

Unsupervised Cross-Modal Distillation for Thermal Infrared Tracking

Jingxian Sun^{*†}

ASGO, School of Computer Science,
Northwestern Polytechnical
University
Xi'an, China
jingxiansun@mail.nwpu.edu.cn

Lichao Zhang^{*}

Aeronautics Engineering College, Air
Force Engineering University
Xi'an, China
lichao.zhang@outlook.com

Yufei Zha^{‡†}

ASGO, School of Computer Science,
Northwestern Polytechnical
University
Xi'an, China
yufeizha@nwpu.edu.cn

Abel Gonzalez-Garcia

wrnch
Montreal, Canada
abel.gonzalezgarcia@wrnch.ai

Peng Zhang[†]

ASGO, School of Computer Science,
Northwestern Polytechnical
University
Xi'an, China
zh0036ng@nwpu.edu.cn

Wei Huang

School of Information Engineering,
Nanchang University
Nan Chang, China
huangwei@ncu.edu.cn

Yanning Zhang[†]

ASGO, School of Computer Science,
Northwestern Polytechnical
University
Xi'an, China
ynzhang@nwpu.edu.cn

ABSTRACT

The target representation learned by convolutional neural networks plays an important role in Thermal Infrared (TIR) tracking. Currently, most of the top-performing TIR trackers are still employing representations learned by the model trained on the RGB data. However, this representation does not take into account the information in the TIR modality itself, limiting the performance of TIR tracking.

To solve this problem, we propose to distill representations of the TIR modality from the RGB modality with Cross-Modal Distillation (CMD) on a large amount of unlabeled paired RGB-TIR data. We take advantage of the two-branch architecture of the baseline tracker, *i.e.* DiMP, for cross-modal distillation working on two components of the tracker. Specifically, we use one branch as a teacher module to distill the representation learned by the model into the other branch. Benefiting from the powerful model in the RGB modality,

the cross-modal distillation can learn the TIR-specific representation for promoting TIR tracking. The proposed approach can be incorporated into different baseline trackers conveniently as a generic and independent component. Furthermore, the semantic coherence of paired RGB and TIR images is utilized as a supervised signal in the distillation loss for cross-modal knowledge transfer. In practice, three different approaches are explored to generate paired RGB-TIR patches with the same semantics for training in an unsupervised way. It is easy to extend to an even larger scale of unlabeled training data. Extensive experiments on the LSOTB-TIR dataset and PTB-TIR dataset demonstrate that our proposed cross-modal distillation method effectively learns TIR-specific target representations transferred from the RGB modality. Our tracker outperforms the baseline tracker by achieving absolute gains of 2.3% Success, 2.7% Precision, and 2.5% Normalized Precision respectively. Code and models are available at <https://github.com/zhanglichao/cmdTIRtracking>.

^{*}Both authors contributed equally to this research.

[†]National Engineering Laboratory for Integrated Aero-Space-Ground-Ocean Big Data Application Technology (ASGO).

[‡]Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

MM '21, October 20–24, 2021, Virtual Event, China

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8651-7/21/10...\$15.00

<https://doi.org/10.1145/3474085.3475387>

CCS CONCEPTS

• Computing methodologies → Unsupervised learning; Tracking; Neural networks.

KEYWORDS

Unsupervised learning; Knowledge distillation; TIR tracking; Convolutional neural network

ACM Reference Format:

Jingxian Sun, Lichao Zhang, Yufei Zha, Abel Gonzalez-Garcia, Peng Zhang, Wei Huang, and Yanning Zhang. 2021. Unsupervised Cross-Modal Distillation for Thermal Infrared Tracking. In *Proceedings of the 29th ACM*

1 INTRODUCTION

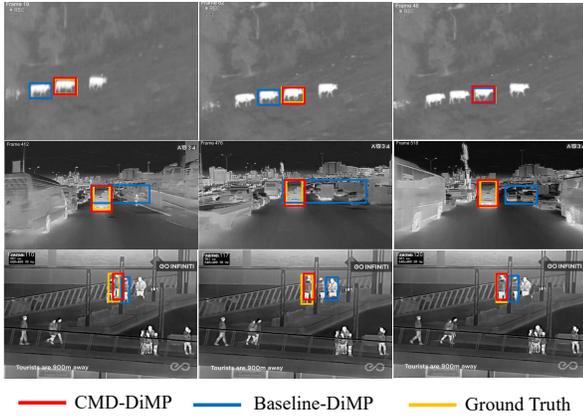


Figure 1: Qualitative comparison between the proposed method and baseline tracker DiMP [4] on LSOTB-TIR dataset [33]. By training with cross-modal distillation, the tracker can effectively track the targets on different challenging scenes.

Thermal Infrared (TIR) tracking [15, 17, 31, 34] aims to locate a target by using videos of the thermal infrared modality, where the initial location of the target is given in the first frame. Compared with RGB trackers, TIR trackers can still distinguish the target from the background under some challenging situations, *e.g.* low illumination, shadow, and occlusion, even working well in total darkness where visual cameras have no signal. These advantages make TIR trackers running in a wide range of applications, such as video surveillance, maritime rescue, various defense systems, and driver assistance at night [26]. However, one key issue of TIR trackers is how to learn powerful features for representing the target efficiently in order to deal with various difficulties specific to the TIR modality, such as thermal crossover, intensity variation, and distractors.

In early work, the hand-crafted features [15, 44], *e.g.* Histogram of Oriented Gradients (HOG) [7], Harr-like feature [29], are employed to represent targets in TIR trackers. These trackers exploit classic learning paradigms, such as multiple instance learning [?], discriminative correlation filter [1], and low-rank sparse learning [19] for TIR tracking. Recently, Convolutional Neural Networks (CNN) have also been introduced to represent the target to improve the performance of TIR trackers [31, 32]. The CNN features extracted by the pre-trained networks, *e.g.* VGGNet [38], ResNet [18], are integrated into the trackers followed by the correlation filter (CF) [20], structural support vector machine [16] or Siamese networks [3] for enhancing the target representation. The results show that the representation with CNN features is discriminative to distinguish the target from distractors in the background, compared with hand-crafted features.

However, both hand-crafted and off-the-shelf features are not optimal for the TIR target, which limits the performance of TIR tracking. In practice, TIR images lack color information and rich

texture features compared with RGB images. The hand-crafted features are designed according to the characteristics of RGB images, while the pre-trained models are derived from large-scale RGB data. They have exclusively used the RGB modality, ignoring the differences between RGB and TIR modalities. As a result, this gap in the representation degrades the discriminability of the tracker for identifying the target from the background in the TIR modality. A TIR-specific representation needs to be specifically tailored for the TIR modality in order to maximally leverage its characteristics. Additionally, the lack of large-scale annotated TIR data makes it impossible to train the networks from scratch. Labeling the TIR data is a time-consuming and laborious work. And at present there is no large-scale TIR data for tracking as normally RGB trackers do [3, 9, 14, 22]. Therefore, in our previous work [45], we collect a large amount of TIR data from other vision tasks, but these TIR data is not annotated for the tracking task. To solve this issue, we propose to use image-to-image translation models to generate synthetic labeled TIR datasets transferring from RGB tracking datasets. Exhaustedly, a lot of extra efforts are still needed for training translation models in this work.

In order to obtain powerful target representations for TIR tracking, we propose to distill the representation of the TIR modality from the RGB modality with Cross-Modal Distillation (CMD) on a large amount of unlabeled paired visible and infrared images. Motivated by the idea of distilling the network knowledge from a teacher model to a student model [21], our method does distill representation knowledge from the RGB modality to the TIR modality. We use DiMP [4] as our baseline tracker which is constructed with the architecture of two branches. We explore distillation operations on Target Center Location (TCL) and Bounding Box Estimation (BBE) in DiMP [4]. As a result, we can obtain the TIR-specific representation guided by the pre-trained model in the RGB modality. Benefiting from the powerful model trained on large-scale labeled RGB data, the learned TIR model can better represent the TIR target. Here, the proposed method is generic and can be applied to different baseline trackers conveniently. In this work, we use our cross-modal distillation method to train the tracker ATOM [9] in section 4.4, and it effectively improves the performance of ATOM [9]

Moreover, an unsupervised training method is proposed to take advantage of the dual-modalities data without any annotations, and this will relieve the dependency on the labeled TIR data. The semantic coherence of the paired RGB and TIR image replaces the manual labels as the ground-truth in the final loss function for model training. This prior is helpful to transfer the model knowledge between the different modalities. In practice, we explore three different approaches (‘center area’, ‘random sampling’, and ‘detection’) to generate paired RGB-TIR patches from the RGB and TIR images as training data. These paired dual-modalities image patches with the same semantics are fed into the network to learn TIR representation under the distillation loss. Here, we do not require any kinds of annotations in the training data, so the whole training procedure can be conducted in an unsupervised manner. Besides, it is easy to extend to an even larger scale of unlabeled training data.

We validate the proposed method on two standard test datasets: LSOTB-TIR dataset [33] and PTB-TIR dataset [30]. Some qualitative results for comparisons between the proposed method and baseline tracker are shown in Fig. 1. Compared with the baseline tracker, we

achieve absolute gains of 2.3% Success, 2.7% Precision and 2.5% Normalized Precision respectively on LSOTB-TIR dataset [33]. These results demonstrate that our proposed Cross-Modal Distillation (CMD) method effectively learns TIR-specific target representations transferred from the RGB modality. The contributions of our work are as follows:

- A representation transferring approach called Cross-Modal Distillation (CMD) is proposed to distill a TIR-specific representation from the RGB modality on a large amount of unlabeled paired RGB-TIR data. This benefits from the powerful model trained on large-scale labeled RGB data. The proposed approach can be incorporated into different baseline trackers conveniently due to its generality and independence.
- An unsupervised manner is proposed without any annotation of the target for training. During training, three different approaches are explored to generate paired RGB-TIR patches with the same semantics automatically. It is easy to extend to an even larger scale of unlabeled training data.
- We conduct extensive experiments on two benchmarks to verify the effectiveness of the proposed method. The results demonstrate that our algorithm achieves a significant improvement against SOTA methods on the TIR tracking challenge.

The remainder of the paper is structured as follows. In section 2, we briefly discuss related works. In section 3, we describe the cross-modal distillation modules to learn the TIR-specific representation transferred from the pre-trained model on the RGB data. In section 4, extensive experiments are carried out on two standard thermal infrared tracking datasets. Finally, we conclude our work and propose future research plans in section 5.

2 RELATED WORK

In this section, we will introduce the works closely related to our study in this paper. More references about multi-modal tracking can be seen in the surveys [46, 47].

TIR Tracking. The hand-crafted features, such as edges, motion features, and HOG, were integrated discriminative correlation filter with scale estimation [17] or spatial regularization [15] for TIR tracking. Their favorable performance was mainly due to the robust feature representation and online learning. The CNN features extracted from the common networks were used to replace the hand-crafted features for target representation in TIR tracking [34]. Recently, MLSSNet [32] trained a multi-level similarity-based Siamese network on an RGB and TIR dataset simultaneously. A multi-task matching framework [31] was proposed to learn deep features in the levels of inter-class and intra-class respectively. These kinds of features complemented each other and recognized TIR objects in the levels of inter-class and intra-class respectively. These feature models were learned and jointly optimized on the TIR tracking task. Besides, the Siamese networks were used to extract the features from the network trained on a large amount of synthetic TIR images for TIR tracking [45].

Unlike previous works, we expect to learn the TIR-specific discriminative representation transferred from the RGB modality by training the network on the paired RGB and TIR images.

Knowledge Distillation for Vision Tasks. The main idea of knowledge distillation [21] was that the student model mimics the

teacher model in order to obtain a competitive or even a superior performance, which benefits the deployment of deep neural networks in mobile devices and embedded systems. In fact, it was important to transfer knowledge between different modalities, because the data or labels of some modalities might not be available during training. The idea is to transfer the annotation or label data through pair-wise samples and this has been widely used for cross-modal applications [40]. Tian *et al.* [40] proposed a contrastive loss to transfer pair-wise relationship across different modalities, while GANs were employed to perform cross-modal distillation among the missing and available modalities [37]. In addition, in the field of visual question answering, the knowledge from trilinear interaction teacher model with image-question-answer as inputs was distilled into the learning of a bilinear interaction student model with image-question as inputs [12]. Besides, lots of cross-modal distillation methods [6, 24] also transferred the knowledge among multiple domains.

Unlike the above applications, we focused on TIR tracking, which lacked large-scale annotated TIR data for training. To overcome this problem, cross-modal distillation (CMD) was introduced to transfer representations from the RGB modality to the TIR modality in this study. Specifically, the knowledge of the RGB model was transferred to the TIR model through unsupervised learning, and it benefited to improve the performance of the TIR tracking task.

3 PROPOSED METHOD

In this section, we propose to transfer the representation of the pre-trained model from RGB modality to TIR modality by using a large amount of paired RGB-TIR data in an unsupervised way. Unlike the classic KD work in [21] which uses two independent models for distillation happening on intermediate features or final outputs, our proposed cross-modal distillation adopts representations of the RGB modality as powerful supervision signals in one branch to guide the representation learning in the TIR modality with the other branch.

3.1 Overview

The training pipeline of the proposed unsupervised training of cross-modal distillation for TIR tracking is shown in Figure 2. In order to transfer the representation from the RGB modality to the TIR modality, cross-modal distillation modules are discussed in Fig. 2 (e) following after the backbone network shown in Fig. 2 (d).

Here, the proposed cross-modal distillation modules are constructed by the convolutional layer, pooling layer, fully connected layer, distillation operation, and so on and contains two components for distillation: the Target Center Location (TCL) distillation for discriminating the target from the background and Bounding Box Estimation (BBE) distillation for fine-tuning the bounding box of the target. Both these two distillations are guided by representations of the RGB modality from one branch to learn the TIR-specific representations with the other branch under the distillation loss. Unlike the training procedure of RGB trackers, our approach is trained to process the information from both RGB and TIR modalities simultaneously, which are denoted as red and blue lines respectively in Fig. 2. While the green line represents the procedure of the cross-modal distillation in Fig. 2. The input of the proposed model architecture are paired dual-modalities images. Before the

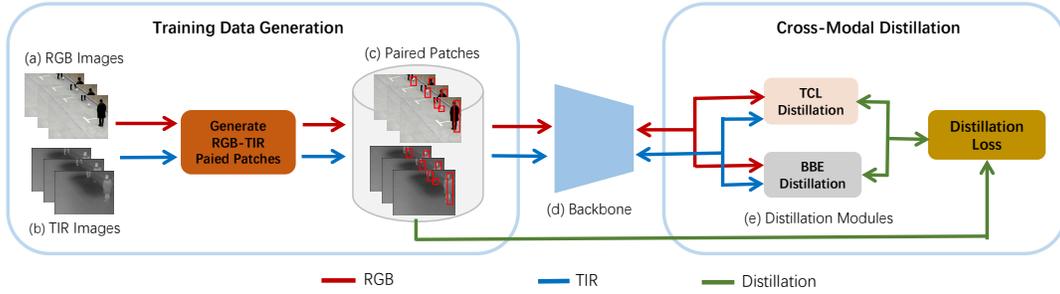


Figure 2: Unsupervised training of Cross-Modal Distillation (CMD) for TIR tracking. Only the paired RGB images shown in Fig. 2 (a) and TIR images shown in Fig. 2 (b) are employed as training data. Before network training, the paired dual-modalities patches in Fig. 2 (c) are generated automatically as the input of the network for cross-modal distillation. The CNN features are extracted by backbone network on both RGB and TIR modalities in Fig. 2 (d), and then they are fed into the cross-modal distillation shown in Fig. 2 (e) with target center location (TCL) distillation and bounding box estimation (BBE) distillation. Here, the red line and blue line denote the processing flow of RGB modality and TIR modality, respectively. The cross-modal distillation flow is expressed as the green line. The line with the reverse arrow indicates the training process.

training of the network, three different approaches are explored to generate paired patches shown in Fig. 2 (c) and they are from the RGB images as in Fig. 2 (a) and TIR images as in Fig. 2 (b). After that, patches’ coordinates in the image and pseudo Gaussian distribution are utilized to construct the annotated label for calculating the loss during the training of the cross-modal distillation.

Mathematically, we denote the RGB image as I^v and TIR image as I^t , respectively. The bounding box is described as $s = (c_x/w, c_y/h, \log w, \log h) \in \mathbb{R}^4$, where (c_x, c_y) is image coordinates of the bounding box center, the width and height of box are denoted as w and h . The CNN features extracted by the backbone network are written as $x^v \in \chi$ for RGB modality and $x^t \in \chi$ for TIR modality. The training dataset can be denoted as $S_{train} = \{x_i^v, x_i^t\}_{i=1}^N$, where N is the volume of the training dataset.

Our goal is to learn a cross-modal distillation $\mathbb{D}(\phi, \theta)$ to achieve TIR-specific representations by transferring the RGB model. Specifically, the target center location (TCL) distillation ϕ is used to identify the target from the distractors in the background for locating the position. Furthermore, to refine the bounding box of the results, the bounding box estimation (BBE) distillation θ is used to learn the regression coefficients to fit the target more accurately. Thus, the final distillation loss contains two parts: target center location loss and bounding box estimation loss, which can be written as:

$$\mathcal{L}(\phi, \theta) = \mathcal{L}_{TCL}(\phi) + \lambda \mathcal{L}_{BBE}(\theta), \quad (1)$$

where λ is a regularization coefficient to balance the two losses.

3.2 Cross-Modal Distillation

In this section, we introduce distillation modules between RGB and TIR modalities based on a large amount of paired RGB-TIR data.

Paired Patches Generation: The large amounts of annotated data are an essential factor for training network [11] that has been demonstrated in other vision tasks. But for TIR tracking, lacking of large-scale training data limits the training of the network from scratch. Recent works [31, 32] directly employ the pre-trained model in the RGB modality for the representations of TIR targets. The results show that this off-the-shelf representation is not optimal for the TIR target, as the appearance of the target in different modalities varies largely. The reason owes to the different imaging mechanisms.

Unlike the previous works [31, 32], in this work, we consider using a large amount of unlabeled paired data from RGB and TIR modalities in this study. This avoids a lot of time and efforts to label them for the acquisition of large amounts of infrared data. To better and more effectively utilize the paired RGB-TIR data for training, three methods: ‘center area’, ‘random sampling’, and ‘detection’, are explored to generate the paired patches with same semantics as the input of the network. Specifically, the ‘center area’ is to assign a square region in the center of the image as a fake target, thus obtaining the bounding box in this image. The ‘random sampling’ is to randomly sample several patches from the image and then they can be fed into the network for training.

Additionally, we also employ the object detectors, such as YOLO [5], Faster R-CNN [36], to detect the locations of the targets in the image. After we obtain bounding boxes of ‘objects’ by the above three methods, we can feed them for two applications in the training procedure. On one hand, the paired patches from dual-modalities are used as input of the network to execute forward inference. On the other hand, the corresponding coordinates of the patches are used to construct the distillation losses for back-propagation training.

Target Center Location Distillation: The target center location (TCL) distillation is utilized to locate the target coarsely by predicting the score of heatmap in the image. This distillation pays more attention to robustness than accuracy during tracking. In this distillation, we expect to learn a TIR filter ϕ derived from the TIR images while fitting with the RGB representation with a Gaussian pseudo distribution g . The cross-modal distillation (CMD) loss of target center location can be written as follows:

$$\mathcal{L}_{TCL}(\phi) = \sum_{i=1}^N \|\phi(x_i^t) \otimes x_i^v - g\|^2 + \mu \|\phi\|^2, \quad (2)$$

where the filter is denoted as ϕ , \otimes is the convolution operation and μ is a regularization parameter.

This loss can be minimized to optimize the weights of a linear convolutional layer. It is helpful to identify the target from distractors in the background. During training, this branch is trained by a meta-learning way with the pre-trained RGB model using the above loss in Eq. 2. During the tracking inference, the center position of

the target is optimized by searching the maximum confidence score within a wide search region in the next frame.

Bounding Box Estimation Distillation: The BBE enables the tracker to be wrapped by the box for accuracy improvement of the tracking performance. Here, we denote the state of the bounding box as $s = (c_x/w, c_y/h, \log w, \log h) \in \mathbb{R}^4$. The distillation loss of bounding box estimation can be denoted as:

$$\mathcal{L}_{BBE}(\theta) = \sum_{i=1}^N \|\psi(\theta(x_i^t)) \odot \varphi(x_i^v) - s\|^2 + \nu \|\theta\|^2, \quad (3)$$

where $\psi(\cdot)$ is the fully connection operation, \odot is the pixel-wise multiplication, and φ is pre-trained model in RGB modality.

We obtain the RGB representation by using the φ to extract on the RGB features x^v , and regard it as the teacher vector to guide the student model θ for learning TIR-specific representations under the pseudo label s . During tracking, the BBE distillation is used to find the bounding box fitting the target ultimately by maximizing the intersection over union (IoU) scores.

Unsupervised Training: In this work, we propose a representation transferring mechanism for two-branch trackers by cross-modal distillation learning. Different from previous training approaches, the proposed training mechanism can work on a large amount of unlabeled RGB-TIR data with an unsupervised training manner. Then RGB and TIR CNN features extracted by the backbone network are utilized as the input of the Cross-Modal Distillation (CMD) for representation distillation. Specifically, RGB and TIR features extracted by the backbone are utilized as the input of the cross-modal Distillation (CMD) module. For the distillation mechanism, the RGB representation obtained by the pre-trained model plays the role of the teacher model, while the model for learning the TIR representation is the student model.

As for the consistency of the semantics information between the paired RGB and TIR patches, the distillation training enables the TIR representation of the student model approachable to the RGB representation of the teacher model. Obviously, the training procedure mainly relies on the consistent information of the paired patches, thus the network can be trained in an unsupervised way. That is to say, our proposed method only needs the clean unlabeled paired dual-modalities images, avoiding the time-consuming, labor-intensive, and cumbersome manually labeled data. Furthermore, it is easy to extend to an even larger scale of unlabeled training data.

4 EXPERIMENTS

In this section, we provide the experimental results of the proposed CMD method on LSOTB-TIR dataset [33] and PTB-TIR dataset [30] to verify its effectiveness. Besides, we compare the trackers equipped with our CMD with several state-of-the-art trackers.

4.1 Evaluation Datasets and Protocols

LSOTB-TIR dataset [33] is a large-scale high-diversity TIR tracking benchmark with a total of 1,400 TIR sequences and more than 600K frames. It is annotated with more than 730K bounding boxes in total. The training dataset contains 1,280 sequences with 47 objects classes and over 650k bounding boxes. And it selects 120 sequences, with 22 object classes and more than 82K frames, as the evaluation dataset. At present, it is larger and more diverse than

other existing TIR datasets. We use the Precision, the Normalized Precision and the Success as the metrics for this evaluation dataset. **PTB-TIR dataset [30]** is a TIR pedestrian tracking dataset for the TIR pedestrian tracker evaluation, which includes 60 thermal sequences with manual annotations. Each sequence has nine attribute labels for the attribute based evaluation to ensure the diversity of the dataset, and all of them come from different devices, scenes, and shooting times. The center location error (CLE) and overlap ratio (OR) are exploited as metrics. That is to say, the Precision Plot and Success Plot are used to rank trackers.

Evaluation protocols. We use one-pass evaluation method [43] (OPE) that each tracker is only initialized in the first frame, and is not affected by the true position of the target during the entire tracking process. The performance of the algorithm is evaluated by precision, normalized precision and success rates. The center location error (CLE) refers to the Euclidean distance between the center of the predicted position and the center of the artificial mark. The Precision is the ratio of the number of video frames whose CLE is less than a given threshold to the total number of video frames. A threshold of 20 pixels is usually set as the sorting criterion. As the Precision is affected by the image’s resolution and the size of the bounding box, we normalize the Precision over the size of the bounding box. Then, the area under the curve (AUC) of the Normalized Precision between 0 and 0.5 is used to evaluate the performance of the trackers. The overlap rate (OR) is the ratio of the union and intersection of the predicted target area and the ground truth area. The Success is the ratio of the number of frames with an overlap rate greater than a set threshold ([0, 1]) to the total number of frames. We usually use the area under the curve (AUC) as an indicator to measure the overall effectiveness of the tracking algorithm.

4.2 Implementation Details

For the baseline tracker DiMP [4], we use the default settings, referring to details in the paper [4]. Here, we update Target Center Location (TCL) and Bounding Box Estimation (BBE) modules to carry out the cross-modal distillation during training. TIR images lack details and texture information compared with RGB images, while the TIR target’s appearance changes stably during tracking. Thus, we need to update the tracking model slightly and carefully to adapt to the TIR characteristics. In practice, we reduce the learning rates used to update BBE and TCL, both to $1e-6$ as the rows **B** & **C** in Table 1. Considering the difference of convergence rates between BBE and TCL, we set the learning rates of BBE and TCL as $1e-7$ and $2e-8$ respectively, during the joint distillation as the row **D** in Table 1. For fair comparison, the model is trained for 50 epochs with mini-batch size of 5, and the learning rate is decreased to multiples of 0.5 at every 15 epochs in all experiments. The number of samples in every epoch is 26,000.

From the training datasets, we sample the paired RGB and TIR patches with the same semantics. RGB patches are fed into one branch of the network and the corresponding TIR patches are input for the other branch. Specifically, we input 5 RGB images and 5 TIR images for each branch. In addition, in order to enhance the mutual information between the modalities during the cross-modal distillation, we connect the paired RGB and TIR patches along the horizontal direction in spatial domain as the input of the network.

Table 1: Analysis of our cross-modal distillation on LSOTB-TIR [33] dataset. We evaluate several variants of our proposed method based on DiMP [4]. The best results are highlighted in bold font.

Modules to be Updated	Training Settings	Reference Branch	Test Branch	Success(↑)	Precision(↑)	Normalized Precision(↑)
None(Baseline)	A1. Same as DiMP [4]	RGB	RGB	66.2 (0.0)	78.7 (0.0)	70.7 (0.0)
BBE	B1. Learning on reference branch	TIR(ft)	RGB	67.4 (1.2)	80.4 (1.7)	72.1 (1.4)
	B2. Learning on test branch	RGB	TIR(ft)	67.2 (1.0)	80.1 (1.4)	71.9 (1.2)
	B3. Combination of B1 and B2	TIR	TIR	67.8 (1.6)	80.8 (2.1)	72.5 (1.8)
TCL	C1. Learning on reference branch	TIR(ft)	RGB	67.7 (1.5)	80.7 (2.0)	72.5 (1.8)
BBE and TCL	D1. Joint learning on test branch	RGB	TIR(ft)	67.1 (0.9)	79.7 (1.0)	71.8 (1.1)
	D2. Joint learning on reference branch	TIR(ft)	RGB	67.6 (1.4)	80.6 (1.9)	72.3 (1.6)
	D3. Joint learning	RGB-TIR(ft)	TIR-RGB(ft)	68.0 (1.8)	80.8 (2.1)	72.7 (2.0)
	D4. Joint learning with random sampling	RGB-TIR(ft)	TIR-RGB(ft)	67.3 (1.1)	80.2 (1.5)	72.1 (1.4)
	D5. Joint learning with a detector	RGB-TIR(ft)	TIR-RGB(ft)	67.5 (1.3)	80.4 (1.7)	72.3 (1.6)

The connected patch ‘RGB-TIR’ means that the RGB patch is in the left side of the TIR patch, while vice versa is called ‘TIR-RGB’. Then, they are fed to the two branches for distilling the representation. In the next part, we describe our three approaches to generate the paired patches which are aligned strictly with the same semantic information from above paired RGB and TIR images.

Paired Patches Generation. We use the training data from the work [45] which takes advantage of large-scale paired RGB-TIR data by collecting from several datasets, e.g. KAIST dataset [23]. The whole training data consists of 126,666 paired RGB and TIR images which are unlabeled for tracking.

We propose three approaches to implement the generation of paired patches from the above training data. The first is to assign central regions of the paired images, which is regarded as a fake object, to be the input paired patches. We call this method as ‘center area’. The second is to randomly crop some regions of approximately 1/6 size of the image, and then resize these regions to 100×100 pixels as the patches. We call this method as ‘random sampling’. In order to obtain more accurate location of objects on RGB-TIR datasets, we use an object detector as the third method called ‘detection’. In this way, the size of the paired patches is dependent on the detection results. Specifically, as the detector is trained in RGB modality, it is suitable for the RGB detection. So, we use it to detect objects of RGB images.

4.3 Analysis of Distillation Mechanisms

For the training settings in DiMP [4], the inputs of reference and test branches are temporal patches from the same sequence. That is to say, they aim to make full use of the target’s change in the single modality based on the continuity of time. Different from that, our goal is to learn cross-modal knowledge by transferring the high-level semantics of the same object in the RGB modality to the TIR modality. After the distillation, the input of the tracking process is only from the TIR modality.

Table 1 shows our analysis of the effectiveness of cross-modal distillation to Target Center Location (TCL) and Bounding Box Estimation (BBE) over the baseline tracker DiMP [4]. The tracking results of the baseline tracker are presented as **A1**.

All the experiments mainly focus on updating the parameters in the branch which is input with TIR patches. We train two modules for the DiMP [4], including the BBE and the TCL. The results are reported in terms of Success, Precision, and Normalized Precision. We explicitly show the patches with the corresponding modality attribute for the reference branch and the test branch of the tracker. Here, ‘RGB’ and ‘TIR’ mean that the patch is from a single modality.

Then, they are mixed to do cross-correlation in the Siamese architecture. For example, for training the reference branch, we keep the upper half and the lower half of a mini-batch as patches from RGB modality and TIR modality respectively, namely ‘RGB-TIR’ in the table. Thus the corresponding places in mini-batch for test branch are input with patches from TIR modality and RGB modality respectively, namely ‘TIR-RGB’ in the table. ‘ft’ means parameters of the branch are to be fine-tuned in the network.

Bounding Box Estimation (BBE). For the BBE, we consider three cases, namely ‘B1. Learning on reference branch’, ‘B2. Learning on test branch’ and ‘B3. Combination of B1 and B2’.

- **B1** We input cross-modal patches to the Siamese branches of BBE. Specifically, we input the TIR patches to the reference branch and input the RGB patches to the test branch in the BBE, separately. We only update the parameters in the model of the TIR branch, namely the reference branch. Thus, we use the well-trained model in the tracking process. Compared with the original model, only the reference branch is updated. We input the TIR patches from the testing dataset [33] to the well-trained model. This method improves 1.2% in Success, 1.7% in Precision, and 1.4% in Normalized Precision. These results indicate that the reference branch needs to be trained to fit the TIR modality to obtain a better representation.
- **B2** We flip the input of patches from the two modalities. Specifically, the TIR patches are input to the test branch and the RGB patches are input to the reference branch. Therefore, the parameters of the test branch in the model are updated by the distillation. Similarly, as the **B1** method, we only update the parameters of the TIR branch. Both of the two methods improve the performance of the baseline tracker. Transfer between RGB and TIR modalities for BBE is useful and any branch of BBE has its own function. Besides, the reference branch trained in **B1** plays a little better than the test branch trained in **B2** with 0.2% in terms of Success, 0.3% in terms of Precision, and 0.2% in terms of Normalized Precision.
- **B3** We combine the information from the two modalities. In practice, we borrow the parameters of the reference branch trained in **B1** and also the parameters of the test branch trained in the **B2**. Then we recombine them together as the complete BBE. We achieve the best result and our tracker outperforms baseline tracker 1.6%, 2.1% and 1.8% in terms of Success, Precision, and Normalized Precision respectively.

The results in case (**B1**) achieves better results than that in case (**B2**). We think the reason is that the reference branch contains

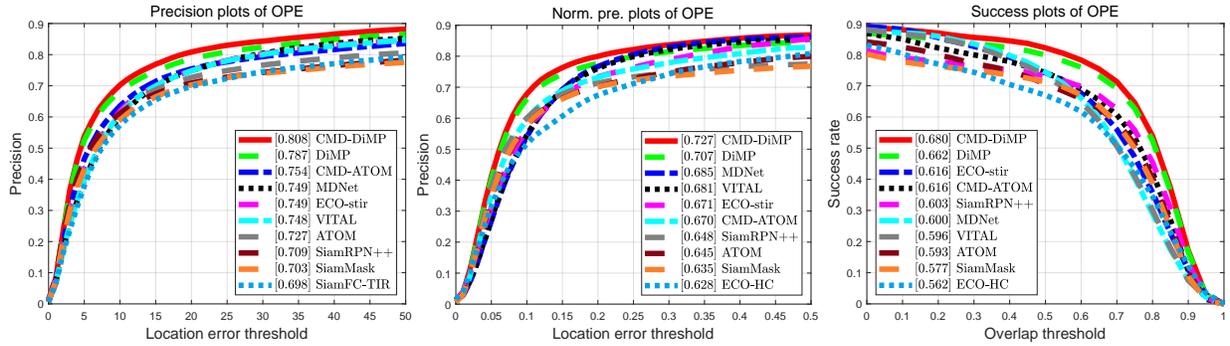


Figure 3: State-of-the-art comparison on LSOTB-TIR dataset [33]. We compare our method with top-10 trackers in terms of Precision, Normalized Precision, and Success. Our proposed approach achieves the best performances on three metrics.

more operations and can transfer more representation knowledge from the RGB model than the test branch. When the reference and test branches (**B3**) are combined, we can achieve the best performance due to the iterative cross-modal distillation operation that can optimize the learned representation of TIR modality.

Target Center Location (TCL). We analyze the cross-modal distillation on TCL as ‘C1. Learning on reference branch’ in Tab. 1. As the model parameters of TCL mostly gather in the reference branch, we only update the parameters in the reference branch. Therefore, the input of the reference branch is the data from the TIR modality and that of the test branch is the data from the RGB modality. From Tab. 1, we can see that the tracker equipped with our cross-modal distillation mechanism performs better than the baseline tracker. With the setting **C1**, we improve the results of the tracker by 1.5% in Success, 2.0% in Precision, and 1.8% in Normalized Precision, showing that we successfully transfer the RGB representation in the test branch to the reference branch in TCL.

By using cross-modal distillation on TCL, we obtain the tracking results which exceed the performances of the tracker equipped with the **B1** and **B2**. We attribute it to the reason that TCL is critical for the location of the target, while BBE is to further refine the target scale based on the location of TCL. All these analyses furthermore demonstrate that our cross-modal distillation can learn more effective TIR representations for TIR images.

Joint Learning of BBE and TCL. Inspired by the combination of the learned reference branch and test branch in BBE as **B3**, we attempt to jointly train BBE and TCL. We consider five cases, i.e. ‘D1. Joint learning on test branch’, ‘D2. Joint learning on reference branch’, ‘D3. Joint learning’, ‘D4. Joint learning with random sampling’ and ‘D5. Joint learning with a detector’, aiming to obtain optimized parameters for each of branches on both TCL and BBE.

We use parameters of the combination models on **B3** and **C1** for the initialization of our model for jointly training. At the beginning, we try to train one branch with TIR input by transferring from the other branch with RGB input for both TCL and BBE simultaneously, i.e. case (**D1**) and case (**D2**). We improve the performances of the baseline tracker in both two cases.

In order to better train the model jointly and avoid the simple combination of two branches in TCL and BBE, we mix patches from both RGB and TIR modalities and then input them to the reference branch and also the test branch as in **D3**. The results of this method achieve 68.0% in Success, 80.8% in Precision, and 72.7%

in Normalized Precision, which are the best results at present. All above cases, i.e. **B1**, **B2**, **B3**, **C1**, **D1**, **D2** and **D3**, use the central regions in the RGB and TIR images for generating the patches.

Besides, we propose randomly assigning different regions in an image to be the input patches as **D4**. Then we can obtain the paired patches under the alignment restriction of the paired images. For this case, we improve the performance by 1.1% on Success, 1.5% on Precision, and 1.4% on Normalized Precision.

Additionally, we explore to use a detector to generate the input patches for distillation training as **D5**. But the classes of generated patches are restricted to the pre-trained detector model, which normally detects objects with specific classes such as ‘person’, ‘car’ etc. Therefore, the results obtained by setting with a detector, namely case (**D5**), are not good enough, compared with results of case (**D3**). As CMD with a detector is a bit better than that with randomly sampling as in Tab. 1, we attribute it to that ‘random sampling’ could not contain any complete object contents of the image.

For now, we can summarize that input patches generated by a detector and random sampling are both restricted to the varieties of the object classes of training data during cross-modal distillation.

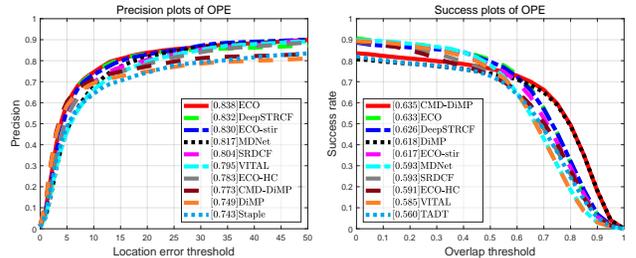


Figure 4: Precision plot and success plot by comparing our tracker with the top-10 trackers on PTB-TIR dataset [30].

4.4 LSOTB-TIR Dataset

In this part, we compare our tracker (CMD-DiMP) against high-quality trackers on the LSOTB-TIR [33] dataset in terms of Success, Precision and Normalized Precision in Fig. 3, including TIR tracker, e.g. ECO-stir [45], and RGB trackers: e.g. MDNet [35], VITAL [39], SiamRPN++ [25], SiamMask [42], SiamFC-TIR [13], ECO-HC [8] and so on. By using our cross-modal training, we improve baseline tracker [4] with absolute gains of 1.8%, 2.1% and 2.0% in terms of Success, Precision and Normalized Precision respectively. These obvious improvements prove that our cross-modal training

Table 2: Attributes analysis of cross-modal distillation on LSOTB-TIR dataset [33]. We evaluate our trackers with several RGB and TIR trackers on Precision, Normalized Precision, and Success (P/NP/S %). The best results are highlighted in bold font.

Attributes Type	Attributes Name	SiamFC-TIR [13]	SiamRPN++ [25]	VITAL [39]	ECO-stir [45]	DiMP [4]	CMD-DiMP
Challenge	Deformation	71.1/61.9/56.3	66.8/60.2/56.6	74.8/64.7/56.9	76.3/65.9/60.1	76.3/67.0/63.3	80.5/70.7/66.8
	Occlusion	65.1/59.2/53.9	63.9/58.6/54.6	72.7/65.9/59.4	71.1/64.9/59.4	73.7/65.7/62.2	76.2/67.8/64.0
	Distractor	63.6/56.8/51.0	64.8/58.1/54.6	71.4/64.9/57.7	74.7/67.2/61.6	75.0/66.4/63.0	76.1/67.3/63.7
	Background clutter	69.2/62.3/54.6	70.4/65.1/60.6	73.0/67.2/58.4	73.9/67.2/61.3	78.9/70.9/66.4	79.8/71.6/67.0
	Out of view	73.0/65.9/58.6	75.0/70.7/63.9	70.6/68.5/59.7	76.7/72.1/66.8	77.4/74.0/67.4	79.7/76.1/69.6
	Scale variation	74.7/68.0/58.4	81.3/74.2/68.2	79.7/74.5/62.1	80.5/74.2/68.0	89.2/81.3/75.1	91.0/83.0/76.7
	Fast motion	72.4/68.8/59.7	74.8/70.3/64.4	72.2/67.8/58.2	73.8/68.0/63.6	82.7/77.1/71.1	86.1/80.1/74.1
	Motion blur	74.0/64.7/57.6	70.7/64.8/58.8	76.8/69.4/59.3	76.7/66.8/61.3	80.3/73.0/67.1	85.6/77.7/71.5
	Thermal crossover	68.0/55.7/51.7	60.0/50.2/48.4	78.0/64.3/58.1	73.1/58.6/54.5	61.9/50.6/50.2	64.7/54.2/52.9
	Intensity variation	85.0/77.7/71.6	83.6/76.7/74.4	74.1/72.5/61.8	76.9/76.1/70.9	91.4/87.6/82.7	91.5/87.7/82.8
	Low resolution	91.4/74.2/65.2	81.2/66.3/62.0	91.1/73.6/60.9	94.1/69.7/64.0	78.1/64.6/60.3	83.1/68.4/64.0
Aspect ratio variation	69.5/59.8/48.9	70.6/65.0/59.4	72.9/63.7/54.2	72.4/58.1/54.9	78.9/72.7/67.2	80.4/70.3/65.0	
Scenario	Vehicle-mounted	74.5/70.1/59.4	86.0/79.2/72.6	83.7/81.5/72.1	84.4/81.1/76.0	91.9/84.2/78.8	96.6/88.3/82.4
	Drone-mounted	68.2/60.2/53.7	64.9/58.4/55.0	69.5/64.1/53.8	69.5/61.5/55.7	75.5/69.7/64.1	73.0/67.4/62.0
	Surveillance	63.9/58.0/53.5	67.0/61.7/57.7	70.0/63.5/57.5	69.4/64.1/59.1	74.5/65.4/62.9	74.8/65.3/62.8
	Hand-held	74.8/64.6/56.3	70.7/64.7/59.8	78.8/68.3/58.8	79.4/66.4/60.3	78.0/69.9/64.4	84.1/75.8/69.9
All	All	69.8/62.4/55.3	70.9/64.8/60.3	74.8/68.1/59.6	74.9/67.1/61.6	78.7/70.7/66.2	80.8/72.7/68.0

can extract more useful and consistent information between RGB and TIR modalities.

Besides, we extend our cross-modal distillation training mechanism to the variant tracker ATOM [9] to prove our method’s generality. For the original ATOM [9] tracker, only the IoUNet component is updated during offline training. Therefore, here we only employ the cross-modal distillation on the Bounding Box Estimation (BBE). We evaluate our tracker (CMD-ATOM) on the LSOTB-TIR dataset [33]. And we improve the baseline tracker ATOM [9] with absolute gains of 2.3%, 2.7%, and 2.5% in Success, Precision, and Normalized Precision respectively, shown in Fig. 3.

4.5 PTB-TIR Dataset

We evaluate our tracker (CMD-DiMP) on the PTB-TIR dataset [30] using the two evaluation metrics, i.e. Precision and Success. We compare our tracker with some trackers on PTB-TIR dataset [30], including TIR tracker ECO-stir [45], and RGB trackers, such as ECO [8], DeepSTRCF [27], MDNet [35], SRDCF [10], VITAL [39], Staple [2], TADT [28], MCCT [41], ECO-HC [8] and so on.

As shown in Fig. 4, our tracker gets 0.773 in terms of Precision and 0.635 in terms of Success, which improves 2.4% and 1.7% for the baseline tracker separately. We win first place in the Success plot and get a lower score in the Precision plot. We observe that the predicted bounding boxes will be expanded to a larger area for containing more similar areas, such as the distractors or background clutter, which reduces the Precision drastically. More importantly, as the BBE pays more attention to the object’s scale changes, there is an obvious improvement when the overlap threshold is bigger than 0.6 in the Success plot of OPE.

4.6 Attribute Analysis on LSOTB-TIR Dataset

The LSOTB-TIR [33] dataset provides 16 attributes to be evaluated, including 12 challenges and 4 scenarios. We compare our tracker (CMD-DiMP) with TIR trackers, ECO-stir [45] and SiamFC-TIR [13], and RGB trackers, VITAL [39], SiamRPN++ [25] and DiMP [4] in terms of Precision, Normalized Precision and Success (P/NP/S %). Table 2 shows that our tracker outperforms the above five trackers for most of the challenges except thermal crossover, low resolution and drone-mounted.

We improve the performance by about 1%-6% in most of attributes compared with DiMP [4]. Compared with other four trackers, we obtain significant progress in most challenges. For attributes like deformation, background clutter, and aspect ratio variation, we get over 6% absolute gains in terms of Success score. Even more, we make progress of 8% in terms of Success score for attributes of scale variation and intensity variation. And for challenges such as fast motion and motion blur, our tracker achieves 10% improvement on Success score. Besides, for vehicle-mounted and drone-mounted, the improvement on Success score made by our tracker is up to 6% and that for hand-held is near to 10%.

5 CONCLUSIONS

In this paper, we propose to learn TIR-specific representations by distilling the pre-trained model of RGB modality to TIR modality on unlabeled RGB-TIR datasets, called cross-modal distillation (CMD). During distillation, representations from the RGB modality model are used as supervised signals to guide the learning of TIR-specific representations for the TIR modality. We take advantage of two branches of the network to deal with data from different modalities for transferring cross-modal knowledge. Moreover, we explore the prior semantic consistency implied in the training data structure itself to automatically generate paired RGB-TIR patches for training. Unlike the normal RGB tracking datasets, the training data for our method can be without any annotation about the targets. Our cross-modal distillation mechanism reduces the dependency of training networks on labeled data and increases the extensibility with a larger amount of unlabeled data in the future. Through this mechanism, we leverage the pre-trained RGB models which are derived from large-scale annotated RGB data, to train models specific for representing targets in TIR modality in an unsupervised way. Extensive experimental results on two standard TIR datasets demonstrate that our proposed CMD effectively improves the performance of the baseline tracker for TIR tracking.

ACKNOWLEDGMENTS

This work is supported in part by National Natural Science Foundation of China (61773397, 62006245, 61971352, 61862043, U19B2037).

REFERENCES

- [1] C.S. Asha and A.V. Narasimhadhan. 2017. Robust infrared target tracking using discriminative and generative approaches. *Infrared Physics Technology* 85 (2017), 114–127.
- [2] Luca Bertinetto, Jack Valmadre, Stuart Golodetz, Ondrej Miksik, and Philip H. S. Torr. 2016. Staple: Complementary Learners for Real-Time Tracking. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 1401–1409.
- [3] Luca Bertinetto, Jack Valmadre, João F. Henriques, Andrea Vedaldi, and Philip H. S. Torr. 2016. Fully-Convolutional Siamese Networks for Object Tracking. In *European Conference Computer Vision*. 850–865.
- [4] Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte. 2019. Learning Discriminative Model Prediction for Tracking. In *The IEEE Conference on International Conference on Computer Vision*. 6181–6190.
- [5] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. 2020. YOLOv4: Optimal Speed and Accuracy of Object Detection. *CoRR* abs/2004.10934 (2020).
- [6] Yun-Chun Chen, Yen-Yu Lin, Ming-Hsuan Yang, and Jia-Bin Huang. 2019. CrDoCo: Pixel-Level Domain Transfer With Cross-Domain Consistency. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 1791–1800.
- [7] Navneet Dalal and Bill Triggs. 2005. Histograms of Oriented Gradients for Human Detection. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 886–893.
- [8] Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. 2017. ECO: Efficient Convolution Operators for Tracking. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 6931–6939.
- [9] Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. 2019. ATOM: Accurate Tracking by Overlap Maximization. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 4660–4669.
- [10] Martin Danelljan, Gustav Häger, Fahad Shahbaz Khan, and Michael Felsberg. 2015. Learning Spatially Regularized Correlation Filters for Visual Tracking. In *The IEEE Conference on International Conference on Computer Vision*. 4310–4318.
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. 2009. ImageNet: A large-scale hierarchical image database. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 248–255.
- [12] Tuong Do, Huy Tran, Thanh-Toan Do, Erman Tjiputra, and Quang D. Tran. 2019. Compact Trilinear Interaction for Visual Question Answering. In *The IEEE Conference on International Conference on Computer Vision*. 392–401.
- [13] Xingping Dong and Jianbing Shen. 2018. Triplet Loss in Siamese Network for Object Tracking. In *European Conference Computer Vision*. 472–488.
- [14] Heng Fan, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Hexin Bai, Yong Xu, Chunyuan Liao, and Haibin Ling. 2019. LaSOT: A High-Quality Benchmark for Large-Scale Single Object Tracking. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 5374–5383.
- [15] Michael Felsberg, Amanda Berg, Gustav Häger, Jörgen Ahlberg, Matej Kristan, Jiri Matas, Ales Leonardis, Luka Cehovin, Gustavo Fernández, Tomás Vojir, Georg Nebel, and Roman P.flugfelder. 2015. The Thermal Infrared Visual Object Tracking VOT-TIR2015 Challenge Results. In *The IEEE Conference on International Conference on Computer Vision*. 639–651.
- [16] Peng Gao, Yipeng Ma, Ke Song, Chao Li, Fei Wang, and Liyi Xiao. 2018. Large Margin Structured Convolution Operator for Thermal Infrared Object Tracking. In *International Conference on Pattern Recognition*. 2380–2385.
- [17] Erhan Gundogdu, Aykut Koc, Berkan Solmaz, Riad I. Hammoud, and A. Aydin Alatan. 2016. Evaluation of Feature Channels for Correlation-Filter-Based Visual Object Tracking in Infrared Spectrum. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 290–298.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 770–778.
- [19] Yujie He, Min Li, Jinli Zhang, and Junping Yao. 2016. Infrared Target Tracking Based on Robust Low-Rank Sparse Learning. *IEEE Geoscience and Remote Sensing Letters* 13, 2 (2016), 232–236.
- [20] João F. Henriques, Rui Caseiro, Pedro Martins, and Jorge Batista. 2015. High-Speed Tracking with Kernelized Correlation Filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37, 3 (2015), 583–596.
- [21] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2014. Distilling the Knowledge in a Neural Network. (2014).
- [22] Lianghua Huang, Xin Zhao, and Kaiqi Huang. 2018. GOT-10k: A Large High-Diversity Benchmark for Generic Object Tracking in the Wild. *CoRR* abs/1810.11981 (2018).
- [23] Soonmin Hwang, Jaesik Park, Namil Kim, Yukyung Choi, and In So Kweon. 2015. Multispectral pedestrian detection: Benchmark dataset and baseline. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 1037–1045.
- [24] Jogendra Nath Kundu, Nishank Lakkakula, and Venkatesh Babu Radhakrishnan. 2019. UM-Adapt: Unsupervised Multi-Task Adaptation Using Adversarial Cross-Task Distillation. In *The IEEE Conference on International Conference on Computer Vision*. 1436–1445.
- [25] Bo Li, Wei Wu, Qiang Wang, Fangyi Zhang, Junliang Xing, and Junjie Yan. 2019. SiamRPN++: Evolution of Siamese Visual Tracking With Very Deep Networks. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 4282–4291.
- [26] Chenglong Li, Nan Zhao, Yijuan Lu, Chengli Zhu, and Jin Tang. 2017. Weighted Sparse Representation Regularized Graph Learning for RGB-T Object Tracking. In *ACM International Conference on Multimedia*. 1856–1864.
- [27] Feng Li, Cheng Tian, Wangmeng Zuo, Lei Zhang, and Ming-Hsuan Yang. 2018. Learning Spatial-Temporal Regularized Correlation Filters for Visual Tracking. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 4904–4913.
- [28] Xin Li, Chao Ma, Baoyuan Wu, Zhenyu He, and Ming-Hsuan Yang. 2019. Target-Aware Deep Tracking. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 1369–1378.
- [29] Ying Li, Pengcheng Li, and Qiang Shen. 2014. Real-time infrared target tracking based on l1 minimization and compressive features. *Applied Optics* 53, 28 (2014), 6518–6526.
- [30] Qiao Liu, Zhenyu He, Xin Li, and Yuan Zheng. 2020. PTB-TIR: A Thermal Infrared Pedestrian Tracking Benchmark. *IEEE Transactions on Multimedia* 22, 3 (2020), 666–675.
- [31] Qiao Liu, Xin Li, Zhenyu He, Nana Fan, Di Yuan, Wei Liu, and Yongsheng Liang. 2020. Multi-Task Driven Feature Models for Thermal Infrared Tracking. In *AAAI Conference on Artificial Intelligence*. 11604–11611.
- [32] Qiao Liu, Xin Li, Zhenyu He, Nana Fan, Di Yuan, and Hongpeng Wang. 2020. Learning Deep Multi-Level Similarity for Thermal Infrared Object Tracking. *IEEE Transactions on Multimedia* PP (2020), 1–1.
- [33] Qiao Liu, Xin Li, Zhenyu He, Chenglong Li, Jun Li, Zikun Zhou, Di Yuan, Jing Li, Kai Yang, Nana Fan, and Feng Zheng. 2020. LSOTB-TIR: A Large-Scale High-Diversity Thermal Infrared Object Tracking Benchmark. In *ACM International Conference on Multimedia*. 3847–3856.
- [34] Qiao Liu, Xiaohuan Lu, Zhenyu He, Chunkai Zhang, and Wen-Sheng Chen. 2017. Deep convolutional neural networks for thermal infrared object tracking. *Knowledge Based Systems* 134 (2017), 189–198.
- [35] H. Nam and B. Han. 2016. Learning Multi-domain Convolutional Neural Networks for Visual Tracking. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 4293–4302.
- [36] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. 2017. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39, 6 (2017), 1137–1149.
- [37] Siddharth Roheda, Benjamin S. Riggan, Hamid Krim, and Liyi Dai. 2018. Cross-Modality Distillation: A Case for Conditional Generative Adversarial Networks. In *IEEE International Conference on Acoustics, Speech and Signal Processing*. 2926–2930.
- [38] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *International Conference on Learning Representations*.
- [39] Yibing Song, Chao Ma, Xiaohu Wu, Lijun Gong, Linchao Bao, Wangmeng Zuo, Chunhua Shen, Rynson W. H. Lau, and Ming-Hsuan Yang. 2018. VITAL: Visual Tracking via Adversarial Learning. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 8990–8999.
- [40] Yonglong Tian, Dilip Krishnan, and Phillip Isola. 2020. Contrastive Representation Distillation. In *International Conference on Learning Representations*.
- [41] Ning Wang, Wengang Zhou, Qi Tian, Richang Hong, Meng Wang, and Houqiang Li. 2018. Multi-Cue Correlation Filters for Robust Visual Tracking. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 4844–4853.
- [42] Qiang Wang, Li Zhang, Luca Bertinetto, Weiming Hu, and Philip H. S. Torr. 2019. Fast Online Object Tracking and Segmentation: A Unifying Approach. In *The IEEE Conference on Computer Vision and Pattern Recognition*. 1328–1338.
- [43] Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. 2015. Object Tracking Benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37, 9 (2015), 1834–1848.
- [44] Xianguo Yu, Qifeng Yu, Yang Shang, and Hongliang Zhang. 2017. Dense structural learning for infrared object tracking at 200+ Frames per Second. *Pattern Recognition Letters* 100 (2017), 152–159.
- [45] Lichao Zhang, Abel Gonzalez-Garcia, Joost van de Weijer, Martin Danelljan, and Fahad Shahbaz Khan. 2019. Synthetic Data Generation for End-to-End Thermal Infrared Tracking. *IEEE Transactions on Image Processing* 28, 4 (2019), 1837–1850.
- [46] Pengyu Zhang, Dong Wang, and Huchuan Lu. 2020. Multi-modal Visual Tracking: Review and Experimental Comparison. *CoRR* abs/2012.04176 (2020).
- [47] Xingchen Zhang, Ping Ye, Henry Leung, Ke Gong, and Gang Xiao. 2020. Object fusion tracking based on visible and infrared images: A comprehensive review. *Information Fusion* 63 (2020), 166–187.