# DomainSiam: Domain-Aware Siamese Network for Visual Object Tracking

Mohamed H. Abdelpakey [(✉)] and Mohamed S. Shehata

Memorial University of Newfoundland,
St. John's, NL A1B 3X5, Canada
mha241@mun.ca

**Abstract.** Visual object tracking is a fundamental task in the field of computer vision. Recently, Siamese trackers have achieved *state-of-the-art* performance on recent benchmarks. However, Siamese trackers do not fully utilize semantic and objectness information from pre-trained networks that have been trained on the image classification task. Furthermore, the pre-trained Siamese architecture is sparsely activated by the category label which leads to unnecessary calculations and overfitting. In this paper, we propose to learn a Domain-Aware, that is fully utilizing semantic and objectness information while producing a class-agnostic using a ridge regression network. Moreover, to reduce the sparsity problem, we solve the ridge regression problem with a differentiable weighted-dynamic loss function. Our tracker, dubbed *DomainSiam*, improves the feature learning in the training phase and generalization capability to other domains. Extensive experiments are performed on five tracking benchmarks including OTB2013 and OTB2015 for a validation set; as well as the VOT2017, VOT2018, LaSOT, TrackingNet, and GOT10k for a testing set. *DomainSiam* achieves a *state-of-the-art* performance on these benchmarks while running at 53 FPS.

**Keywords:** Object tracking · Siamese network · Ridge regression network · Dynamic loss.

## 1 Introduction

Tracking is a fundamental task in many computer vision tasks such as surveillance [18], computer interactions [31] and image understanding [22]. The objective of tracking is to find the trajectory of the object of interest over time. This is a challenge since the object of interest undergoes appearance changes such as occlusions, motion blur, and background cluttering [44,1]. Recent deep trackers such as CFNet [40] and DeepSRDCF [10] use pre-trained networks which have been trained on image classification or object recognition.

In recent years, convolutional neural networks (CNNs) have achieved superior performance against hand-crafted trackers (e.g., CACF [32], SRDCF [11], KCF [16] and SAMF [27]). Siamese trackers such as SiamFC [4], CFNet [40], SiamRPN [23], and DensSiam [30] learn a similarity function to separate the foreground

from its background. However, Siamese trackers do not fully utilize the semantic and objectness information from pre-trained networks that have been trained on image classification. In image classification, the class categories of the objects are pre-defined while in object tracking task the tracker needs to be class-agnostic while benefiting from semantic and objectness information. Moreover, the image classification increases the inter-class differences while forcing the features to be insensitive to intra-class changes [26].

In this paper, we propose DomainSiam to learn Domain-Aware, that is fully utilizing semantic and objectness information from a pre-trained network. DomainSiam consists of the DensSiam with a self-attention module [30] as a backbone network and a regression network to select the most discriminative convolutuinal filters to leverage the semantic and objectness information. Moreover, we develop a differentiable weighted-dynamic domain loss function to train the regression network. The developed loss function is monotonic, dynamic, and smooth with respect to its hyper-parameter, which can be reduced to $l_1$ $or$ $l_2$ during the training phase. On the other hand, the shrinkage loss function [28] is static and it can not be adapted during the training phase. Most regression networks solve the regression problem with static loss such as the closed-form solution if the input to the network is not high-dimensional or minimizing $l_2$. The results will be made available [1].
To summarize, the main contributions of this paper are three-fold.

- A novel architecture is proposed for object tracking to capture the Domain-Aware features with semantic and objectness information. The proposed architecture enables the features to be robust to appearance changes. Moreover, it decreases the sparsity problem as it produces the most important feature space. Consequently, it decreases the overhead calculations.

- Developing a differentiable weighted-dynamic domain loss function specifically for visual object tracking to train the regression network to extract the domain channels that is activated by target category. The developed loss is monotonic with respect to its hyper-parameter, this will be useful in case of high dimensional data and non-convexity. Consequently, this will increase the performance of the tracker.
- The proposed architecture tackles the generalization capability from one domain to another domain (e.g., from ImageNet to VOT datasets).

The rest of the paper is organized as follows. Related work is presented in section 2. Section 3 details the proposed approach. We present the experimental results in Section 4. Finally, section 5 concludes the paper.

## 2   Related work

Recently, Siamese-based trackers have received significant attention especially in realtime tracking. In this section, we firstly introduce the *state-of-the-art*

---

[1] https://vip-mun.github.io/DomainSiam

Siamese-based trackers. Then, we briefly introduce the gradient-based localization guidance.

### Siamese-based Trackers

The first Siamese network was first proposed in [6] for signature verification. In general, Siamese network consists of two branches the target branch and the search branch; both branches share the same parameters. The score map which indicates the position of the object of interest is generated by the last cross-correlation layer. The pioneering work of using Siamese in object tracking is SiamFC [4]. SiamFC searches the target image in the search image. Siamese Instance Search [39] proposed SINT, which has the query branch and the search branch, the backbone of this architecture is inherited from AlexNet [21]. CFNet [40] improved SiamFC by adding a correlation layer to the target branch. SA-Siam [15] proposed two Siamese networks, the first network encodes the semantic information and the second network encodes the appearance model, which is different from our architecture that has only one Siamese network. SiamRPN [24] formulated the tracking problem as a local one-shot detection. SiamRPN consists of a Siamese network as a feature extractor and a region proposal network which includes the classification branch and regression branch. DensSiam [30] used the Densely-Siamese architecture to make the Siamese network deeper while maintaining the performance of the network. DensSiam allows the low-level and high-level features to flow within layers without vanishing gradients problem. Moreover, a self-attention mechanism was integrated to force the network to capture the non-local features. SiamMask [42] used Siamese networks for object tracking using augmentation loss to produce a binary segmentation mask. In addition, the binary segmentation mask locates the object of interest accurately. ATOM [8] proposed the Siamese network with explicit components for target estimation and classification. The component is trained offline to maximize the overlapping between the estimated bounding box and the target. Most Siamese trackers do not fully utilize the semantic and ojectness information from pre-trained networks.

### Gradient-based Localization Guidance

In this category of learning, the objective is to determine the most important channel of the network with respect to the object category. In an object classification task, each category activates a set of certain channels. Grad-CAM [37] used a gradient of any target logit (e.g., "cat") and using this gradient, determined the active category channel for this logit. The work in [46] demonstrated that the global average pooling of the gradients is implicitly acting as Attention for the network; consequently, it can locate the object of interest accurately.

## 3    Proposed Approach

We propose DomainSiam for visual object tracking, the complete pipeline is shown in Fig. 1. The DensSiam with the Self-Attention network is used as a feature extractor, however, in any Siamese network these features do not fully utilize the semantic and objectness information. Furthermore, the channels in Siamese networks are sparsely activated. We use the ridge regression network with a differentiable weighted-dynamic loss function to overcome the previous problems.
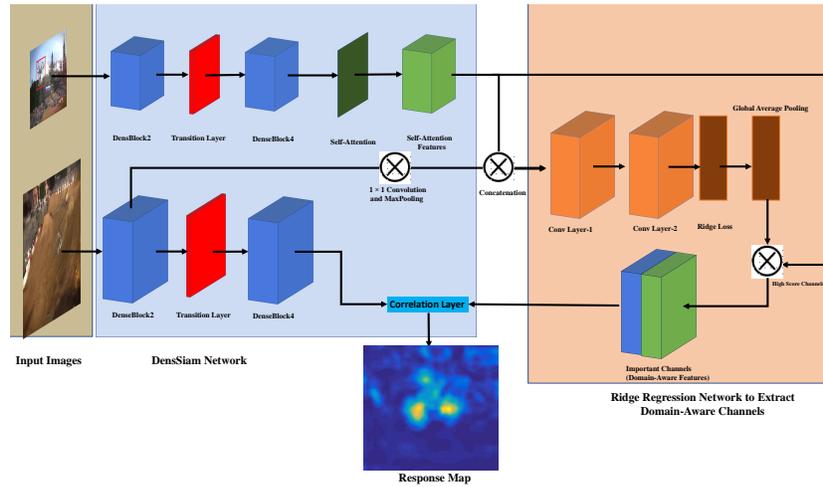


**Fig. 1.** The architecture of DomainSiam tracker. It consists of three blocks, the input images block which includes the target image and search image, DensSiam network with a Self-Attention module at the end of the target branch, and the Ridge Regression Network that highlights the important channels and produces the Domain-Aware features. The response map is produced by the correlation layer which is the final layer. The correlation layer calculates the correlation between the Domain-Aware channels and search branch features which is represented by DenseBlock4.

### 3.1    Ridge Regression Network with Domain-Aware Features

In Fig. 1, the pipeline is divided into three blocks, the input block to the target branch and the search branch, the DensSiam block which has the same architecture in [30], and the ridge regression network. The DensSiam network produces two feature maps for target and search images, respectively. Imbalanced distribution of the training data makes the feature maps produced by Siamese networks less discriminative as there is a high number of easy samples compared to the hard samples. Siamese networks use pre-trained networks which have
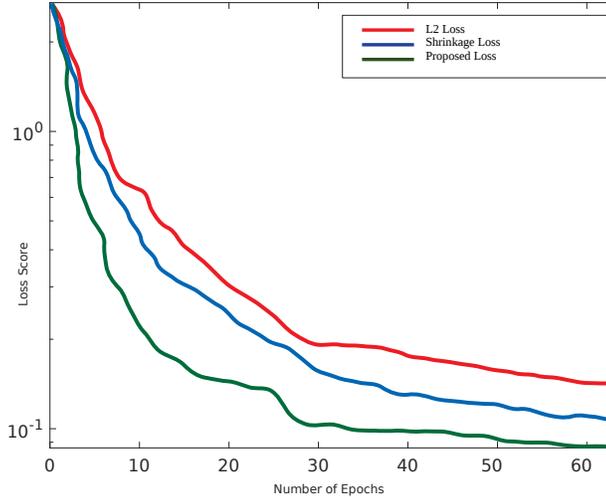
**Fig. 2.** A comparison of convergence speed on $L_2$ loss, Shrinkage loss [28], and our proposed loss function. The average loss is calculated on a batch of eight samples on VOT2018 [19] dataset.

been trained on other tasks (e.g., classification and recognition). These networks increase inter-class differences and is also insensitive to intra-class variations. Consequently, this property decreases the performance of Siamese networks as the tracker needs to be more discriminative to the same object category. Moreover, the pre-trained network is sparsely activated by the object category. In other words, in the feature channels/maps there are only a few active channels that correspond to the object category. The regression network in Fig. 1 highlights the importance of each channel in the feature map to the object of interest and discards the others.

In Fig. 1, the ridge regression network regresses all samples in the input image patch to their soft labels by optimizing the following objective function.

$$\arg\min_{w} \|W * X_{i,j} - Y(i,j)\|^2 + \lambda\|W\|^2 \qquad (1)$$

Where $\| * \|$ denotes the convolution operation, $W$ is the weight of the regression network, $\mathbf{X} \in \mathbb{R}^{N \times D}$ is the input features and $\mathbf{Y} \in \mathbb{R}^{N \times D}$ is the soft label. Gaussian distribution is used as a soft label map and its centre is aligned to the target center and $\lambda > 0$ is the regularization parameter.

$$Y(i,j) = e^{-\frac{i^2+j^2}{2\sigma^2}} \qquad (2)$$

Where $(i, j)$ is the location corresponding to the target location and $\sigma$ is the Gaussian kernel width. The closed-form analytic solution for equation 1 is defined as

$$W = \left(X^\top X + \lambda I\right)^{-1} X^\top Y \tag{3}$$

The optimal solution of $W$ can be achieved by equation 3; however, solving this equation is computationally expensive as $X^\top X \in \mathbb{R}^{D \times D}$. Instead, we use the ridge regression network with the proposed loss function to solve equation 1.

### 3.2   Ridge Regression Optimization

In Fig. 1, the ridge regression network consists of two convolutional layers, ridge loss and the global average pooling. The global average pooling encourages the proposed loss function to localize the object of interest accurately compared to the global max pooling. It is worth mentioning that both global average pooling and global max pooling have similar performances on object classification tasks. As shown in Fig. 1, in the last block, the Domain-Aware feature space is calculated by

$$\delta_i = GAP(\partial L / \partial F_i) \tag{4}$$

Where $\delta$ is the Domain-Aware non-sparse features, $GAP$ is the global average pooling, $L$ is the domain-dynamic loss function which will be discussed later, and $F$ is the input feature channel of the $i^{th}$ channel to the ridge regression network. Let the objective function of the ridge regression network be $x$

$$x = \|W * X_{i,j} - Y(i,j)\|^2 + \lambda\|W\|^2 \tag{5}$$

We propose a differentiable weighted-dynamic loss function for visual object tracking to solve equation 5, inspired by [2] that uses a general loss function for variational autoencoder, monocular depth estimation, and global registration, as follows.

$$L(x, \alpha) = \frac{|\alpha - 2|}{\alpha} e^{ay} \left( \left( \frac{x^2}{|\alpha - 2|} + 1 \right)^{\alpha/2} - 1 \right) \tag{6}$$

where $a \in [0, 1]$ is a hyper-parameter, $y$ is the regression target of a sample, and $\alpha \in \mathbb{R}$ is the parameter that controls the robustness of the loss. The exponent term in this loss function tackles the imbalanced distribution of the training set by assigning a higher weight to hard samples. The imbalanced data occurs when the number of easy samples (background) is extremely higher than the hard samples (foreground).

The advantage of this loss function over equation 1 and equation 3 is it can automatically adjust the robustness during the training phase. This advantage comes from the $\alpha$ parameter. For example, at $\alpha = 2$ the equation 6 becomes $L_2$

$$\lim_{\alpha \to 2} L(x, \alpha) = \frac{e^{ay}}{2} x^2 \tag{7}$$

Similarly, when $\alpha = 1$, the equation 6 becomes $L_1$

$$L(x,\alpha) = (\sqrt{x^2 + 1})e^{ay} - 1 \tag{8}$$

Another advantage of equation 6 is becoming Lorentzia loss function [5] by allowing $\alpha = 0$ as follows

$$\lim_{x \to 0} L(x,\alpha) = \log\left(\frac{1}{2}x^2 + 1\right)e^{ay} \tag{9}$$

As noticed before that the proposed loss function is dynamic, this allows the network to also learn a robust representation. The gradient of the equation 6 with respect to $\alpha$ is always positive. Consequently, this property makes the loss monotonic with respect to $\alpha$ and useful for non-convex optimization.

$$\frac{\partial L}{\partial \alpha}(x,\alpha) \geq 0 \tag{10}$$

The final proposed loss function is given by

$$L(x,\alpha) = \begin{cases} \frac{e^{ay}}{2}x^2 & \text{if } \alpha = 2 \\ \log\left(\frac{1}{2}(x)^2 + 1\right)e^{ay} & \text{if } \alpha = 0 \\ (1 - \exp\left(-\frac{1}{2}(x)^2\right))e^{ay} & \text{if } \alpha = -\infty \\ \frac{|\alpha-2|}{\alpha}e^{ay}\left(\left(\frac{(x)^2}{|\alpha-2|} + 1\right)^{\alpha/2} - 1\right) & \text{otherwise} \end{cases} \tag{11}$$

Fig. 2 shows that the optimization over the proposed loss function achieves faster convergence speed, while in Shrinkage loss function that is proposed in [28] and the original ridge regression loss function 1 ($l_2$), the convergence speed is slower. The importance of each channel in the feature map is calculated by plugging equation 11 into equation 4. It is worth mentioning that the output feature map of the ridge regression network contains only the activated channels that have the most semantic and objectnes information corresponding to the object category. The Domain-Aware features and the feature channels from denseBlock4 are fed into the correlation layer to calculate the similarity and produce the response map.

## 4    Experimental Results

The benchmarks are divided into two categories, the validation set including OTB2013 [43] and OTB2015 [44], and the testing set including VOT2017 [20], VOT2018 [19], and GOT10k [17]. We introduce the implementation details in the next sub-section and then we compare the proposed tracker to the *state-of-the-art* trackers.

**Table 1.** Comparison with the state-of-the-art trackers on VOT2017 in terms of Accuracy (A), expected Average Overlap (EAO), and Robustness (R).

| Tracker | A↑ | EAO↑ | R ↓ | FPS |
|---|---|---|---|---|
| CSRDCF++ | 0.453 | 0.229 | 0.370 | > 25 |
| SAPKLTF | 0.482 | 0.184 | 0.581 | > 25 |
| Staple | 0.530 | 0.169 | 0.688 | > 80 |
| ASMS | 0.494 | 0.169 | 0.623 | > 25 |
| SiamFC | 0.502 | 0.188 | 0.585 | 86 |
| SiamDCF | 0.500 | 0.473 | 0.249 | 60 |
| ECOhc | 0.494 | 0.435 | 0.238 | 60 |
| DensSiam | 0.540 | 0.350 | 0.250 | 60 |
| **DomainSiam(proposed)** | **0.562** | **0.374** | **0.201** | 53 |

### 4.1   Implementation Details

We used the pre-trained DensSiam network (DenseBlock2 and DenseBlock4) that has been trained on Large Scale Visual Recognition Challenge (ILSVRC15) [36]. ILSVRC15 has over 4000 sequences with approximately 1.3 million frames with their labels. The DomainSiam, which has been trained on 1000 classes, can benefit from this class diversity. We implemented DomainSiam in Python using PyTorch framework [35]. Experiments are performed on Linux with a Xeon E5 @2.20 GHz CPU and a Titan XP GPU. Testing speed of DomainSiam is 53 FPS which is beyond realtime speed.

**Training**. The ridge regression network is trained with its proposed loss function separately from the Siamese network with 70 epochs. The highest scores associated with 100 channels are selected as the Domain-Aware features. The training is applied with a momentum of 0.9, the batch size of 8 images, and the learning rate is annealed geometrically at each epoch from $10^{-3}$ to $10^{-8}$.

**Tracking Settings**. The initial scale variation is $O^s$ where $O = 1.0375$ and $s = \{-2, 0, 2\}$. We adopt the target image size of $127 \times 127$ and the search image size of $255 \times 255$ with a linear interpolation to update the scale with a factor of 0.435.

### 4.2   Comparison with the State-of-the-Arts

In this section we use five benchmarks to evaluate DomainSiam against *state-of-the-art* trackers. We use VOT2017 [20], VOT2018 [19], LaSOT [12], TrackingNet [33], and GOT10k [17].

**Results on VOT2017 and VOT2018**

The results on the VOT dataset in Table 1 and Table 2 are given by VOT-Toolkit. DomainSiam outperforms the *state-of-the-art* trackers listed in both tables. It is worth mentioning that DoaminSiam is about 2% higher than the DensSiam tracker in terms of Accuracy (A) and Expected Average Overlap (EAO) in Table 1 while running at 53 frames per second. Table 2 shows that DomainSiam has a robustness of 0.221 which is about 5% higher than the second best tracker

**Table 2.** Comparison with *state-of-the-art* trackers on VOT2018 in terms of Accuracy (A), expected Average Overlap (EAO), and Robustness (R).

| Tracker | A↑ | EAO↑ | R↓ | FPS |
|---|---|---|---|---|
| ASMS [41] | 0.494 | 0.169 | 0.623 | 25 |
| SiamRPN [23] | 0.586 | 0.383 | 0.276 | 160 |
| SA_Siam_R [15] | 0.566 | 0.337 | 0.258 | 50 |
| FSAN [19] | 0.554 | 0.256 | 0.356 | 30 |
| CSRDCF [29] | 0.491 | 0.256 | 0.356 | 13 |
| SiamFC [4] | 0.503 | 0.188 | 0.585 | 86 |
| SAPKLTF [19] | 0.488 | 0.171 | 0.613 | 25 |
| DSiam [14] | 0.215 | 0.196 | 0.646 | 25 |
| ECO [7] | 0.484 | 0.280 | 0.276 | 60 |
| **DomainSiam(proposed)** | **0.593** | **0.396** | **0.221** | 53 |

**Table 3.** Comparisons with *state-of-the-art* trackers on TrackingNet dataset in terms of the Precision (PRE), Normalized Precision (NPRE), and Success.

| Tracker | PRE ↑ | NPRE ↑ | SUC.↑ |
|---|---|---|---|
| Staple_CA [32] | 0.468 | 0.605 | 0.529 |
| BACF [13] | 0.461 | 0.580 | 0.523 |
| MDNet [34] | 0.565 | 0.705 | 0.606 |
| CFNet [40] | 0.533 | 0.654 | 0.578 |
| SiamFC [4] | 0.533 | 0.663 | 0.571 |
| SAMF [27] | 0.477 | 0.598 | 0.504 |
| ECO-HC [7] | 0.476 | 0.608 | 0.541 |
| Staple [3] | 0.470 | 0.603 | 0.528 |
| ECO [7] | 0.492 | 0.618 | 0.554 |
| CSRDCF [29] | 0.480 | 0.622 | 0.534 |
| **DomainSiam(proposed)** | **0.585** | **0.712** | **0.635** |

**Table 4.** Comparison with *state-of-the-art* trackers on LaSOt dataset in terms of the Normalized Precision and Success.

| Tracker | Norm. Prec. (%)↑ | Success (%)↑ |
|---|---|---|
| MDNet [34] | 46.0 | 39.7 |
| DaSiam [47] | 49.6 | 41.5 |
| STRCF [25] | 34.0 | 30.8 |
| SINT [39] | 35.4 | 31.4 |
| StrucSiam [45] | 41.8 | 33.5 |
| SiamFC [4] | 42.0 | 33.6 |
| VITAL [38] | 45.3 | 39.0 |
| ECO [9] | 33.8 | 32.4 |
| DSiam [14] | 40.5 | 33.3 |
| **DomainSiam(proposed)** | 53.7 | 43.6 |

(SiamRPN) while outperforming all other trackers in terms of accuracy and expected average overlap.

**Table 5.** Comparison *state-of-the-art* trackers on GOT10k dataset in terms of Average Overlap (AO), and Success Rates (SR) at overlap thresholds of 0.50 and 0.75.

| TRACKER | DomainSiam (proposed) | CFNet | SiamFC | GOTURN | CCOT | ECO | HCF | MDNet |
|---------|-----------------------|-------|--------|--------|------|-----|-----|-------|
| AO | **0.414** | 0.374 | 0.348 | 0.347 | 0.325 | 0.316 | 0.315 | 0.299 |
| SR(0.50) | **0.451** | 0.404 | 0.353 | 0.375 | 0.328 | 0.309 | 0.297 | 0.303 |
| SR(0.75) | **0.214** | 0.144 | 0.098 | 0.124 | 0.107 | 0.111 | 0.088 | 0.099 |

**Results on TrackingNet Dataset**
This is a large-scale dataset that was collected from YouTube videos. Table 3 shows that DomainSiam outperforms MDNet which is the second best tracker on the TrackingNet dataset. with 2% in terms of precision and about 3% in terms of success. DomainSiam outperforms all other trackers on TrackingNet dataset.
**Results on LaSOT Dataset**
The average sequence length in this dataset is about 2500 frames. Table 4 shows that DomainSiam achieves the best success score with over 2% from the second best tracker (DaSiam). Our tracker significantly outperforms DaSiam with 4% in terms of normalized precision.
**Results on GOT10k Dataset**
This dataset has 180 test sequences. We tested the proposed tracker against 7 trackers as shown in Table 5. DomainSiam outperforms CFNet which is the best tracker in terms of Average Overlap (AO) with 4%. It is worth mentioning that DomainSiam achieves the best performance among all trackers in terms of success Rate (SR) at thresholds of 0.50 and 0.75.

## 5   Conclusions and Future Work

In this paper, we introduced DomainSiam tracker, a Siamese with a ridge regression network to fully utilize the semantic and objectness information for visual object tracking while also producing a class-agnostic. We developed a differentiable weighted-dynamic loss function to solve the ridge regression problem. The developed loss function improves the feature learning as it automatically adjusts the robustness during the training phase. Furthermore, it utilizes the activated channels which correspond to the object category label. The proposed architecture decreases the sparsity problem in Siamese networks and provides an efficient Domain-Aware feature space that is robust to appearance changes. DomainSiam does not need to be re-trained from scratch as the ridge regression network with the proposed loss function is trained separately from the Siamese network. DomainSiam with the proposed loss function exhibits a superior convergence speed compared to other loss functions. The ridge regression network with the proposed loss function can be extended to other tasks such as object detection and semantic segmentation.

# References

1. Alahari, K., Berg, A., Hager, G., Ahlberg, J., Kristan, M., Matas, J., Leonardis, A., Cehovin, L., Fernandez, G., Vojir, T., et al.: The thermal infrared visual object tracking vot-tir2015 challenge results. In: Computer Vision Workshop (ICCVW), 2015 IEEE International Conference on. pp. 639–651. IEEE (2015)
2. Barron, J.T.: A general and adaptive robust loss function. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4331–4339 (2019)
3. Bertinetto, L., Valmadre, J., Golodetz, S., Miksik, O., Torr, P.H.: Staple: Complementary learners for real-time tracking. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1401–1409 (2016)
4. Bertinetto, L., Valmadre, J., Henriques, J.F., Vedaldi, A., Torr, P.H.: Fully-convolutional siamese networks for object tracking. In: European conference on computer vision. pp. 850–865. Springer (2016)
5. Black, M.J., Anandan, P.: The robust estimation of multiple motions: Parametric and piecewise-smooth flow fields. Computer vision and image understanding **63**(1), 75–104 (1996)
6. Bromley, J., Guyon, I., LeCun, Y., Säckinger, E., Shah, R.: Signature verification using a" siamese" time delay neural network. In: Advances in neural information processing systems. pp. 737–744 (1994)
7. Danelljan, M., Bhat, G., Khan, F.S., Felsberg, M.: Eco: Efficient convolution operators for tracking. In: Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, USA. pp. 21–26 (2017)
8. Danelljan, M., Bhat, G., Khan, F.S., Felsberg, M.: Atom: Accurate tracking by overlap maximization. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4660–4669 (2019)
9. Danelljan, M., Bhat, G., Shahbaz Khan, F., Felsberg, M.: Eco: Efficient convolution operators for tracking. In: CVPR (2017)
10. Danelljan, M., Hager, G., Shahbaz Khan, F., Felsberg, M.: Convolutional features for correlation filter based visual tracking. In: Proceedings of the IEEE International Conference on Computer Vision Workshops. pp. 58–66 (2015)
11. Danelljan, M., Hager, G., Shahbaz Khan, F., Felsberg, M.: Learning spatially regularized correlation filters for visual tracking. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 4310–4318 (2015)
12. Fan, H., Lin, L., Yang, F., Chu, P., Deng, G., Yu, S., Bai, H., Xu, Y., Liao, C., Ling, H.: Lasot: A high-quality benchmark for large-scale single object tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5374–5383 (2019)
13. Galoogahi, H.K., Fagg, A., Lucey, S.: Learning background-aware correlation filters for visual tracking. In: Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, USA. pp. 21–26 (2017)
14. Guo, Q., Feng, W., Zhou, C., Huang, R., Wan, L., Wang, S.: Learning dynamic siamese network for visual object tracking. In: Proc. IEEE Int. Conf. Comput. Vis. pp. 1–9 (2017)
15. He, A., Luo, C., Tian, X., Zeng, W.: A twofold siamese network for real-time object tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4834–4843 (2018)
16. Henriques, J.F., Caseiro, R., Martins, P., Batista, J.: High-speed tracking with kernelized correlation filters. IEEE Transactions on Pattern Analysis and Machine Intelligence **37**(3), 583–596 (2015)

17. Huang, L., Zhao, X., Huang, K.: Got-10k: A large high-diversity benchmark for generic object tracking in the wild. arXiv preprint arXiv:1810.11981 (2018)
18. Kendall, A., Grimes, M., Cipolla, R.: Posenet: A convolutional network for real-time 6-dof camera relocalization. In: Computer Vision (ICCV), 2015 IEEE International Conference on. pp. 2938–2946. IEEE (2015)
19. Kristan, M., Leonardis, A., Matas, J., Felsberg, M., Pflugfelder, R., Cehovin Zajc, L., Vojir, T., Bhat, G., Lukezic, A., Eldesokey, A., et al.: The sixth visual object tracking vot2018 challenge results. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 0–0 (2018)
20. Kristan, M., Leonardis, A., Matas, J., Felsberg, M., Pflugfelder, R., Cehovin Zajc, L., Vojir, T., Hager, G., Lukezic, A., Eldesokey, A., et al.: The visual object tracking vot2017 challenge results. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1949–1972 (2017)
21. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems. pp. 1097–1105 (2012)
22. Lenc, K., Vedaldi, A.: Understanding image representations by measuring their equivariance and equivalence. In: Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2015)
23. Li, B., Yan, J., Wu, W., Zhu, Z., Hu, X.: High performance visual tracking with siamese region proposal network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 8971–8980 (2018)
24. Li, B., Yan, J., Wu, W., Zhu, Z., Hu, X.: High performance visual tracking with siamese region proposal network. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
25. Li, F., Tian, C., Zuo, W., Zhang, L., Yang, M.H.: Learning spatial-temporal regularized correlation filters for visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4904–4913 (2018)
26. Li, X., Ma, C., Wu, B., He, Z., Yang, M.H.: Target-aware deep tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1369–1378 (2019)
27. Li, Y., Zhu, J.: A scale adaptive kernel correlation filter tracker with feature integration. In: European Conference on Computer Vision. pp. 254–265. Springer (2014)
28. Lu, X., Ma, C., Ni, B., Yang, X., Reid, I., Yang, M.H.: Deep regression tracking with shrinkage loss. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 353–369 (2018)
29. Lukezic, A., Vojir, T., Zajc, L.C., Matas, J., Kristan, M.: Discriminative correlation filter with channel and spatial reliability. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. vol. 2 (2017)
30. Mohamed, M.M.: Denssiam: End-to-end densely-siamese network with self-attention model for object tracking. In: Advances in Visual Computing: 13th International Symposium, ISVC 2018, Las Vegas, NV, USA, November 19–21, 2018, Proceedings. vol. 11241, p. 463. Springer (2018)
31. Molchanov, P., Yang, X., Gupta, S., Kim, K., Tyree, S., Kautz, J.: Online detection and classification of dynamic hand gestures with recurrent 3d convolutional neural network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4207–4215 (2016)
32. Mueller, M., Smith, N., Ghanem, B.: Context-aware correlation filter tracking. In: Proc. IEEE Conf. Comput. Vis. Pattern Recognit.(CVPR). pp. 1396–1404 (2017)

33. Muller, M., Bibi, A., Giancola, S., Alsubaihi, S., Ghanem, B.: Trackingnet: A large-scale dataset and benchmark for object tracking in the wild. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 300–317 (2018)
34. Nam, H., Han, B.: Learning multi-domain convolutional neural networks for visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4293–4302 (2016)
35. Paszke, A., Gross, S., Chintala, S., Chanan, G.: Pytorch: Tensors and dynamic neural networks in python with strong gpu acceleration. PyTorch: Tensors and dynamic neural networks in Python with strong GPU acceleration **6** (2017)
36. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. International Journal of Computer Vision **115**(3), 211–252 (2015)
37. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 618–626 (2017)
38. Song, Y., Ma, C., Wu, X., Gong, L., Bao, L., Zuo, W., Shen, C., Lau, R.W., Yang, M.H.: Vital: Visual tracking via adversarial learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 8990–8999 (2018)
39. Tao, R., Gavves, E., Smeulders, A.W.: Siamese instance search for tracking. In: Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on. pp. 1420–1429. IEEE (2016)
40. Valmadre, J., Bertinetto, L., Henriques, J., Vedaldi, A., Torr, P.H.: End-to-end representation learning for correlation filter based tracking. In: Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on. pp. 5000–5008. IEEE (2017)
41. Vojir, T., Noskova, J., Matas, J.: Robust scale-adaptive mean-shift for tracking. Pattern Recognition Letters **49**, 250–258 (2014)
42. Wang, Q., Zhang, L., Bertinetto, L., Hu, W., Torr, P.H.: Fast online object tracking and segmentation: A unifying approach. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1328–1338 (2019)
43. Wu, Y., Lim, J., Yang, M.H.: Online object tracking: A benchmark. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2411–2418 (2013)
44. Wu, Y., Lim, J., Yang, M.H.: Object tracking benchmark. IEEE Transactions on Pattern Analysis and Machine Intelligence **37**(9), 1834–1848 (2015)
45. Zhang, Y., Wang, L., Qi, J., Wang, D., Feng, M., Lu, H.: Structured siamese network for real-time visual tracking. In: Proceedings of the European conference on computer vision (ECCV). pp. 351–366 (2018)
46. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2016)
47. Zhu, Z., Wang, Q., Li, B., Wu, W., Yan, J., Hu, W.: Distractor-aware siamese networks for visual object tracking. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 101–117 (2018)