times h \times 2}^{cls-all} &= \sum_{l=3}^5 \alpha_l P_l^{cls}, \\ P_{w \times h \times 4}^{reg-all} &= \sum_{l=3}^5 \beta_l P_l^{reg}, \end{aligned} \quad (2)$$

where  $\alpha_l$  and  $\beta_l$  are the weights corresponding to each map and are optimized together with the network. By combining the classification map and the regression map independently, the classification module and the regression module can focus on the domains they need.

### 3.4. Ground-truth and Loss

**Classification Labels and Regression Targets.** As shown in Figure 3, the target on each search patch is marked with a ground-truth bounding box. The width, height, top-left corner, center point and bottom-right corner of the ground-truth bounding box are represented by  $g_w, g_h, (g_{x1}, g_{y1})$ ,



$(g_{x_c}, g_{y_c})$  and  $(g_{x_2}, g_{y_2})$ , respectively. With  $(g_{x_c}, g_{y_c})$  as the center and  $\frac{q_w}{2}, \frac{q_h}{2}$  as the axes length, we can get the ellipse  $E_1$ :

$$\frac{(p_i - g_{x_c})^2}{(\frac{q_w}{2})^2} + \frac{(p_j - g_{y_c})^2}{(\frac{q_h}{2})^2} = 1. \quad (3)$$

With  $(g_{x_c}, g_{y_c})$  as the center and  $\frac{q_w}{4}, \frac{q_h}{4}$  as the axes length, we can get the ellipse  $E_2$ :

$$\frac{(p_i - g_{x_c})^2}{(\frac{q_w}{4})^2} + \frac{(p_j - g_{y_c})^2}{(\frac{q_h}{4})^2} = 1. \quad (4)$$

If the location  $(p_i, p_j)$  falls within the ellipse  $E_2$ , it is assigned with a positive label, and if it falls outside the ellipse  $E_1$ , it is assigned with a negative label, and if it falls between the ellipses  $E_2$  and  $E_1$ , ignore it. The location  $(p_i, p_j)$  with a positive label is used to regress the bounding box, and the regression targets can be formulated as:

$$\begin{aligned} d_l &= p_i - g_{x_1}, \\ d_t &= p_j - g_{y_1}, \\ d_r &= g_{x_2} - p_i, \\ d_b &= g_{y_2} - p_j, \end{aligned} \quad (5)$$

where  $d_l, d_t, d_r, d_b$  are the distances from the location to the four sides of the bounding box, as shown in Figure 3.

**Classification Loss and Regression Loss.** We define our multi-task loss function as follows:

$$L = \lambda_1 L_{cls} + \lambda_2 L_{reg}, \quad (6)$$

where  $L_{cls}$  is the cross entropy loss,  $L_{reg}$  is the IoU (Intersection over Union) Loss. We do not search for the hyper-parameters of Eq.6, and simply set  $\lambda_1 = \lambda_2 = 1$ . Similar to GIoU [33], we define IoU loss as:

$$L_{IoU} = 1 - IoU, \quad (7)$$

where  $IoU$  represents the area ratio of intersection to union of the predicted bounding box and the ground-truth bounding box. The location  $(p_i, p_j)$  with a positive label is within the ellipse  $E_2$  and the regression value is greater than 0, so  $0 < IoU \leq 1$ , then  $0 \leq L_{IoU} < 1$ . The IoU loss can make  $d_l, d_t, d_r, d_b$  jointly be regressed.

### 3.5. Training and Inference

**Training.** Our entire network can be trained end-to-end on large-scale datasets. We train SiamBAN with image pairs sampled on videos or still images. The training sets include ImageNet VID [34], YouTube-BoundingBoxes [30], COCO [25], ImageNet DET [34], GOT10k [15] and LaSOT [9]. The size of a template patch is  $127 \times 127$  pixels, while the size of a search patch is  $255 \times 255$  pixels. Also, although our negative samples are much less than anchor-based trackers

[21, 20], negative samples are still much more than positive samples. Therefore we collect at most 16 positive samples and 48 negative samples from one image pair.

**Inference.** During inference, we crop the template patch from the first frame and feed it to the feature extraction network. The extracted template features are cached, so we do not have to calculate them in subsequent tracking. For subsequent frames, we crop the search patch and extract feature based on the target position of the previous frame, and then perform prediction in the search region to get the total classification map  $P_{w \times h \times 2}^{cls-all}$  and regression map  $P_{w \times h \times 2}^{reg-all}$ . Afterward, we can get prediction boxes by the following equation:

$$\begin{aligned} p_{x_1} &= p_i - d_l^{reg}, \\ p_{y_1} &= p_j - d_t^{reg}, \\ p_{x_2} &= p_i + d_r^{reg}, \\ p_{y_2} &= p_j + d_b^{reg}, \end{aligned} \quad (8)$$

where  $d_l^{reg}, d_t^{reg}, d_r^{reg}$  and  $d_b^{reg}$  denote the prediction values of the regression map,  $(p_{x_1}, p_{y_1})$  and  $(p_{x_2}, p_{y_2})$  are the top-left corner and bottom-right corner of the prediction box.

After prediction boxes are generated, we use the cosine window and scale change penalty to smooth target movements and changes [21], then the prediction box with the best score is selected and its size is updated by linear interpolation with the state in the previous frame.

## 4. Experiments

### 4.1. Implementation Details

We initialize our backbone networks with the weights pre-trained on ImageNet [34] and the parameters of the first two layers are frozen. Our network is trained with stochastic gradient descent (SGD) with a minibatch of 28 pairs. We train a total of 20 epochs, using a warmup learning rate of 0.001 to 0.005 in the first 5 epochs and a learning rate exponentially decayed from 0.005 to 0.00005 in the last 15 epochs. In the first 10 epochs, we only train the box adaptive heads, and in the last 10 epochs fine-tuned the backbone network with one-tenth of the current learning rate. Weight decay and momentum are set as 0.0001 and 0.9. Our approach is implemented in Python using PyTorch on a PC with Intel Xeon(R) 4108 1.8GHz CPU, 64G RAM, Nvidia GTX 1080Ti.

### 4.2. Comparison with State-of-the-art Trackers

We compare our **SiamBAN** tracker with the state-of-the-art trackers on six tracking benchmarks. Our tracker achieves state-of-the-art results and run at 40 FPS.

**VOT2018 [17].** We evaluate our tracker on the Visual Object Tracking challenge 2018 (VOT2018) consisting of 60 sequences. The overall performance of the tracker

	DRT [36]	RCO [17]	UPDT [3]	SiamRPN [21]	MFT [17]	LADCF [45]	ATOM [5]	SiamRPN++ [20]	DiMP [2]	Ours
EAO↑	0.355	0.376	0.379	0.384	0.386	0.389	0.401	0.417	0.441	0.452
Accuracy↑	0.518	0.507	0.536	0.588	0.505	0.503	0.590	0.604	0.597	0.597
Robustness↓	0.201	0.155	0.184	0.276	0.140	0.159	0.203	0.234	0.152	0.178

Table 1. Detailed comparisons on VOT2018. The best two results are highlighted in red and blue fonts. DiMP is the ResNet-50 version (DiMP-50), the same below.

	SA_SIAM_R [18]	SiamCRF_RT [18]	SPM [39]	SiamRPN++ [20]	SiamMask [42]	ARTCS [18]	SiamDW_ST [49]	DCFST [18]	DiMP [2]	Ours
EAO↑	0.252	0.262	0.275	0.285	0.287	0.287	0.299	0.317	0.321	0.327
Accuracy↑	0.563	0.549	0.577	0.599	0.594	0.602	0.600	0.585	0.582	0.602
Robustness↓	0.507	0.346	0.507	0.482	0.461	0.482	0.467	0.376	0.371	0.396

Table 2. Detailed comparisons on VOT2019 real-time experiments.

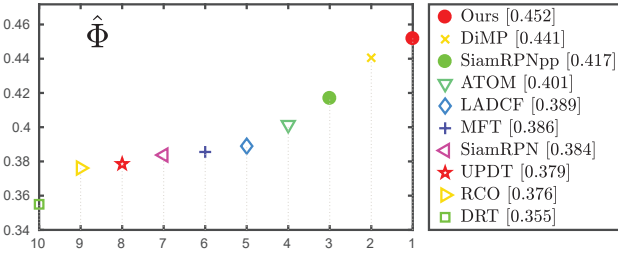


Figure 4. Expected averaged overlap performance on VOT2018. SiamRPNpp is SiamRPN++, the same below.

is evaluated using the EAO (Expected Average Overlap), which combines accuracy (average overlap during successful tracking periods) and robustness (failure rate). Table 1 shows the comparison with almost all the top-performing trackers in the VOT2018. Among previous approaches, DiMP [2] achieves the best EAO and SiamRPN++ [20] achieves the best accuracy, they all use ResNet-50 to extract feature. DiMP has the same accuracy as our tracker, and although its failure rate is slightly lower than ours, our EAO is slightly better, without any online update. Compared with SiamRPN++, our tracker achieves similar accuracy, but the failure rate decreases by 23.9% and EAO increases by 8.4%. Among these trackers, our tracker has the highest EAO and ranks second in terms of accuracy. This shows that our tracker not only accurately estimates the target’s location but also maintain good robustness.

**Comparison of attributes on VOT2018.** All sequences of VOT2018 are per-frame annotated by the following visual attributes: camera motion, illumination change, occlusion, size change, and motion change. Frames that do not correspond to any of the five attributes are represented as unassigned. We compare the EAO of the visual attributes

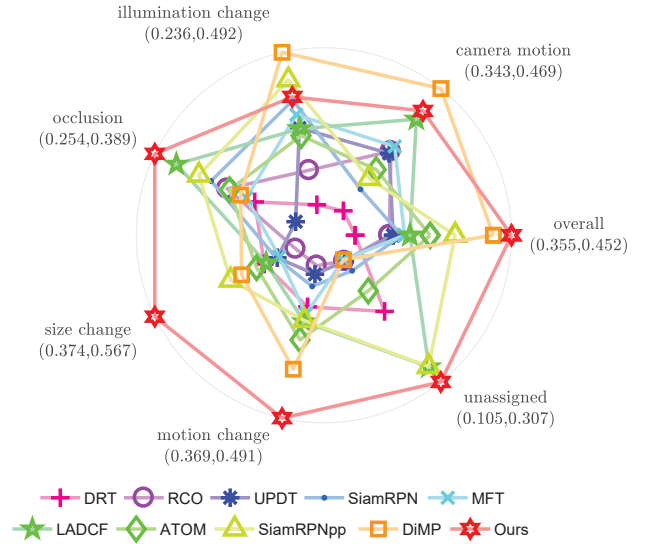


Figure 5. Comparison of EAO on VOT2018 for the following visual attributes: camera motion, illumination change, occlusion, size change and motion change. Frames that do not correspond to any of the five attributes are marked as unassigned. The values in parentheses indicate the EAO range of each attribute and overall of the trackers.

of the top-performing trackers. As shown in Figure 5, our tracker ranks first on attributes of occlusion, size change, and motion change, and ranks second and third on attributes of camera motion and illumination. This shows that our tracker is robust to occlusion, size changes, and motion changes in the target while having the ability to cope with camera motion and illumination changes.

**VOT2019 [18].** We evaluate our tracker on Visual Object Tracking challenge 2019 (VOT2019) real-time ex-

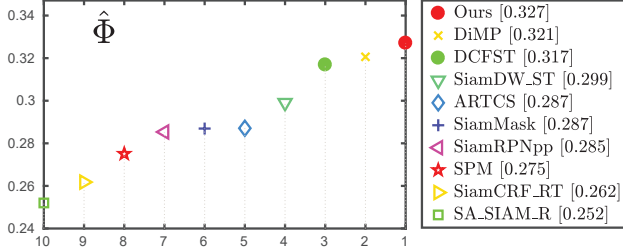


Figure 6. Expected averaged overlap performance on VOT2019.

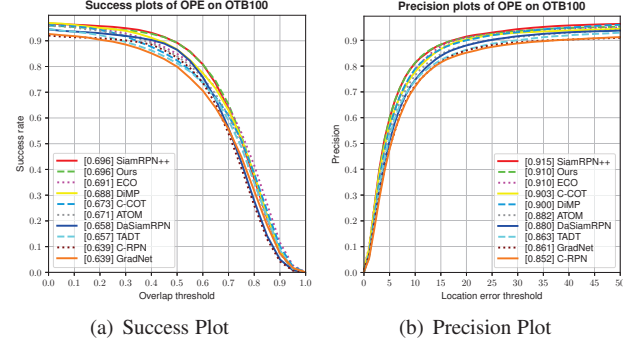


Figure 7. Success and precision plots on OTB100.

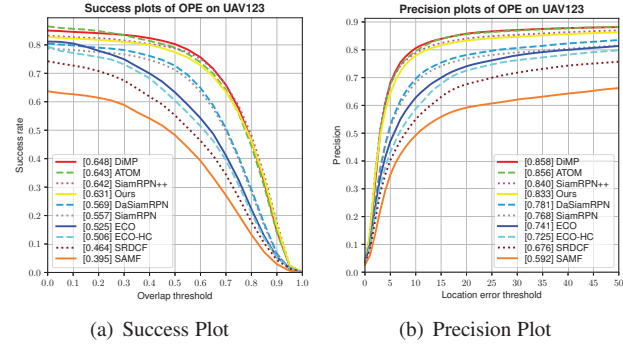


Figure 8. Success and precision plots on UAV123.

periments. The VOT2019 sequences were replaced by 20% compared to the VOT2018. Table 2 shows the results presented in terms of EAO, robustness, and accuracy. SiamMargin [18] achieves a lower failure rate through on-line updates, but our accuracy is higher than it. Although SiamRPN++ achieves similar accuracy to our tracker, our failure rate is 17.8% lower than it and achieves 14.7% relative gain in EAO. Among these trackers, our tracker has the highest accuracy and EAO. This shows that our method can accurately estimate the target bounding box.

**OTB100 [43].** OTB100 is a widely used public tracking benchmark consisting of 100 sequences. Our SiamBAN tracker is compared with numerous state-of-the-art trackers including SiamRPN++ [20], ECO [6], DiMP [2], C-COT

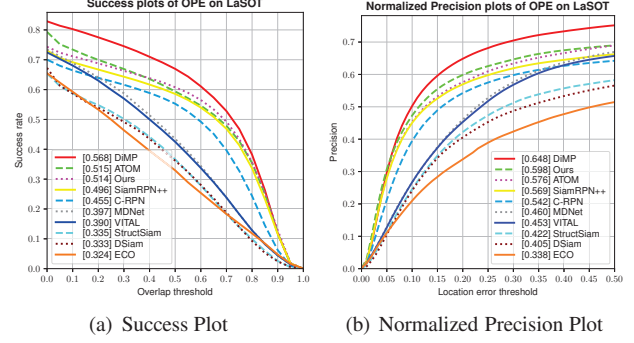


Figure 9. Success and normalized precision plots on LaSOT.

	MDNet	ECO	C-COT	UPDT	ATOM	DiMP	Ours
	[28]	[6]	[8]	[3]	[5]	[2]	
AUC↑	0.422	0.466	0.488	0.537	0.584	0.620	0.594

Table 3. Comparison with State-of-the-art trackers on the NFS dataset in terms of AUC.

[8], ATOM [5], DaSiamRPN [52], TADT [23], C-RPN [10], GradNet [22]. Figure 7 illustrates the success and precision plots of the compared trackers. Prior to SiamRPN++, due to the limited representation capabilities of shallow networks, the Siamese network based [52] trackers achieves sub-optimal performance on the OTB100. After using ResNet-50 as the feature extraction network, SiamRPN++ achieves leading results. Compared to SiamRPN++, achieves similar results with a simpler design.

**NFS [16].** The NFS dataset consists of 100 videos (380K frames) captured from real-world scenes with higher frame rate cameras. We evaluate our tracker in the 30FPS version of the dataset, AUC are shown in Table 3. It can be seen that our tracker ranks second and improved by 40.8% compared to the best tracker in the NFS paper.

**UAV123 [27].** The UAV123 is a new aerial video benchmark and dataset, which contains 123 sequences captured from a low-altitude aerial perspective. The benchmarks can be used to assess whether the tracker is suitable for deployment to a UAV in real-time scenarios. We compare our tracker with other 9 state-of-art real-time trackers, including DiMP [2], ATOM [5], SiamRPN++ [20], DaSiamRPN [52], SiamRPN [21], ECO [6], ECO-HC [6], SRDCF [7], SAMF [24]. Figure 8 shows the success and precision plots. Our tracker achieves state-of-the-art score.

**LaSOT [9].** LaSOT is a high-quality, large-scale dataset with a total of 1,400 sequences. Compared to the previous dataset, LaSOT has longer sequences with an average sequence length of more than 2,500 frames. Each sequence has various challenges from the wild where the target may disappear and reappear in the view, which tests the ability of

the tracker to re-track the target. We evaluate our tracker on the test set consisting of 280 videos with trackers including DiMP [2], ATOM [5], SiamRPN++ [20], C-RPN [10], MD-Net [28], VITAL [35], StructSiam [48], DSiam[11], ECO [6]. The results including success plots and normalized precision plots are illustrated in Figure 9. Our tracker ranks third in terms of AUC, second in terms of normalized precision and 5.1% higher than SiamRPN++.

### 4.3. Ablation Study

**Discussion on Multi-level Prediction.** To explore the role of different level features and the effect of aggregation of multi-level features, we have performed an ablation study on multi-layer prediction. It can be found from Table 4 that when only single-layer feature are used, *conv4* performs best. Compared with the single-layer features, when using the aggregation of the two-layer features, the performance has been improved, and the performance of *conv4* and *conv5* aggregation is the best. After aggregating three layers of features, our tracker achieves the best results.

**Discussion on Sample Label Assignment.** The sample label assignment plays a key role in the performance of a tracker. However, many Siamese network based trackers [1, 38, 11] do not pay enough attention to it. For example, SiamFC considers the elements of the score map within the radius  $R$  of the center to be positive samples. The label assignment method only considers the center position of the target, ignoring the size of the target. Intuitively, the sample label assignment should be different for targets of different sizes and shapes. Therefore, our label assignment also takes into account the target scale and aspect ratio. It is worth noting that we also set the buffer to ignore the ambiguous samples. The specific is in Section 3.4.

To illustrate the advantages of our label assignment method, we conduct comparative experiments with the other two label assignments. As shown in Figure 10, for convenience, we refer to these three types of labels as ellipse labels, circle labels and rectangle labels. For fair comparison, we define circles  $C_1$ ,  $C_2$  and rectangles  $R_1$ ,  $R_2$  in a similar way to define ellipses  $E_1$ ,  $E_2$ . Specifically, with  $(g_{xc}, g_{yc})$  as the center and  $\frac{\sqrt{g_w \times g_h}}{2}$ ,  $\frac{\sqrt{g_w \times g_h}}{4}$  as the radius, we can get the circles  $C_1$  and  $C_2$ . The rectangle  $R_1$  is the same position and size as the ground-truth bounding box. The center of the rectangle  $R_2$  is  $(g_{xc}, g_{yc})$ , and the sides length is  $\frac{g_w}{2}$ ,  $\frac{g_h}{2}$ .

As shown in Table 4, under the same number of iterations and training dataset, SiamBAN performs better than SiamBAN with circle labels and SiamBAN with rectangle labels. We believe that the reason is that ellipse labels can more accurately label positive and negative samples than circular labels and rectangular labels so that the trained tracker can more accurately distinguish the foreground-background and is more robust.

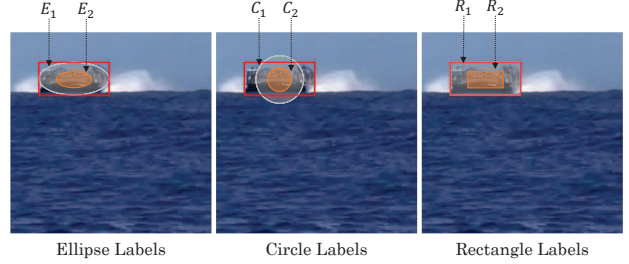


Figure 10. Three sample label assignment methods: ellipse labels, circle labels, rectangle labels.  $E_1$ ,  $E_2$ ,  $C_1$ ,  $C_2$ ,  $R_1$ ,  $R_2$  represent ellipse  $E_1$ , ellipse  $E_2$ , circle  $C_1$ , circle  $C_2$ , rectangle  $R_1$ , rectangle  $R_2$ , respectively.

L3	L4	L5	Circle	Rectangle	Ellipse	AUC
✓					✓	0.675
	✓				✓	0.683
		✓			✓	0.662
✓	✓				✓	0.687
✓		✓			✓	0.681
	✓	✓			✓	0.689
✓	✓	✓	✓			0.686
✓	✓	✓		✓		0.688
✓	✓	✓			✓	<b>0.696</b>

Table 4. Quantitative comparison results of our tracker and its variants with different detection heads and different label assignment methods on OTB100. L3, L4, L5 represent *conv3*, *conv4*, *conv5*, respectively. Circle, Rectangle, Ellipse represent circle labels, rectangle labels, ellipse labels, respectively.

## 5. Conclusions

In this paper, we exploit the expressive power of the fully convolutional network and propose a simple yet effective visual tracking framework named SiamBAN, which does not require a multi-scale searching schema and pre-defined candidate boxes. SiamBAN directly classifies objects and regresses bounding boxes in a unified network. Therefore, the visual tracking problem becomes a classification-regression problem. Extensive experiments on six visual tracking benchmarks demonstrate that SiamBAN achieves state-of-the-art performance and runs at 40 FPS, confirming its effectiveness and efficiency.

**Acknowledgments.** This work was supported by the National Natural Science Foundation of China (No. 61972167, 61772494, 61872112), the Fundamental Research Funds for the Central Universities (No. 30918014108), and the Open Project Program of the National Laboratory of Pattern Recognition (NLPR) (No. 202000012).



## References

- [1] Luca Bertinetto, Jack Valmadre, Joao F Henriques, Andrea Vedaldi, and Philip HS Torr. Fully-convolutional siamese networks for object tracking. In *European Conference on Computer Vision*, pages 850–865. Springer, 2016. 2, 8
- [2] Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte. Learning discriminative model prediction for tracking. In *IEEE International Conference on Computer Vision*, 2019. 6, 7, 8
- [3] Goutam Bhat, Joakim Johnander, Martin Danelljan, Fahad Shahbaz Khan, and Michael Felsberg. Unveiling the power of deep tracking. In *European Conference on Computer Vision*, pages 483–498, 2018. 1, 6, 7
- [4] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):834–848, 2017. 3
- [5] 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*, pages 4660–4669, 2019. 6, 7, 8
- [6] 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*, pages 6638–6646, 2017. 1, 7, 8
- [7] Martin Danelljan, Gustav Hager, Fahad Shahbaz Khan, and Michael Felsberg. Learning spatially regularized correlation filters for visual tracking. In *IEEE international conference on computer vision*, pages 4310–4318, 2015. 7
- [8] Martin Danelljan, Andreas Robinson, Fahad Shahbaz Khan, and Michael Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. In *European Conference on Computer Vision*, pages 472–488. Springer, 2016. 7
- [9] 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. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5374–5383, 2019. 2, 5, 7
- [10] Heng Fan and Haibin Ling. Siamese cascaded region proposal networks for real-time visual tracking. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 7952–7961, 2019. 7, 8
- [11] Qing Guo, Wei Feng, Ce Zhou, Rui Huang, Liang Wan, and Song Wang. Learning dynamic siamese network for visual object tracking. In *IEEE International Conference on Computer Vision*, pages 1763–1771, 2017. 2, 8
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 2, 3
- [13] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. MobileNets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. 2, 3
- [14] Lichao Huang, Yi Yang, Yafeng Deng, and Yinan Yu. DenseBox: Unifying landmark localization with end to end object detection. *arXiv preprint arXiv:1509.04874*, 2015. 2
- [15] 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. 2, 5
- [16] Hamed Kiani Galoogahi, Ashton Fagg, Chen Huang, Deva Ramanan, and Simon Lucey. Need for speed: A benchmark for higher frame rate object tracking. In *IEEE International Conference on Computer Vision*, pages 1125–1134, 2017. 2, 7
- [17] Matej Kristan, Ales Leonardis, Jiri Matas, Michael Felsberg, Roman Pflugfelder, Luka Cehovin Zajc, Tomas Vojir, Goutam Bhat, Alan Lukezic, Abdelrahman Eldesokey, et al. The sixth visual object tracking VOT2018 challenge results. In *European Conference on Computer Vision*, pages 0–0, 2018. 2, 5, 6
- [18] Matej Kristan, Jiri Matas, Ales Leonardis, Michael Felsberg, Roman Pflugfelder, Joni-Kristian Kamarainen, Luka Cehovin Zajc, Ondrej Drbohlav, Alan Lukezic, Amanda Berg, et al. The seventh visual object tracking VOT2019 challenge results. In *IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. 2, 6, 7
- [19] Hei Law and Jia Deng. CornerNet: Detecting objects as paired keypoints. In *European Conference on Computer Vision*, pages 734–750, 2018. 2
- [20] Bo Li, Wei Wu, Qiang Wang, Fangyi Zhang, Junliang Xing, and Junjie Yan. SiamRPN++: Evolution of siamese visual tracking with very deep networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4282–4291, 2019. 1, 2, 3, 4, 5, 6, 7, 8
- [21] Bo Li, Junjie Yan, Wei Wu, Zheng Zhu, and Xiaolin Hu. High performance visual tracking with siamese region proposal network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 8971–8980, 2018. 1, 2, 3, 4, 5, 6, 7
- [22] Peixia Li, Boyu Chen, Wanli Ouyang, Dong Wang, Xiaoyun Yang, and Huchuan Lu. GradNet: Gradient-guided network for visual object tracking. In *IEEE International Conference on Computer Vision*, pages 6162–6171, 2019. 7
- [23] Xin Li, Chao Ma, Baoyuan Wu, Zhenyu He, and Ming-Hsuan Yang. Target-aware deep tracking. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1369–1378, 2019. 7
- [24] Yang Li and Jianke Zhu. A scale adaptive kernel correlation filter tracker with feature integration. In *European Conference on Computer Vision*, pages 254–265. Springer, 2014. 7
- [25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014. 2, 5
- [26] Chao Ma, Jia-Bin Huang, Xiaokang Yang, and Ming-Hsuan Yang. Hierarchical convolutional features for visual tracking.

- In *IEEE international conference on computer vision*, pages 3074–3082, 2015. 4
- [27] Matthias Mueller, Neil Smith, and Bernard Ghanem. A benchmark and simulator for uav tracking. In *European Conference on Computer Vision*, pages 445–461. Springer, 2016. 2, 7
- [28] Hyeonseob Nam and Bohyung Han. Learning multi-domain convolutional neural networks for visual tracking. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4293–4302, 2016. 7, 8
- [29] Dale Purves, George J Augustine, David Fitzpatrick, William C Hall, Anthony-Samuel LaMantia, James O McNamara, and Leonard E White. *Neuroscience*. 4th, volume 857. Oxford University Press, 2008. 2
- [30] Esteban Real, Jonathon Shlens, Stefano Mazzocchi, Xin Pan, and Vincent Vanhoucke. YouTube-BoundingBoxes: A large high-precision human-annotated data set for object detection in video. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5296–5305, 2017. 2, 5
- [31] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 779–788, 2016. 2
- [32] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 2
- [33] Hamid Rezaatoughi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 658–666, 2019. 5
- [34] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 2, 5
- [35] Yibing Song, Chao Ma, Xiaohe Wu, Lijun Gong, Linchao Bao, Wangmeng Zuo, Chunhua Shen, Rynson WH Lau, and Ming-Hsuan Yang. VITAL: Visual tracking via adversarial learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 8990–8999, 2018. 8
- [36] Chong Sun, Dong Wang, Huchuan Lu, and Ming-Hsuan Yang. Correlation tracking via joint discrimination and reliability learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 489–497, 2018. 6
- [37] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. FCOS: Fully convolutional one-stage object detection. In *IEEE International Conference on Computer Vision*, pages 9627–9636, 2019. 2, 3
- [38] Jack Valmadre, Luca Bertinetto, João Henriques, Andrea Vedaldi, and Philip HS Torr. End-to-end representation learning for correlation filter based tracking. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2805–2813, 2017. 4, 8
- [39] Guangting Wang, Chong Luo, Zhiwei Xiong, and Wenjun Zeng. SPM-Tracker: Series-parallel matching for real-time visual object tracking. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3643–3652, 2019. 2, 6
- [40] Panqu Wang, Pengfei Chen, Ye Yuan, Ding Liu, Zehua Huang, Xiaodi Hou, and Garrison Cottrell. Understanding convolution for semantic segmentation. In *2018 IEEE winter conference on applications of computer vision*, pages 1451–1460. IEEE, 2018. 3
- [41] Qiang Wang, Zhu Teng, Junliang Xing, Jin Gao, Weiming Hu, and Stephen Maybank. Learning attentions: residual attentional siamese network for high performance online visual tracking. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4854–4863, 2018. 2
- [42] Qiang Wang, Li Zhang, Luca Bertinetto, Weiming Hu, and Philip HS Torr. Fast online object tracking and segmentation: A unifying approach. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1328–1338, 2019. 2, 3, 6
- [43] Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. Object tracking benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1834–1848, 2015. 2, 7
- [44] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1492–1500, 2017. 2, 3
- [45] 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 object tracking. *IEEE Transactions on Image Processing*, 2019. 6
- [46] Ze Yang, Shaohui Liu, Han Hu, Liwei Wang, and Stephen Lin. RepPoints: Point set representation for object detection. In *IEEE International Conference on Computer Vision*, pages 9657–9666, 2019. 2, 3
- [47] Jiahui Yu, Yuning Jiang, Zhangyang Wang, Zhimin Cao, and Thomas Huang. UnitBox: An advanced object detection network. In *24th ACM international conference on Multimedia*, pages 516–520. ACM, 2016. 2
- [48] Yunhua Zhang, Lijun Wang, Jinqing Qi, Dong Wang, Mengyang Feng, and Huchuan Lu. Structured siamese network for real-time visual tracking. In *European Conference on Computer Vision*, pages 351–366, 2018. 8
- [49] Zhipeng Zhang and Houwen Peng. Deeper and wider siamese networks for real-time visual tracking. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4591–4600, 2019. 2, 3, 6
- [50] Xingyi Zhou, Jiacheng Zhuo, and Philipp Krahenbuhl. Bottom-up object detection by grouping extreme and center points. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 850–859, 2019. 2
- [51] Chenchen Zhu, Yihui He, and Marios Savvides. Feature selective anchor-free module for single-shot object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 840–849, 2019. 2, 3
- [52] Zheng Zhu, Qiang Wang, Bo Li, Wei Wu, Junjie Yan, and Weiming Hu. Distractor-aware siamese networks for visual object tracking. In *European Conference on Computer Vision*, pages 101–117, 2018. 1, 4, 7