

# DEEPMIX: ONLINE AUTO DATA AUGMENTATION FOR ROBUST VISUAL OBJECT TRACKING

Ziyi Cheng<sup>1\*</sup>, Xuhong Ren<sup>2\*</sup>, Felix Juefei-Xu<sup>3</sup>, Wanli Xue<sup>2†</sup>, Qing Guo<sup>4†</sup>, Lei Ma<sup>5,6,1</sup>, Jianjun Zhao<sup>1</sup>

<sup>1</sup>Kyushu University, Japan    <sup>2</sup>Tianjin University of Technology, China

<sup>3</sup>Alibaba Group, USA    <sup>4</sup>Nanyang Technological University, Singapore

<sup>5</sup>University of Alberta, Canada    <sup>6</sup>Alberta Machine Intelligence Institute (Amii), Canada

## ABSTRACT

Online updating of the object model via samples from historical frames is of great importance for accurate visual object tracking. Recent works mainly focus on constructing effective and efficient updating methods while neglecting the training samples for learning discriminative object models, which is also a key part of a learning problem. In this paper, we propose the *DeepMix* that takes historical samples' embeddings as input and generates augmented embeddings online, enhancing the state-of-the-art online learning methods for visual object tracking. More specifically, we first propose the *online data augmentation* for tracking that online augments the historical samples through object-aware filtering. Then, we propose *MixNet* which is an offline trained network for performing online data augmentation within one-step, enhancing the tracking accuracy while preserving high speeds of the state-of-the-art online learning methods. The extensive experiments on three different tracking frameworks, *i.e.*, DiMP, DSiam, and SiamRPN++, and three large-scale and challenging datasets, *i.e.*, OTB-2015, LaSOT, and VOT, demonstrate the effectiveness and advantages of the proposed method.

**Index Terms**— data augmentation, online updating, visual object tracking, deepmix, mixnet

## 1. INTRODUCTION

Visual object tracking (VOT) is one of the most widely studied computer vision approaches that can produce the trajectory of the moving object from a sequence of frames. It has seen ubiquitous applications ranging from navigation for robots, intelligent video surveillance, smart logistics, robotics for manufacturing, *etc.* Based on the frame processing procedure, the VOT can be divided into online tracking and offline tracking. For online tracking, only the current frame and the previous frames can be used to determine the tracking result for the current frame, and the tracking results for the previous frames, once computed, can no longer be altered based on later frames.

\*Ziyi Cheng and Xuhong Ren are co-first authors and contribute equally to this work. †Wanli Xue and Qing Guo are corresponding authors (xue-wanli@email.tjut.edu.cn and tsingqguo@ieee.org).

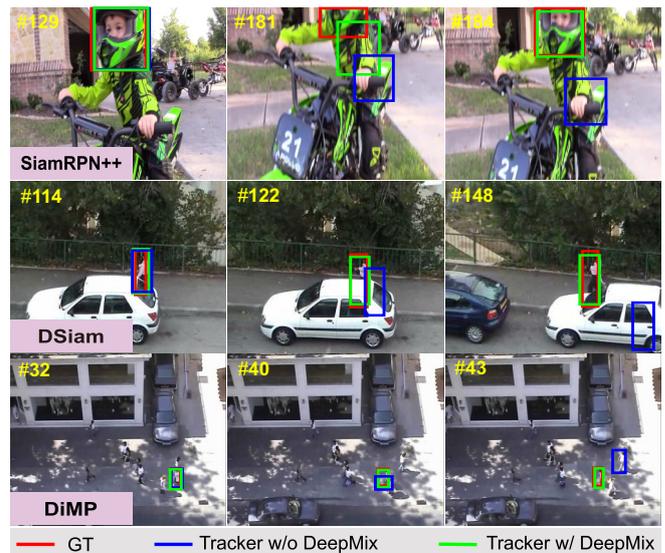


Fig. 1: Examples of three state-of-the-art trackers, *e.g.*, SiamRPN++ [1], DSiam [2, 3], and DiMP [4] with or without the proposed *DeepMix*.

For offline tracking, the tracking results can be produced after having access to all the frames. Needless to say, the settings of online tracking make it more suitable for real-world applications and deployment.

In the early days, VOT is implemented using correlation filters (CFs) with notable works such as the minimum output sum of squared error (MOSSE) CF [5], kernelized correlation filter (KCF) [6], spatially regularized discriminative correlation filter [7, 8, 9, 10, 11], *etc.* With the fast-paced development of deep learning approaches that are catered towards computer vision problems, deep learning based VOT has gained tremendous popularity. Methods such as SiamFC [12] and subsequent SiamRPN [13], DiMP [14], DSiam [2], DaSiamRPN [15], SiamRPN++ [16], *etc.*, have progressively produced better results on various VOT benchmarks [17, 18, 19, 20]. We think that online updating of the object model via samples from historical frames is of great importance for accurate visual object tracking. These mentioned recent works mainly focus on constructing powerful deep backbones or effective and efficient online updating methods while neglecting the training samples for learning discriminative object models, which is

also a key part of a learning problem in our opinion.

In this paper, we propose the *DeepMix* that takes historical samples’ embeddings as input and generates augmented embeddings online, enhancing the state-of-the-art online learning methods for visual object tracking. More specifically, we first propose the *online data augmentation* for tracking that online augments the historical samples through object-aware filtering. Then, we propose *MixNet* which is an offline trained network for performing online data augmentation within one-step, enhancing the tracking accuracy while preserving high speeds of the state-of-the-art online learning methods. *MixNet* predicts different convolution parameters dramatically for object and background regions, respectively, according to the input training samples, thus is able to generate effective training samples for online data augmentation. As shown in Fig. 1, our *DeepMix* let three state-of-the-art trackers, *i.e.*, DiMP [14], DSiam [2, 3], and SiamRPN++ [1], localize objects more accurately. The extensive experiments on the above three different tracking frameworks and three large-scale and challenging datasets, *i.e.*, OTB-2015, LaSOT, and VOT, further demonstrate the effectiveness and advantages of the proposed method.

## 2. RELATED WORK

### 2.1. General Visual Object Tracking

In visual object tracking task, Siamese network-based methods [12, 2, 15, 16, 21] have been a popular solution. Among them, SiamFC [12] is an initial implementation, which extracts the deep features from template and search region and performs cross-correlation to predict object position. SiamRPN [13] proposes to add a bounding box regressor to the Siamese Network as done in the object detection task.

Another effective solution treats tracking as a online learning problem. They [22, 14, 23, 24, 25] train a classifier from past frames with an online update strategy and distinguish the target from background via the classifier. DiMP [14] applies meta-learning formulation by collecting historical frame feature representation to classify objects. ROAM [26] offline trains an LSTM [27] to generate the adaptive learning rate to enhance online training. Several works also focus on the training sample management during the online tracking. [28] manages samples by assigning different weights to training samples. ECO [29] employs a Gaussian mixture model (GMM) to choose more distinguishing samples. UpdateNet [30] uses history templates to generating a more robust template through a CNN which is similar to our *DeepMix*. However, it aims to update effective templates. Although above works have shown great advantages of managing training samples, they ignore how to make more effective use of online samples, *e.g.*, data augmentation. In this work, we focus on the data augmentation for the state-of-the-art trackers during tracking.

### 2.2. Data Augmentation Methods

Data augmentation is an important method to improve generalization performance of deep models. Early works [31, 32, 33] employ cropping, horizontal and vertical flips, and rotation to generate diverse data. Recently, AutoAugment [34] automatically searches for augmentation policies given a predefined set of transformations, which needs a great quantity of training time. Several studies reduce the search costs significantly, such as Population-based augmentation (PBA) [35] and fast AutoAugment (FAA) [36]. These works focus on augmenting a single image.

Recent works [37, 38, 39] propose to mix multiple images for data augmentation, inspiring our research on mixing multiple training samples for tracking. Specifically, Mixup [38] utilizes an element-wise combination of two images. CutMix [39] replaces a portion of an image with the contents of another image. AugMix [37] merges multiple images where each image is processed by several augment operations randomly, making the trained model see diverse samples. All of these methods focus on image-level augmentation and do not consider the speed for online data. In contrast to previous works, we propose the *DeepMix* allowing efficient online data augmentation for visual object tracking.

## 3. METHOD

In this section, we first discuss the background and motivation of this work and formulate the online data augmentation for tracking in Sec. 3.1 and 3.2. Then, we propose the *MixNet* in Sec. 3.3 to realize effective and efficient online data augmentation for tracking. Finally, we detail how to embed *MixNet* into state-of-the-art trackers (*i.e.*, SiamRPN++ [16], DSiam [2, 3], and DiMP [14]) in Sec. 3.4.

### 3.1. Background and Motivation

Given a live video  $\mathcal{V} = \{\mathbf{I}_t\}_{t=1}^T$  having  $T$  frames and the object bounding box annotated at the first frame (*i.e.*,  $\mathbf{b}_1$ ), we aim to estimate the object’s position and size at the subsequent  $T - 1$  frames. Most of the state-of-the-art methods complete this task by maintaining an object model for matching it with the subsequent frames. In general, we formulate the object localization at  $t$ -th frame as

$$\mathbf{p}_t = \arg \max_{\mathbf{p}} \mathbf{M}_t[\mathbf{p}] = \arg \max_{\mathbf{p}} \text{loc}(\varphi(\mathbf{I}_t), \theta_t), \quad (1)$$

where the  $\mathbf{M}_t$  is a heat map whose the maximum (*i.e.*,  $\mathbf{M}_t[\mathbf{p}_t]$ ) indicates the object’s position in the frame  $\mathbf{I}_t$ , and it can be calculated by the  $\varphi(\mathbf{I}_t)$  and  $\theta_t$  where  $\varphi(\cdot)$  is the backbone network for extracting embedding.

The object model  $\theta_t$  determines the localization accuracy, which is initialized at the first frame and updated at subsequent frames. For example, in the popular Siamese network-based trackers [12, 16], the object model is constructed by using the

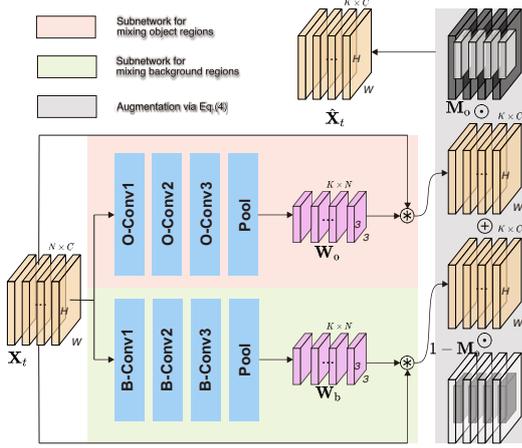


Fig. 2: Architecture of MixNet that contains two subnetworks to estimate filtering parameters for mixing object and background regions, respectively.

embedding of the object at the first frame, *i.e.*,  $\theta_t = \varphi(\mathbf{I}_1)$ , and the localization is implemented by using cross-correlation, *i.e.*,  $\text{loc}(\varphi(\mathbf{I}_t), \theta_t) = \varphi(\mathbf{I}_t) * \varphi(\mathbf{I}_1)$ . More recently, [14] proposes DiMP that uses an online updated classifier for localization (*i.e.*, the  $\text{loc}(\cdot)$  is set as a convolution layer).

The state-of-the-art trackers (*e.g.*, DiMP [14] and DSiam [2]) online update the object model to adapt object and background appearance variation. In general, we formulate the object model updating via

$$\theta_{t+1} = \text{update}(\theta_t, \mathcal{X}_t) \quad (2)$$

where  $\mathcal{X}_t$  denotes the set of training samples that are cropped from the historical frames in which the objects are previously detected. For example, DSiam [2] proposes to online update the object model of the Siamese network via a transformation that is learned from the previous frame. DiMP [14] updates the classifier’s parameters through a pre-trained model predictor that takes 50 samples from previous frames as inputs.

Note that, the updating process is a typical learning module and recent works have demonstrated that data augmentation is of great importance for enhancing the accuracy of image classification under various interferences [40]. Following similar ideas, we aim to explore how to online augment effective training samples, *i.e.*,  $\mathcal{X}_t$ , for the state-of-the-art trackers with existing updating methods, *i.e.*, online transformation in DSiam [2] and model predictor in DiMP [14].

### 3.2. Online Data Augmentation for Tracking

A simple way for online data augmentation is to borrow the existing techniques and conduct augmentation on the collected historical training samples, *i.e.*,  $\mathcal{X}_t$ , through

$$\hat{\mathcal{X}}_t = \text{aug}(\mathcal{X}_t, \mathcal{T}) \quad (3)$$

where  $\mathcal{T}$  denotes the set of sample-level transformations (*e.g.*, adding noise, blur, and rain) and  $\hat{\mathcal{X}}_t$  is the augmented samples. The state-of-the-art data augmentation techniques are usually

employed in the offline training process with random and diverse degradation-related sample-level transformations [40]. For example, AugMix [40] conducts multiple random augmentations for a raw sample and mixes them up for training the image classifiers. However, these methods could not adapt to the online data augmentation for visual object tracking directly due to the following reasons: ❶ the random sample-level transformations would generate new samples that require high time costs to extract their deep features, slowing down the trackers significantly. ❷ the transformations are based on degradation factors (*e.g.*, noise, blur, fog, rain, *etc.*) that can corrupt the raw samples, leading to the less discriminative object model.

To address above challenges, we propose the online data augmentation on embeddings of training samples. Specifically, we first extract embeddings of training samples in  $\mathcal{X}_t$  and get  $\{\varphi(\mathbf{I}_i) \in \mathbb{R}^{C \times W \times H} | \mathbf{I}_i \in \mathcal{X}_t\}$ . Then, we concatenate all embeddings and obtain a tensor  $\mathbf{X}_t \in \mathbb{R}^{N \times C \times W \times H}$  where  $N$  is the number of training samples in  $\mathcal{X}_t$ . Our goal is to map the tensor  $\mathbf{X}_t$  to a new counterpart denoted as  $\hat{\mathbf{X}}_t \in \mathbb{R}^{K \times C \times W \times H}$  that can be fed into existing updating modules to produce a more effective object model. Note that, performing augmentation on the embedding level is much more efficient than that on the sample level, which alleviates the first challenge. In terms of the second challenge, we mix embeddings of all samples with the guidance of previously localization results. Intuitively, the interested object might be at any position in the scene during the video capturing process and it is reasonable to augment the training samples by putting the object to possible background regions. To this end, given  $\mathbf{X}_t$  and the detected bounding boxes  $\{\mathbf{b}_i \in \mathbb{R}^{4 \times 1} | i = 1, \dots, N\}$  of  $N$  training samples, we can split the samples to object and background regions and mixing them up to produce new samples. We formulate this process by

$$\hat{\mathbf{X}}_t = (\mathbf{W}_o \circledast \mathbf{X}_t) \odot \mathbf{M}_o + (\mathbf{W}_b \circledast \mathbf{X}_t) \odot (1 - \mathbf{M}_o), \quad (4)$$

where  $\mathbf{M}_o \in \mathbb{R}^{N \times C \times W \times H}$  are binary masks for the  $N$  training samples where the elements within the object regions are set to one while others are zero. The object regions are obtained according to the detection results  $\{\mathbf{b}_i \in \mathbb{R}^{4 \times 1} | i = 1, \dots, N\}$  and the ‘ $\odot$ ’ denotes the element-wise multiplication. Besides, the ‘ $\circledast$ ’ denotes the convolution layer while the tensor  $\mathbf{X}_t$  is filtered by  $\mathbf{W}_{\{o \text{ or } b\}} \in \mathbb{R}^{K \times N \times 3 \times 3}$  with the padding hyper-parameter to be one. The variable  $K$  denotes the number of samples in the output tensor. More specifically, for the  $c$ th channel of  $\mathbf{X}_t$  (*i.e.*,  $\mathbf{X}_t[c] \in \mathbb{R}^{N \times W \times H}$ ), we filter it with  $\mathbf{W}_{\{o \text{ or } b\}}$  and get the  $c$ th channel of  $\hat{\mathbf{X}}_t[c] \in \mathbb{R}^{K \times W \times H}$ . Intuitively,  $\mathbf{W}_{\{o \text{ or } b\}}$  indicates how to fuse the  $N$  training samples and get  $K$  new samples. The  $\mathbf{W}_o$  and  $\mathbf{W}_b$  take the charge of mixing object and background regions, respectively.

However, to make above method work, we should consider the following issues: ❶ how to estimate  $\mathbf{W}_o$  and  $\mathbf{W}_b$  to fit different cases? ❷ how to make the above module to be efficient? To alleviate these issues, we propose the *MixNet* that is able to produce augmented data in one step.

**Table 1:** Comparison with SOTA Trackers under OPE setup.

Dataset	OTB-2015		LaSOT	
	AUC	Prec	Success	Prec
Metrics				
MLT [41]	0.611	-	0.368	-
GradNet [42]	0.639	0.861	0.365	0.351
ATOM [43]	0.667	0.879	0.514	0.505
SiamDW [41]	0.674	-	0.384	-
POST [44]	0.678	0.907	0.481	0.463
CRPN [45]	0.675	-	0.455	-
ASRCF [46]	0.692	0.922	0.359	0.337
MAML [24]	0.712	-	0.523	-
DSiam	0.646	0.845	0.438	0.431
DSiam-DeepMix	<b>0.658</b>	<b>0.861</b>	<b>0.439</b>	<b>0.431</b>
SiamRPN++	0.650	0.853	0.447	0.446
SiamRPN++-DeepMix	<b>0.663</b>	<b>0.870</b>	<b>0.459</b>	<b>0.463</b>
DiMP	0.660	0.859	0.532	0.532
DiMP-DeepMix	<b>0.683</b>	<b>0.890</b>	<b>0.536</b>	<b>0.538</b>

### 3.3. MixNet for Efficient Online Data Augmentation

We propose *MixNet* that takes the  $\mathbf{X}_t$  as the input and predict the kernels  $\mathbf{W}_o$  and  $\mathbf{W}_b$  that are suitable for  $\mathbf{X}_t$ . We can use the pre-trained *MixNet* to generate the kernels in a one-step way, leading to efficient online data augmentation. We show the architecture of *MixNet* in Fig. 2. Specifically, *MixNet* contains two sub-networks for predicting the  $\mathbf{W}_o$  and  $\mathbf{W}_b$ , respectively. The two sub-networks share the same architecture but have independent parameters. The architecture has three convolution layers with a kernel size of  $3 \times 3$  and an averaging pooling layer. Note that, we can embed *MixNet* into diverse tracking frameworks by properly setting the input and output channels for them.

### 3.4. Implementation for SOTA Trackers

In this part, we detail the way of using our method for three state-of-the-art trackers, *i.e.*, SiamRPN++ [13, 1], DSiam [2, 3], and DiMP [14]. Simply, we can embed *DeepMix* into these trackers by transforming their training samples and get  $\hat{\mathbf{X}}_t$ . Then, we mix it with the raw training samples (*i.e.*,  $\mathbf{X}_t$ ) by  $\alpha_1 \hat{\mathbf{X}}_t + \alpha_2 \mathbf{X}_t$  and feed it into the updating modules. For all examples, we fix  $\alpha_1 = 0.05$  and  $\alpha_2 = 0.8$  and discuss its influence in Sec. 4.3.

**DSiam and SiamRPN++ with DeepMix.** We collect historical samples ( $\mathbf{X}_t$  size is  $15 \times 256 \times 29 \times 29$ ) to generate the kernel ( $\mathbf{W}_{\{o \text{ or } b\}}$  size is  $15 \times 1 \times 3 \times 3$ ) and then filter with the features (its size is  $1 \times 256 \times 29 \times 29$ ) of the current frame. Finally, *DeepMix* output the new samples ( $\hat{\mathbf{X}}_t$  size is  $1 \times 256 \times 29 \times 29$ ).

**DiMP with DeepMix.** DiMP stores samples to train the classifier, we directly input these samples ( $\mathbf{X}_t$  size is  $50 \times 256 \times 22 \times 22$ ) and obtain new samples with the same size as the Fig. 2 illustrate. We directly train the *MixNet* along with its embedded trackers. It means that we can simply use the original training program and training data of the targeted trackers to train their own *MixNets*. We make some minor

**Table 2:** Comparison with SOTA Trackers on VOT2018

Metrics	EAO	Accuracy	Robustness
DaSiamRPN [15]	0.383	0.586	0.276
SiamMask [47]	0.387	0.642	0.295
MAML [24]	0.392	0.635	0.220
UpdateNet [30]	0.393	-	-
ATOM [43]	0.401	0.590	0.204
SiamDW [48]	0.405	0.597	0.234
DSiam	0.266	0.577	0.421
DSiam-DeepMix	<b>0.287</b>	<b>0.58</b>	<b>0.407</b>
SiamRPN++	0.348	0.583	0.29
SiamRPN++-DeepMix	<b>0.405</b>	<b>0.597</b>	<b>0.234</b>
DiMP	0.214	0.578	0.553
DiMP-DeepMix	<b>0.234</b>	<b>0.612</b>	<b>0.51</b>

modifications to *MixNet* to adapt to different trackers.

**Training details.** For DSiam and SiamRPN++, the original training program uses a pair of images (*i.e.*, template and search region) as a training sample. For each sample, We apply data augmentation strategies on a template to construct a training set containing 15 samples. We input them to *MixNet* and generate a new sample to mix with the original template. We implement the SGD optimizer with a weight decay of 0.0005, base lr of 0.005, and momentum of 0.9. We train the *MixNet* for 40 epochs and 6000 samples per epoch.

In terms of DiMP, its training program pick up 3 images from each video as training samples for its model predictor. In order to match our *MixNet* during testing, we change it to 50 images from each video for *MixNet*. We apply SGD optimizer with weight decay of 0.0005 for all parameter layers, base lr of 0.005, and momentum of 0.9. We train the *MixNet* for 50 epochs and 1000 videos per epoch.

## 4. EXPERIMENTS

### 4.1. Setups

**Datasets, metrics, and baseline.** We evaluate *DeepMix* on three tracking benchmarks: VOT-2018 [18] (60 videos, 356 frames average length), LaSOT [19] (280 videos, 2448 frames average length), OTB-2015 [17] (100 videos, 590 frames average length). VOT-2018 implements a reset-based evaluation that once the object is lost, the tracker is restarted with the ground truth box five frames later and gets a penalty. The main evaluation criterion is expected average overlap (EAO) [49]. The higher EAO indicates better performance. OTB-2015 and LaSOT only give the trackers the ground-truth of initial frame and obtain bounding box sequence, terms *one-pass evaluation*. We use the AUC, which represents the area under the curve of the success plot, to evaluate the performance of trackers.

We compare *DeepMix*-based trackers with six top methods on VOT2018, *i.e.*, DaSiamRPN [15], SiamMask [47], MAML [24], UpdateNet [30], ATOM [43], SiamDW [48]. We also choose five excellent trackers on OTB2015, *i.e.*, MLT [41], GradNet [42], POST [44], CRPN [45], ASRCF [46].

**Table 3:** Ablation analysis of different backbones

Metrics	AUC	Precision	NormPrecision
DiMP (ResNet18)	0.660	0.859	0.807
DiMP (ResNet18)-DeepMix	<b>0.683</b>	<b>0.890</b>	<b>0.834</b>
DiMP (ResNet50)	0.684	0.894	0.842
DiMP (ResNet50)-DeepMix	<b>0.694</b>	<b>0.900</b>	<b>0.852</b>

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

We compare DeepMix with the state-of-the-art methods on three challenging tracking benchmarks. Respectively, the backbones of DiMP, DSiam and SiamRPN++ are ResNet18 [50], AlexNet [32] and MobileV2 [51] on account of DeepMix’s extreme improvement for a simple network.

**OTB-2015 and LaSOT.** We report results on the OTB-2015 and LaSOT datasets in Table 1. Results show that: ❶ DeepMix improves all of its original trackers. ❷ DeepMix has a significant improvement of 2.3 percentage points of AUC and raises the ranking of DiMP from fifth to third compared with other trackers on OTB-2015. DiMP-DeepMix also achieves the top success score on LaSOT. ❸ DiMP collects historical samples’ embeddings and trains the predictor online. Our MixNet is designed for data augmentation, thus DeepMix has better compatibility with DiMP and achieves more improvement than DSiam and SiamRPN++.

**VOT2018.** We report results in Table 2. DeepMix with three trackers still is in effect. Unlike testing on other datasets, DiMP collects 250 samples for online training when testing on VOT2018 and costs extreme memory with DeepMix. Therefore, we still implement the same hyperparameters as evaluating on OTB-2015. Although it has achieved much lower score than its report, the results also prove the effectiveness of DeepMix. SiamRPN-DeepMix obtains a striking 0.057 improvement on EAO. Even though it uses MobileV2 as the backbone, achieves state-of-the-art on VOT2018.

## 4.3. Ablation study

**DeepMix with different backbones.** In order to validate the generality of DeepMix with different backbones, we present the result on OTB-2015 dataset in Table 3. It shows that: ❶ DeepMix can take stable effect for any network architecture. ❷ DeepMix has more improvement on ResNet18-based than ResNet50-based model. It probably because more powerful networks are less dependent on DeepMix.

**Validation of MixNet.** We implement a naive data augmentation method as the baseline to validate the effectiveness of the MixNet. That is, we calculate the filtering parameters, *i.e.*,  $\mathbf{W}_o$  and  $\mathbf{W}_b$ , by online optimizing an objective function via the gradient descent to replace the proposed MixNet. Specifically, we define an objective function that is the  $L_2$  distance between the predicted heat map and a Gaussian map having the highest score on the detected position. Then, at  $t$ th frame, we can minimize the objective function by tuning the  $\mathbf{W}_o$  and  $\mathbf{W}_b$  via the gradient descent for ten iterations. We denote this method as ‘DeepMix-Opt’ and compare it with the

**Table 4:** Comparing three variants of our method, *i.e.*, DeepMix-Opt, DeepMix-single and DeepMix, on DiMP tracker and OTB-2015 to validate the effectiveness of MixNet.

Metrics	AUC	Prec.	NormPrec.	FPS
DiMP	0.660	0.859	0.807	27.0
DiMP-DeepMix-Opt	<b>0.667</b>	<b>0.873</b>	<b>0.816</b>	<b>10</b>
DiMP-DeepMix-single	<b>0.676</b>	<b>0.884</b>	<b>0.831</b>	<b>26.5</b>
DiMP-DeepMix	<b>0.683</b>	<b>0.890</b>	<b>0.834</b>	<b>26.0</b>

final version DeepMix based on the DiMP tracker and OTB-2015 dataset. As shown in Table 4, DeepMix-Opt via online iterative optimization can also enhance the tracking accuracy but immensely increase the computational cost, slowing down the DiMP significantly.

As Fig. 2 shows, MixNet has two branches and outputs two filters ( $\mathbf{W}_o$  and  $\mathbf{W}_b$ ). We test another version of MixNet: keep only one branch and output one convolution kernel, then filters with samples  $\mathcal{X}_t$ , regardless of object or background. We term it as DeepMix-single. As shown in Table 4, DeepMix-single outperforms DiMP with a competitive speed, it still is weaker than the final version DeepMix. Therefore, learning different patterns of objects or backgrounds is an important method for online data augmentation. In contrast, DeepMix with the MixNet achieves much higher accuracy improvement with only 1 FPS speed decrease, demonstrating the effectiveness and advantages of the proposed MixNet.

## 5. CONCLUSION

In this work, we have taken a deep dive into the data augmentation aspect for improving online visual object tracking, a long-overlooked facet in this domain. Specifically, we have proposed the DeepMix as a complete pipeline that takes historical samples’ embeddings as input and generates augmented online, thus enhancing the state-of-the-art online learning methods for visual object tracking. To this end, we have proposed the *online data augmentation* for tracking that online augments the historical samples through object-aware filtering. Then, we have further proposed the *MixNet* which is an offline trained deep neural network for performing online data augmentation within one-step, for boosting the tracking accuracy while preserving high speeds of the state-of-the-art online learning methods. We have carried out extensive experiments on three different tracking frameworks, *i.e.*, DiMP, DSiam, and SiamRPN++, and on three large-scale and challenging datasets, *i.e.*, OTB-2015, LaSOT, and VOT. The experimental results have demonstrated and verified the effectiveness and advantages of the proposed method.

In the future, we plan to further study the online data augmentation for visual object tracking by considering different degradation, *e.g.*, motion blur [52], rain [53, 54], illumination variation [55, 56], *etc.*

**Acknowledgement.** This work was supported in part by the National Natural Science Foundation of China under Grant 61906135, 62020106004, and 92048301, Tianjin Science and Technology Plan Project under Grant 20JCQNJC01350, JSPS KAKENHI Grant

No.20H04168, 19K24348, 19H04086, and JST-Mirai Program Grant No.JPMJMI18BB and JPMJMI20B8, Japan, and Natural Science Foundation of Tianjin under Grant KJZ40420200017. This work was also supported by the Canada CIFAR AI program.

## 6. REFERENCES

- [1] 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, pp. 4282–4291.
- [2] Qing Guo, Wei Feng, Ce Zhou, Rui Huang, Liang Wan, and Song Wang, "Learning dynamic siamese network for visual object tracking," in *ICCV*, 2017, pp. 1763–1771.
- [3] Qing Guo, Xiaofei Xie, Felix Juefei-Xu, Lei Ma, Zhongguo Li, Wanli Xue, Wei Feng, and Yang Liu, "SPARK: Spatial-aware Online Incremental Attack Against Visual Tracking," in *ECCV*, 2020.
- [4] Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte, "Learning discriminative model prediction for tracking," in *ICCV*, 2019, pp. 6181–6190.
- [5] David S Bolme, J Ross Beveridge, Bruce A Draper, and Yui Man Lui, "Visual object tracking using adaptive correlation filters," in *CVPR*, 2010, pp. 2544–2550.
- [6] João F Henriques, Rui Caseiro, Pedro Martins, and Jorge Batista, "High-speed tracking with kernelized correlation filters," *IEEE TPAMI*, vol. 37, no. 3, pp. 583–596, 2014.
- [7] Martin Danelljan, Gustav Hager, Fahad Shahbaz Khan, and Michael Felsberg, "Learning spatially regularized correlation filters for visual tracking," in *ICCV*, 2015, pp. 4310–4318.
- [8] Ce Zhou, Qing Guo, Liang Wan, and Wei Feng, "Selective object and context tracking," in *ICASSP*, 2017, pp. 1947–1951.
- [9] Pengyu Zhang, Qing Guo, and Wei Feng, "Fast spatially-regularized correlation filters for visual object tracking," in *Pacific Rim International Conference on Artificial Intelligence*. Springer, Cham, 2018, pp. 57–70.
- [10] Wei Feng, Ruize Han, Qing Guo, Jianke Zhu, and Song Wang, "Dynamic saliency-aware regularization for correlation filter-based object tracking," *IEEE TIP*, vol. 28, no. 7, pp. 3232–3245, 2019.
- [11] Qing Guo, Ruize Han, Wei Feng, Zhihao Chen, and Liang Wan, "Selective spatial regularization by reinforcement learned decision making for object tracking," *IEEE TIP*, vol. 29, pp. 2999–3013, 2020.
- [12] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. S. Torr, "Fully-convolutional siamese networks for object tracking," in *arXiv preprint arXiv:1606.09549*, 2016.
- [13] Bo Li, Junjie Yan, Wei Wu, Zheng Zhu, and Xiaolin Hu, "High performance visual tracking with siamese region proposal network," in *CVPR*, 2018, pp. 8971–8980.
- [14] Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte, "Learning discriminative model prediction for tracking," in *ICCV*, 2019, pp. 6181–6190.
- [15] Zheng Zhu, Qiang Wang, Bo Li, Wei Wu, Junjie Yan, and Weiming Hu, "Distractor-aware siamese networks for visual object tracking," in *ECCV*, 2018, pp. 101–117.
- [16] 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, pp. 4277–4286.
- [17] Y. Wu, J. Lim, and M.-H. Yang, "Object tracking benchmark," *IEEE TPAMI*, vol. 37, no. 9, pp. 1834–1848, 2015.
- [18] 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 *ECCV*, 2018, pp. 0–0.
- [19] H. Fan, L. Lin, F. Yang, P. Chu, G. Deng, S. Yu, H. Bai, Y. Xu, C. Liao, and H. Ling, "Lasot: A high-quality benchmark for large-scale single object tracking," in *CVPR*, 2019, pp. 5369–5378.
- [20] Qing Guo, Wei Feng, Ruijun Gao, Yang Liu, and Song Wang, "Exploring the effects of blur and deblurring to visual object tracking," *IEEE TIP*, vol. 30, pp. 1812–1824, 2021.
- [21] Zedu Chen, Bineng Zhong, Guorong Li, Shengping Zhang, and Rongrong Ji, "Siamese box adaptive network for visual tracking," in *CVPR*, 2020, pp. 6668–6677.
- [22] Qing Guo, Wei Feng, Ce Zhou, and Bin Wu, "Structure-regularized compressive tracking," 2016.
- [23] Martin Danelljan, Luc Van Gool, and Radu Timofte, "Probabilistic regression for visual tracking," in *CVPR*, 2020, pp. 7183–7192.
- [24] Guangting Wang, Chong Luo, Xiaoyan Sun, Zhiwei Xiong, and Wenjun Zeng, "Tracking by instance detection: A meta-learning approach," in *CVPR*, 2020, pp. 6288–6297.
- [25] Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte, "Know your surroundings: Exploiting scene information for object tracking," *arXiv preprint arXiv:2003.11014*, 2020.
- [26] Tianyu Yang, Pengfei Xu, Runbo Hu, Hua Chai, and Antoni B Chan, "Roam: Recurrently optimizing tracking model," in *CVPR*, 2020, pp. 6718–6727.
- [27] Sepp Hochreiter and Jürgen Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [28] M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg, "Adaptive decontamination of the training set: A unified formulation for discriminative visual tracking," in *CVPR*, 2016.
- [29] M. Danelljan, G. Bhat, F. S. Khan, and M. Felsberg, "Eco: Efficient convolution operators for tracking," in *CVPR*, 2017, pp. 6931–6939.
- [30] Lichao Zhang, Abel Gonzalez-Garcia, Joost van de Weijer, Martin Danelljan, and Fahad Shahbaz Khan, "Learning the model update for siamese trackers," in *ICCV*, 2019, pp. 4010–4019.
- [31] Dan Ciregan, Ueli Meier, and Jürgen Schmidhuber, "Multi-column deep neural networks for image classification," in *CVPR*, 2012, pp. 3642–3649.
- [32] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.
- [33] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger, "On calibration of modern neural networks," *arXiv preprint arXiv:1706.04599*, 2017.
- [34] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le, "Autoaugment: Learning augmentation policies from data," *arXiv preprint arXiv:1805.09501*, 2018.
- [35] Daniel Ho, Eric Liang, Xi Chen, Ion Stoica, and Pieter Abbeel, "Population based augmentation: Efficient learning of augmentation policy schedules," in *ICML*, 2019, pp. 2731–2741.
- [36] Sungbin Lim, Ildoo Kim, Taesup Kim, Chiheon Kim, and Sungwoong Kim, "Fast autoaugment," in *NeurIPS*, 2019, pp. 6665–6675.
- [37] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan, "Augmix: A simple data processing method to improve robustness and uncertainty," *arXiv preprint arXiv:1912.02781*, 2019.
- [38] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz, "mixup: Beyond empirical risk minimization," *arXiv preprint arXiv:1710.09412*, 2017.
- [39] Sangdoon Yoo, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo, "Cutmix: Regularization strategy to train strong classifiers with localizable features," in *ICCV*, 2019, pp. 6023–6032.
- [40] Dan Hendrycks\*, Norman Mu\*, Ekin Dogus Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan, "Augmix: A simple method to improve robustness and uncertainty under data shift," in *ICLR*, 2020.
- [41] Janghoon Choi, Junseok Kwon, and Kyoung Mu Lee, "Deep meta learning for real-time target-aware visual tracking," in *ICCV*, 2019, pp. 911–920.
- [42] Peixia Li, Boyu Chen, Wanli Ouyang, Dong Wang, Xiaoyun Yang, and Huchuan Lu, "Gradnet: Gradient-guided network for visual object tracking," in *ICCV*, 2019, pp. 6162–6171.
- [43] Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg, "Atom: Accurate tracking by overlap maximization," in *CVPR*, 2019, pp. 4655–4664.
- [44] Ning Wang, Wengang Zhou, Guojun Qi, and Houqiang Li, "Post: Policy-based switch tracking," in *AAAI*, 2020, pp. 12184–12191.
- [45] H. Fan and H. Ling, "Siamese cascaded region proposal networks for real-time visual tracking," in *CVPR*, 2019, pp. 7944–7953.
- [46] Kenan Dai, Huchuan Lu Dong Wang, Chong Sun, and Jianhua Li, "Visual tracking via adaptive spatially-regularized correlation filters," in *CVPR*, 2019.
- [47] Qiang Wang, Li Zhang, Luca Bertinetto, Weiming Hu, and Philip H.S. Torr, "Fast online object tracking and segmentation: A unifying approach," in *CVPR*, 2019, pp. 1328–1338.
- [48] Zhipeng Zhang and Houwen Peng, "Deeper and wider siamese networks for real-time visual tracking," in *CVPR*, 2019.
- [49] Matej Kristan, Jiri Matas, Ales Leonardis, Michael Felsberg, Luka Cehovin, Gustavo Fernandez, Tomas Vojir, Gustav Hager, Georg Nebehay, and Roman Pflugfelder, "The visual object tracking vot2015 challenge results," in *ICCV*, 2015, pp. 1–23.
- [50] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Identity mappings in deep residual networks," in *ECCV*, 2016, pp. 630–645.
- [51] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *CVPR*, 2018, pp. 4510–4520.
- [52] Qing Guo, Felix Juefei-Xu, Xiaofei Xie, Lei Ma, Jian Wang, Bing Yu, Wei Feng, and Yang Liu, "Watch out! motion is blurring the vision of your deep neural networks," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [53] Liming Zhai, Felix Juefei-Xu, Qing Guo, Xiaofei Xie, Lei Ma, Wei Feng, Shengchao Qin, and Yang Liu, "It's raining cats or dogs? adversarial rain attack on dnn perception," *arXiv preprint arXiv:2009.09205*, 2020.
- [54] Qing Guo, Jingyang Sun, Felix Juefei-Xu, Lei Ma, Xiaofei Xie, Wei Feng, and Yang Liu, "Efficientderain: Learning pixel-wise dilation filtering for high-efficiency single-image deraining," in *AAAI*, 2021.
- [55] Lan Fu, Changqing Zhou, Qing Guo, Felix Juefei-Xu, Hongkai Yu, Wei Feng, Yang Liu, and Song Wang, "Auto-exposure fusion for single-image shadow removal," *arXiv preprint arXiv:2103.01255*, 2021.
- [56] Binyu Tian, Qing Guo, Felix Juefei-Xu, Wen-Le Chan, Yupeng Cheng, Xiaohong Li, Xiaofei Xie, and Shengchao Qin, "Bias field poses a threat to dnn-based x-ray recognition," *arXiv preprint arXiv:2009.09247*, 2020.