

# Dynamic Fusion Network for RGBT Tracking

Jingchao Peng, Haitao Zhao, Zhengwei Hu

Automation Department, School of Information Science and Engineering,  
East China University of Science and Technology, Shanghai, China

## Abstract

For both visible and infrared images have their own advantages and disadvantages, RGBT tracking has attracted more and more attention. The key points of RGBT tracking lie in feature extraction and feature fusion of visible and infrared images. Current RGBT tracking methods mostly pay attention to both *individual features* (features extracted from images of a single camera) and *common features* (features extracted and fused from an RGB camera and a thermal camera), while pay less attention to the different and dynamic contributions of individual features and common features for different sequences of registered image pairs. This paper proposes a novel RGBT tracking method, called Dynamic Fusion Network (DFNet), which adopts a two-stream structure, in which two non-shared convolution kernels are employed in each layer to extract individual features. Besides, DFNet has shared convolution kernels for each layer to extract common features. Non-shared convolution kernels and shared convolution kernels are adaptively weighted and summed according to different image pairs, so that DFNet can deal with different contributions for different sequences. DFNet has a fast speed, which is 28.658 FPS. The experimental results show that when DFNet only increases the Mult-Adds of 0.02% than the non-shared-convolution-kernel-based fusion method, Precision Rate (PR) and Success Rate (SR) reach 88.1% and 71.9% respectively.

## 1 Introduction

Object tracking is a popular computer vision task, whose purpose is to continuously track the position of the object in the subsequent frames when given in the first frame. Tracking in a complex visual scenery, including rain, smoke, or night, is one of the most difficult computer vision tasks [1, 2], especially for visible-image-based tracking [2, 3]. However, infrared sensors can work around the clock, infrared has a strong ability to penetrate smoke, which can supplement the deficiencies of visible images in bad visual conditions [4, 5, 6, 7]. Therefore, RGBT tracking has attracted more and more attention.

Since 2018, due to the powerful learning ability, Deep Learning (DL) models, especially Convolutional Neural Networks (CNN), are widely used to address RGBT tracking [8, 9, 10, 11, 12, 13, 14, 15]. DL-based tracking methods have demonstrated their capabilities over traditional fusion tracking methods [16, 17, 18, 19] or other tracking methods (e.g., sparse representation-based methods [20, 21, 22], and graph-based methods [23, 24, 25, 26]). The advantage of DL-based tracking methods is their ability to learn more effective and robust features than hand-crafted features [27, 28, 3]. DL-based tracking methods can be divided into pixel-level [8, 9, 10], feature-level [11, 12, 13, 14], and decision-level [15] fusion tracking. Compared with the pixel-level fusion method, the feature-level fusion method has lower requirements for image registration and can tolerate a certain amount of noise [12, 13]. Compared with the decision-level fusion method, it has lower computational complexity and faster speed [29, 3]. Recently research works of DL-based RGBT tracking mainly focus on feature-level fusion [3].

Due to visible light reflection and infrared radiation have different imaging properties, visible and infrared images have different *individual features* [30], which can be used to track objects based on single-modal images. In visible-image-based tracking, objects can be distinguished through rich textures and different colors. While in infrared-image-based tracking, objects can be distinguished by high-contrast light-dark changes that reflect the heat of the object. In order to fully utilize the individual features from the two different modalities, feature-level fusion methods are often adopted, in which two Convolutional Neural Networks (CNNs) were often employed to handle visible and infrared images, respectively. For example, Zhang et al. [11] utilized two different CNNs to respectively extract individual features from visible and infrared images. In their work the visible and infrared features were concatenated and sent to follow layers for tracking the object. ConvNet [13] and SiamFT [12] employ fusion sub-networks to select discriminative features after extracting the individual features. DSiamMFT [31] and FANet [32] focus on multi-layer feature fusion to achieve more effective hierarchical feature aggregation. For simplicity, this paper denotes the basic feature-level fusion method without any bells and whistles, which only uses two different CNNs to respectively handle visible and infrared images, as non-shared-convolution-kernel-based fusion method.

In addition to individual features, since visible and infrared images are shot in the same scene and are used to track the same object, there are *common features* in the two modalities [33]. Common features reflect the size, location, contour, and so on, which are also important information in object tracking [34]. When individual features are not enough to achieve good tracking performances, it is necessary to use common features such as the semantics of the object, and other characteristics of the object at the corresponding position of the visible and infrared images for object tracking. Both MANet [14], CAT [30], IVFuseNet [34], and SiamIVFN [29] use a shared convolution kernel to extract common features. Their experimental results show that the shared-convolution-kernel-based fusion methods can extract common features that are more informative than the non-shared-convolution-kernel-based fusion method. But in their networks, the contributions of individual features and common fea-

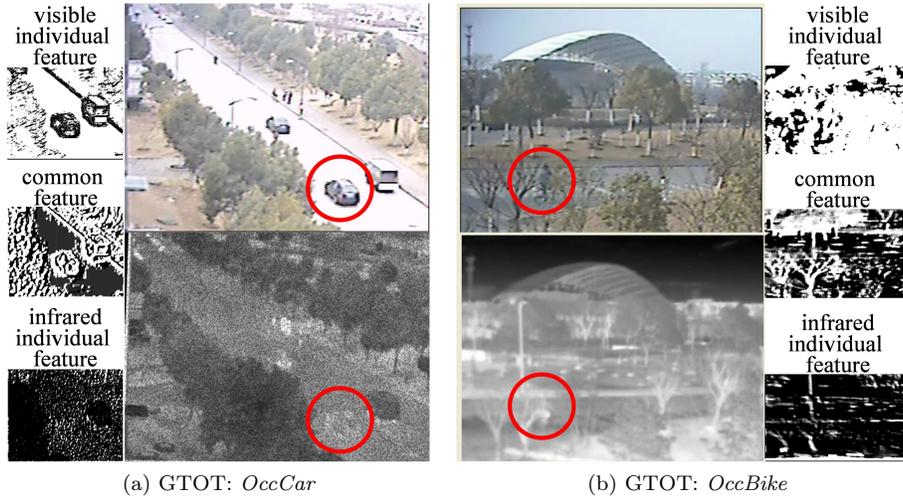


Figure 1: Two registered image pairs with different contributions of individual features and common features. The two images in the first row are visible images and the other two images in the second row are the corresponding registered infrared images. The images in (a) are from GTOT: *OccCar*. The visible image is clear but the infrared image is not clear because of noise, so the tracking task should pay more attention to the individual features of the visible image. The images in (b) are from GTOT: *OccBike*. the target is not easy to detect in both visible and infrared images due to the messy background. The tracking task should pay more attention to the common features.

tures are prefixed and have no consideration of adaption to the registered image pairs captured in different scenes.

However, the contributions of individual features and common features are not always fixed. Visible images are greatly affected by the illumination and prone to overexposure or underexposure. Infrared images are easy to be interfered with by the external scenery and internal systems and prone to noise [35, 36]. In other words, the reliability of different modalities is not always fixed. For different reliable modalities, individual features and common features contribute to different degrees. When one modality is reliable, the individual features of the reliable modality contribute more, as shown in Figure 1 (a). When it is impossible to track based on single-modal images, more attention needs to be paid to common features, as shown in Figure 1 (b). Therefore, the tracker needs to adaptively calculate different contributions of individual features and common features in different scenes.

To solve the performance limitation of the network in changing scenes, dynamic convolution has natural advantages. The concept of dynamic convolution (e.g., CondConv [37], Dynamic Convolution [38], and WeightNet [38]) usually

adopts the method of attention over convolution kernels. Dynamic convolution has been applied in scene segmentation, scene synthesizing, image inpainting, biomedical imaging, and so on [39, 40, 41, 42]. Due to aggregating multiple convolution kernels adapted to each input, dynamic convolution has more representation power without increasing the width and depth of the network. The aggregation of multiple convolution kernels in convolution kernel space makes it possible to make full use of multiple convolution kernels only by one convolution operation. Therefore, dynamic convolution is computationally efficient. But dynamic convolution is designed for integration into existing CNN architectures, cannot aggregate individual features and common features in fusion tasks.

Motivated by the above analysis, we propose a novel RGBT tracking method called dynamic fusion network (DFNet). DFNet adopts a two-stream structure, which has non-shared convolution kernels to extract individual features. One CNN is utilized to extract features from infrared images, and the other one is for handle visible images. Besides, DFNet has shared convolution kernels to extract common features. DFNet adaptively merges the shared convolution kernels and the non-shared convolution kernels in convolution kernel space through the dynamic convolution. To satisfy strict latency requirements for object tracking, DFNet only needs two convolution operations in each layer to extract the individual and common features of visible and infrared images. Since the weights of shared and non-shared convolution kernels are dynamically computed, it can adaptively calculate the contributions of individual features and common features to different scenes.

Specifically, the proposed method has the following advantages:

1. DFNet has shared kernels and non-shared kernels that separately extract the common features and individual features. DFNet has a strong representation power.
2. DFNet adaptively calculates the contributions of individual features and common features according to different registered image pairs.
3. By fusing multiple kernels in convolution kernel space, DFNet boosts the PR/SR by 1.1%/0.9% with only 0.02% additional Mult-Adds.

## 2 Related Work

Section 1 overviews DL-based RGBT tracking. This section focuses on three most related works to ours: ConvNet [13], MANet [14], IVFuseNet [34]. Their simplified feature extraction layer diagrams are shown in Figure 2.

### 2.1 ConvNet

ConvNet [13] uses different convolutional networks to extract the individual features of visible and infrared images and then fuses them, its feature extraction layer can be expressed as:

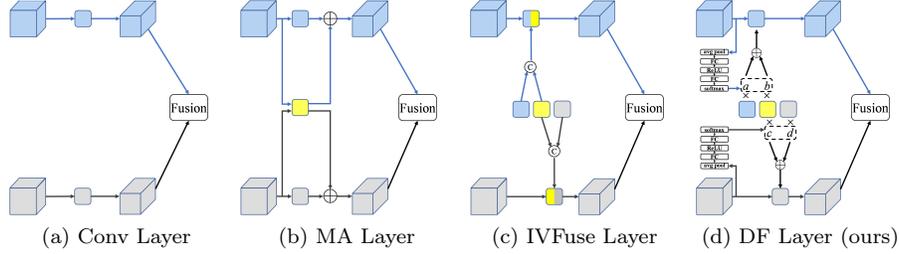


Figure 2: Simplified feature extraction layer diagrams of (a) ConvNet. The blue and gray branches respectively handle visible and infrared images. (b) MANet. The yellow block represents the shared convolution kernel which is used to extract common features. (c) IVFuseNet. The shared and the non-shared convolution kernels are concatenated to form a new convolution kernel for feature extraction. (d) DFNet (ours) The shared and the non-shared convolution kernels are weighted and summed to form a new convolution kernel for feature extraction. The weights are calculated based on the input.

Table 1: Expressions of different fusion methods.  $W_{RGB}$  and  $W_T$  represent the non-shared convolution kernel for RGB features and the convolution kernel for thermal features respectively.  $W_{share}$  represents the shared convolution kernel.  $F_{RGB}$  and  $F_T$  represent the input of visible and infrared branch respectively.  $*$  represents convolution operation.  $\sigma(\cdot)$  represents activation function.

Fusion method	Expression
ConvNet	$\sigma \left( \begin{bmatrix} W_{RGB} & W_T \end{bmatrix} * \begin{bmatrix} F_{RGB} \\ F_T \end{bmatrix} \right)$ $= \begin{bmatrix} \sigma(W_{RGB} * F_{RGB}) & \sigma(W_T * F_T) \end{bmatrix}$
MANet	$\begin{bmatrix} \sigma W_{RGB} & \sigma W_{share} & \sigma W_T \end{bmatrix} * \left( \begin{bmatrix} 1 \\ 1 & 1 \\ & & 1 \end{bmatrix} * \begin{bmatrix} F_{RGB} \\ F_T \end{bmatrix} \right)$ $= \begin{bmatrix} \sigma(W_{RGB} * F_{RGB}) + \sigma(W_{share} * F_{RGB}) & \sigma(W_{share} * F_T) + \sigma(W_T * F_T) \end{bmatrix}$
IVFuseNet	$\sigma \left( \left( \begin{bmatrix} W_{RGB} & W_{share} & W_T \end{bmatrix} \begin{bmatrix} 1 \\ 1 & 1 \\ & & 1 \end{bmatrix} \right) * \begin{bmatrix} F_{RGB} \\ F_T \end{bmatrix} \right)$ $= \begin{bmatrix} \sigma \left( \begin{bmatrix} W_{RGB} & W_{share} \end{bmatrix} * F_{RGB} \right) & \sigma \left( \begin{bmatrix} W_{share} & W_T \end{bmatrix} * F_T \right) \end{bmatrix}$
DFNet	$\sigma \left( \left( \begin{bmatrix} W_{RGB} & W_{share} & W_T \end{bmatrix} \begin{bmatrix} a & c \\ b & d \end{bmatrix} \right) * \begin{bmatrix} F_{RGB} \\ F_T \end{bmatrix} \right)$ $= \begin{bmatrix} \sigma \left( (aW_{RGB} + bW_{share}) * F_{RGB} \right) & \sigma \left( (cW_{share} + dW_T) * F_T \right) \end{bmatrix}$

$$\begin{aligned} & \sigma \left( \begin{bmatrix} W_{RGB} & W_T \end{bmatrix} * \begin{bmatrix} F_{RGB} \\ F_T \end{bmatrix} \right) \\ &= \begin{bmatrix} \sigma(W_{RGB} * F_{RGB}) & \sigma(W_T * F_T) \end{bmatrix} \end{aligned} \quad (1)$$

where  $W_{RGB}$  and  $W_T$  represent the convolution kernel for RGB features and the convolution kernel for thermal features respectively,  $F_{RGB}$  and  $F_T$  represent the input of visible and infrared branch respectively,  $*$  represents convolution operation,  $\sigma(\cdot)$  represents activation function, such as ReLU.

In ConvNet, different convolution kernels are used to extract individual features from visible and infrared images. Then these two features are fused and sent to domain-specific layers for binary classification and identification of the target. Besides, ConvNet designs a fusion sub-network, which adaptively fuses two individual features to removing redundant noise. The feature extraction process performed two convolution operations in one layer, therefore the speed of ConvNet is fast. ConvNet focuses on individual features but does not fully consider the common features.

## 2.2 MANet

Li, et al. [14] argue that common features of visible and infrared images is crucial to the effectiveness of the fusion. Therefore, MANet employs a shared convolution kernel to extract the common features of visible and infrared images. The feature extraction layer of MANet can be expressed as:

$$\begin{aligned} & \begin{bmatrix} \sigma W_{RGB} & \sigma W_{share} & \sigma W_T \end{bmatrix} * \left( \begin{bmatrix} 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} F_{RGB} \\ F_T \end{bmatrix} \right) \\ &= \begin{bmatrix} \sigma(W_{RGB} * F_{RGB}) + \sigma(W_{share} * F_{RGB}) & \sigma(W_{share} * F_T) + \sigma(W_T * F_T) \end{bmatrix} \end{aligned} \quad (2)$$

where  $W_{share}$  represents the shared convolution kernel.

Before respective convolution operations, visible and infrared features must both undergo a shared convolution operation, which uses the same convolution kernel. It is worth noting that MANet fuses shared and non-shared features in the feature space. Please note that the activation function  $\sigma(\cdot)$  is not a linear operation, and we have:

$$\sigma(W_{RGB} * F_{RGB}) + \sigma(W_{share} * F_{RGB}) \neq \sigma((W_{RGB} + W_{share}) * F_{RGB}) \quad (3)$$

$$\sigma(W_{share} * F_T) + \sigma(W_T * F_T) \neq \sigma((W_{share} + W_T) * F_T) \quad (4)$$

Therefore, four convolution operations are needed. The complexity of the operation is large, which is not conducive to the real-time requirements of the tracking task. In addition, the weights of the shared and non-shared convolution kernel are equal, so that they cannot be adjusted in real-time in the face of different contributions of individual features and common features.

### 2.3 IVFuseNet

Unlike the fusion of shared and non-shared features in feature space, IVFuseNet [34] merges the shared and non-shared convolution kernels in convolution kernel space. IVFuseNet concatenates two small-sized convolution kernels, one of them is a shared convolution kernel. The visible and infrared images are respectively convolved with different spliced convolution kernels. The feature extraction layer of IVFuseNet can be expressed as:

$$\sigma \left( \left( \begin{bmatrix} W_{RGB} & W_{share} & W_T \end{bmatrix} \begin{bmatrix} 1 & & & \\ & 1 & 1 & \\ & & & 1 \end{bmatrix} * \begin{bmatrix} F_{RGB} & \\ & F_T \end{bmatrix} \right) \right) \quad (5)$$

$$= \begin{bmatrix} \sigma \left( \begin{bmatrix} W_{RGB} & W_{share} \end{bmatrix} * F_{RGB} \right) & \sigma \left( \begin{bmatrix} W_{share} & W_T \end{bmatrix} * F_T \right) \end{bmatrix}$$

Since the shared and non-shared convolution kernels are fused in convolution kernel space, IVFuseNet only needs to perform two convolution operations. However, due to the shared and non-shared convolution kernels are concatenated in advance, the size of the convolution kernel is smaller than that of MANet, which means IVFuseNet has weak representation power than MANet. For example, in MANet, the size of the shared and non-shared convolution kernel in the first layer is  $96 \times 3 \times 7 \times 7$  and  $96 \times 3 \times 3 \times 3$ ; while in IVFuseNet, the size of the shared and non-shared convolution kernel in the corresponding layer is  $24 \times 3 \times 7 \times 7$  and  $72 \times 3 \times 3 \times 3$ . Besides, the channel size of the shared and non-shared convolution kernel needs to be prefixed, and the coupling rate cannot be adjusted in real-time in the face of different contributions of individual features and common features.

We summarize the related works below:

1. the speed of ConvNet is fast, but ConvNet does not have shared convolution kernel to extract common features.
2. Although MANet has both shared and non-shared convolution kernels, the speed of MANet is much slower than that of ConvNet. Moreover, MANet has no design to deal with the different contributions of the individual features and common features.
3. IVFuseNet have both shared and non-shared convolution kernels, and the speed of IVFuseNet is fast. However, compared with MANet, IVFuseNet has weak representation power. Moreover, IVFuseNet also has no procedure to handle the different contributions of the individual features and common features.

## 3 Our Method

In this section, we will introduce a novel RGBT tracking method called dynamic fusion network (DFNet). We first introduce the dynamic fusion layer. Dynamic

convolution is used in the convolution kernel space to fuse shared and non-shared convolution kernels. Then we use the dynamic fusion layer as the basic module to construct DFNet for RGBT tracking. DFNet has the advantages of MANet and IVFuseNet, which has shared convolution kernels to extract common features. We highlight the differences of the network structures between DFNet and the related models in Table 1. Due to the fusion of shared and non-shared convolution kernels in convolution kernel space, DFNet has high inference efficiency. Besides, adaptive convolutional features can be extracted in the face of changes in the scene because of its dynamic nature.

### 3.1 Dynamic Fusion Layer

Dynamic fusion layer fuses the shared convolution kernel and non-shared convolution kernels in convolution kernel space. That is, the convolution kernels are merged, then the convolution operation is performed:

$$\sigma \left( \left( \left[ \begin{array}{ccc} W_{RGB} & W_{share} & W_T \end{array} \right] \left[ \begin{array}{cc} a & \\ b & c \\ & d \end{array} \right] \right) * \left[ \begin{array}{cc} F_{RGB} & \\ & F_T \end{array} \right] \right) \quad (6)$$

$$= \left[ \begin{array}{cc} \sigma((aW_{RGB} + bW_{share}) * F_{RGB}) & \sigma((cW_{share} + dW_T) * F_T) \end{array} \right]$$

Its structure diagram is shown in Figure 2 (d). In feature extraction, dynamic fusion layer only needs to perform two convolution operations on the visible and infrared inputs respectively to obtain common features and individual features, which greatly reduces computational cost.

The fusion of shared and non-shared convolution kernels is a weighted addition:

$$\begin{cases} \widetilde{W}_{RGB} = aW_{RGB} + bW_{share} \\ \widetilde{W}_T = cW_{share} + dW_T \end{cases} \quad (7)$$

s.t.  $0 < \{a, b, c, d\} < 1$   $a + b = 1$   $c + d = 1$

In this way, the size of the convolution kernels is not changed, and no additional artificially set parameters are introduced. The weights  $a$ ,  $b$ ,  $c$ , and  $d$  are adaptive, which can be obtained from the input through Global Average Pooling (GAP), two-layer full connection (FC), and softmax:

$$\begin{cases} [a, b] = \mathcal{F}(F_{RGB}) \\ [c, d] = \mathcal{F}(F_T) \end{cases} \quad (8)$$

$$\mathcal{F}(X) = \text{Softmax} \circ \text{FC} \circ \text{ReLU} \circ \text{GAP}(X)$$

The weights are input-dependent so that the dynamic fusion layer can use different convolution kernels for different image pairs, which enhances the expressive ability of the model. Compared with the convolution layer without weights, the dynamic fusion layer only increases the Multi-Adds of 0.02%, which can guarantee the real-time performance of the network. For details, please refer to Section 5.4.

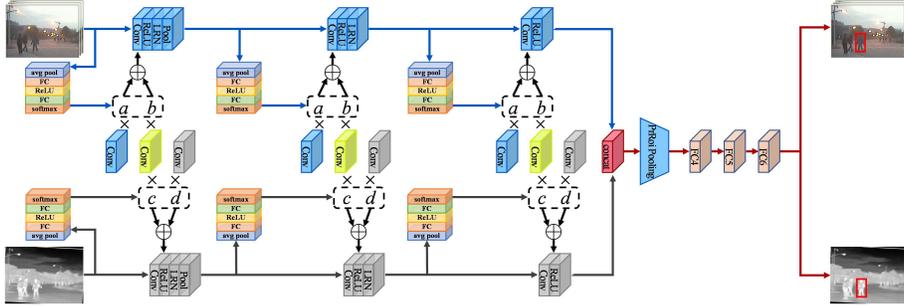


Figure 3: The overall architecture of DFNet. DFNet uses a multi-domain learning framework. Three dynamic fusion layers are used to extract and fuse the features of visible and infrared images. PrRoiPooling is used to unify the features into a  $3 \times 3$  size. Three fully connected layers are used to determine whether the candidate is the object or background.

In the dynamic fusion layer, the convolution kernel is updated through the back-propagation algorithm. In each iteration, the convolution kernel of visible and infrared features iterates as follows:

$$\widetilde{W}_{RGB}^{(i)} = a \left( W_{RGB}^{(i-1)} + \text{lr} \frac{\partial L}{\partial W_{RGB}^{(i-1)}} \right) + b \left( W_{share}^{(2i-1)} + \text{lr} \frac{\partial L}{\partial W_{share}^{(2i-1)}} \right) \quad (9)$$

$$\widetilde{W}_T^{(i)} = c \left( W_{share}^{(2i-2)} + \text{lr} \frac{\partial L}{\partial W_{share}^{(2i-2)}} \right) + d \left( W_T^{(i-1)} + \text{lr} \frac{\partial L}{\partial W_T^{(i-1)}} \right) \quad (10)$$

where  $i$  is the iteration numbers, lr is the learning rate, and  $L$  is the loss function. In each iteration, the non-shared convolution kernels are updated once, and the shared convolution kernel are updated twice.

### 3.2 The Architecture

The overall architecture of DFNet is shown in Figure 3. The features of visible and infrared images are first extracted and fused through three dynamic fusion layers. After PrRoiPooling [43], the features of different sizes are unified into  $3 \times 3$ . Then, features enter three fully connected networks to determine the object or background. At the end of the tracking process, DFNet takes the candidate with the highest network output score as the object:

$$x_t^* = \arg \max \mathcal{F} (x_t^i) \quad (11)$$

where  $x_t^i$  represents the  $i$ -th candidate frame in the  $t$ -th frame,  $\mathcal{F}(\cdot)$  represents the score of network output, and  $x_t^*$  represents the final object result of the  $t$ -th frame.

Inspired by RT-MDNet [44], DFNet adopts a multi-domain learning framework. During training, all video sequences share three dynamic fusion layers, FC4 and FC5. Each video sequence uses a domain-specific FC6. During testing, the multiple domain-specific FC6s are replaced with a reinitialized FC6.

## 4 Implementation Details

We train and test DFNet on the PyTorch platform with i7-10700K CPU and TITAN RTX GPU. We will introduce the details of training and online tracking process in this section.

### 4.1 Offline Training

The pre-trained network in VGG-M [45] is adopted to initialize the model and use ImageNet [46] and RGBT (GTOT [47] or RGBT234 [48]) mixed dataset to train DFNet. We train the network using stochastic gradient descent with momentum. The momentum is set to 0.9, the learning rate is set to  $1e-4$ , and the weight decay is set to  $5e-4$ . The number of epochs is set to 60.

### 4.2 Online Tracking

In the online tracking phase, we initialize the model with the trained three dynamic fusion layers, FC4, and FC5. We reinitialize a new FC6 and use the first frame to train FC6. Specifically, we collect 500 positive samples ( $\text{IOU} > 0.7$ ) and 5000 negative samples ( $\text{IOU} < 0.3$ ) from the first frame as training samples, and use stochastic gradient descent with momentum for training. The momentum is set to 0.9, the learning rate is set to  $1e-4$ , and the weight decay is set to  $5e-4$ .

In the follow-up tracking phase, three dynamic fusion layers are fixed, but FC4, FC5, and FC6 are fine-tuned online. We collect 50 positive samples and 200 negative samples, perform long-term updates every 10 frames, and perform short-term updates when tracking fails. The learning rate of FC6 is set to  $1e-3$ , and the learning rate of FC4 and FC5 is set to  $5e-4$ . At time  $t$ , we use a Gaussian sampler to collect 256 candidates around the object position in the previous frame and calculate their respective classification scores through the network. Then, the Multi-layer Perceptron (MLP) is used to regress the average value of the top five bounding boxes of the classification score to obtain the final bounding box.

## 5 Experiment

### 5.1 Dataset and Evaluation Matrix

We use two RGBT datasets, GTOT [47] and RGBT234, [48] to compare DFNet with other tracking methods. We use ImageNet and GTOT mixed dataset as the training set when evaluating on RGBT234; we use ImageNet and RGBT234

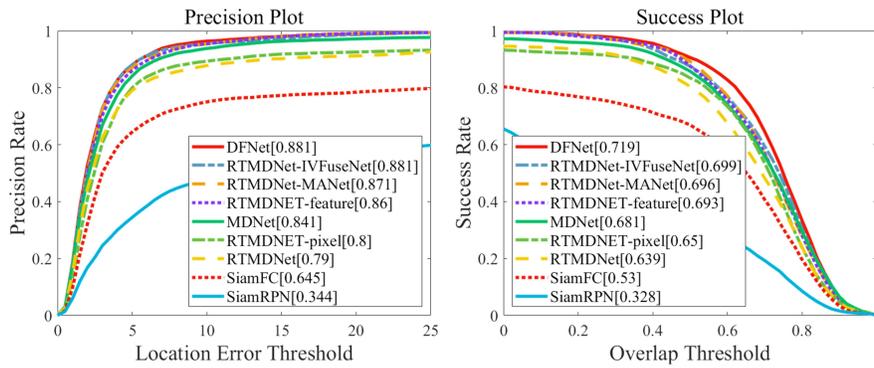
Table 2: RGBT234 dataset PR/SR scores based on attributes. The best, second-best, and third-best PR/SR are shown in red, yellow, and blue.

Tracker	SiamFC	SiamRPN	MDNet	RTMDNet	RTMDNet -pixel	RTMDNet -feature	RTMDNet -IVFuseNet	RTMDNet -MANet	DFNet
BC	0.496/0.333	0.187/0.116	0.683/0.462	0.630/0.402	0.582/0.375	0.705/0.439	0.659/0.397	0.697/0.437	0.714/0.452
CM	0.564/0.407	0.321/0.226	0.689/0.497	0.637/0.438	0.626/0.429	0.690/0.467	0.663/0.438	0.676/0.448	0.692/0.471
DEF	0.591/0.431	0.281/0.212	0.685/0.497	0.654/0.451	0.611/0.419	0.679/0.466	0.651/0.448	0.682/0.445	0.661/0.462
FM	0.518/0.374	0.276/0.155	0.690/0.448	0.679/0.404	0.595/0.330	0.637/0.365	0.618/0.356	0.621/0.374	0.640/0.378
HO	0.521/0.367	0.261/0.164	0.654/0.459	0.634/0.422	0.586/0.369	0.621/0.403	0.592/0.381	0.644/0.409	0.641/0.412
LI	0.495/0.356	0.231/0.154	0.674/0.451	0.605/0.391	0.609/0.404	0.763/0.504	0.742/0.492	0.756/0.497	0.789/0.528
LR	0.603/0.404	0.295/0.159	0.734/0.502	0.683/0.447	0.727/0.464	0.794/0.492	0.787/0.471	0.730/0.446	0.797/0.496
MB	0.554/0.405	0.310/0.209	0.702/0.517	0.669/0.467	0.633/0.442	0.676/0.470	0.670/0.450	0.635/0.433	0.702/0.489
NO	0.765/0.564	0.404/0.282	0.862/0.636	0.842/0.576	0.828/0.557	0.859/0.582	0.856/0.564	0.868/0.569	0.871/0.599
PO	0.629/0.446	0.275/0.188	0.810/0.567	0.754/0.513	0.714/0.508	0.856/0.569	0.838/0.554	0.817/0.544	0.857/0.575
SV	0.634/0.461	0.308/0.210	0.767/0.549	0.747/0.508	0.697/0.461	0.751/0.499	0.743/0.494	0.762/0.505	0.749/0.501
TC	0.681/0.488	0.233/0.155	0.801/0.585	0.763/0.551	0.727/0.487	0.771/0.520	0.737/0.490	0.743/0.490	0.796/0.543
ALL	0.610/0.435	0.295/0.197	0.756/0.536	0.718/0.485	0.691/0.458	0.761/0.504	0.741/0.485	0.756/0.493	0.772/0.513

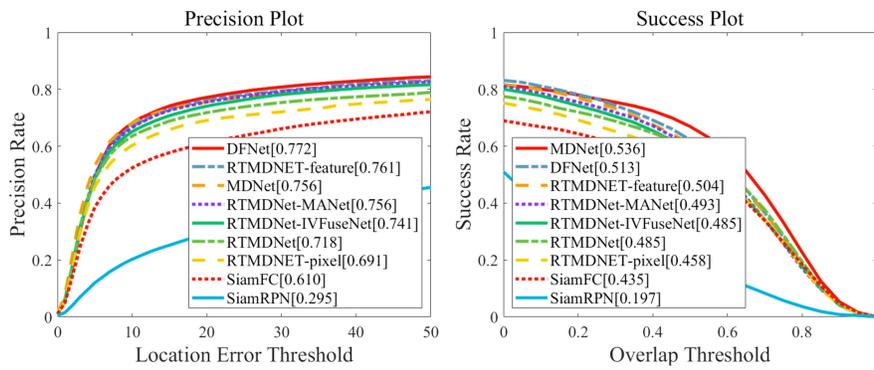
mixed dataset as the training set when evaluating on GTOT. The GTOT and RGBT234 datasets have 50 and 234 RGBT sequences of image pairs aligned in space and time, respectively. In one-pass evaluation (OPE), we use Precision Rate (PR) and Success Rate (SR) as evaluation indicators to evaluate the tracking results. PR refers to the proportion of frames whose difference between the output position and the ground truth bounding box is within the threshold. The thresholds of GTOT is set to 5, The thresholds of RGBT234 is set to 20. SR is the proportion of frames where the overlap ratio between the output position and the ground truth bounding box is greater than the threshold. The area under the curves (AUC) is employed to calculate the SR score.

## 5.2 Comparison with Other Methods

We compared DFNet with visible tracking methods (MDNet [49], RT-MDNet [44], SiamFC [50], and SiamRPN [51]) and fusion tracking methods (pixel-level fusion [8, 9, 10], feature-level fusion [11, 12], MANet [14], and IVFuseNet [34]). To be fair, all fusion methods have been implemented on the RT-MDNet tracking framework, represented below as RTMDNet-pixel, RTMDNet-feature, RTMDNet-MANet, and RTMDNet-IVFuseNet, respectively. The overall tracking performance is shown in Figure 4. For all the indicators of these two benchmarks, our DFNet has clearly outperformed other tracking methods. Specifically, on the GTOT benchmark, our DFNet reached 88.1%/71.9% on PR/SR. While on the RGBT234 benchmark, our DFNet reached PR/SR 77.2%/51.3%. To further show the effectiveness of DFNet, we list the performance of each attribute of the RGBT234 dataset. The specific tracking results are shown in Table 2. It can be concluded from the table that our proposed DFNet outperforms other trackers in 8 cases with higher PR.



(a) comparison on GTOT



(b) comparison on RGBT234

Figure 4: Overall performance compared with other trackers on GTOT (a) and RGBT234 (b).

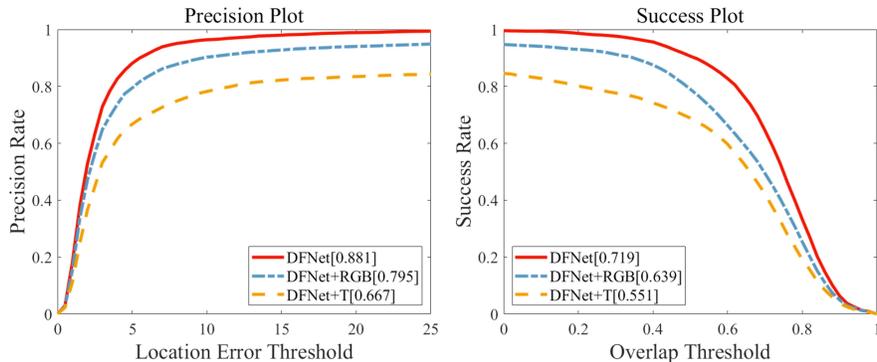


Figure 5: Comparison of visible, infrared, and fusion tracking

## 5.3 Ablation Study

### 5.3.1 The Importance of Fusion

In order to show the importance of fusion for tracking, we compared the tracking performance of DFNet+RGB, DFNet+T, and DFNet. DFNet+RGB and DFNet+T respectively indicate that DFNet solely relies on visible or infrared images for tracking. The tracking performance of DFNet is shown in Figure 5. The PR/SR of RGBT tracking are 8.6%/8.0% higher than tracking using visible image alone, and 21.4%/16.8% higher than tracking using infrared image alone. Experimental results show that the performance of DFNet is significantly better than that of methods based on single-modal images.

### 5.3.2 Dynamic Fusion at Different Layers

We perform a number of ablation studies on RT-MDNet to verify the performance of the dynamic fusion layer at three different layers. The results are shown in Table 3. It can be found that the more the dynamic fusion layer is used, the better the performance is. Using dynamic fusion layers for all three layers produces the best results. And the later the dynamic fusion layer used in the network, the better the performance is.

We visualized the weights of the dynamic fusion layer in the order of the videos in GTOT, as shown in Figure 6. In Figure 6 (a), the blue solid line represents the average of the visible non-shared convolution kernel weights, and the blue shading represents the range of the visible non-shared convolution kernel weights. The red solid line represents the average of the visible shared convolution kernel weights, and the red shade represents the range of the visible shared convolution kernel weights. In Figure 6 (b), the red solid line represents the average of the infrared non-shared convolution kernel weights, and the red shade represents the range of the infrared non-shared convolution kernel weights. The blue solid line represents the average of the infrared shared convolution kernel weights, and the blue shading represents the range of the infrared shared

Table 3: Dynamic fusion at different layers in RT-MDNet.  $\checkmark$  indicates this layer uses a dynamic fusion layer to replace the vanilla convolution, while - indicates not.

Network	C1	C2	C3	PR	SR
Feature-level fusion	-	-	-	0.860	0.693
	$\checkmark$	-	-	0.860	0.689
Feature-level fusion	-	$\checkmark$	-	0.863	0.695
+	-	-	$\checkmark$	0.866	0.691
Dynamic fusion layer	$\checkmark$	$\checkmark$	-	0.867	0.699
	$\checkmark$	-	$\checkmark$	0.866	0.699
	-	$\checkmark$	$\checkmark$	0.872	0.702
DFNet	$\checkmark$	$\checkmark$	$\checkmark$	0.881	0.709

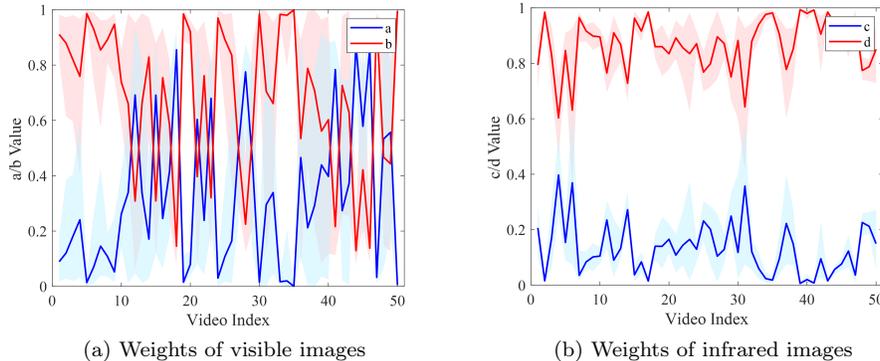


Figure 6: Weights for different video sequences across the GTOT dataset in DFNet.

convolution kernel weights. It can be found that the dynamic fusion layer can calculate different weights according to videos. In this way, the dynamic fusion layer makes the fusion tracker adaptively calculate the contributions of individual features and common features.

In addition, we visualized the weights of the dynamic fusion layer in two video sequences, as shown in Figure 7. Figure 7 (a) is from *OccCar-2*. It can be found that in this video sequence, the contributions of individual features and common features are also different. At the beginning of the video, the car is clear in the visible images, and the contributions of individual features of visible images are large. When the car is blocked by leaves, visible images cannot clearly distinguish the car, so the contributions of the individual features reduce, while the contributions of the common features increase. As the car comes out of the leaves, the contributions of individual features of visible images increase. Figure 7 (b) is from *FastMotorNig*. When the bicycle is blocked by a street light, the

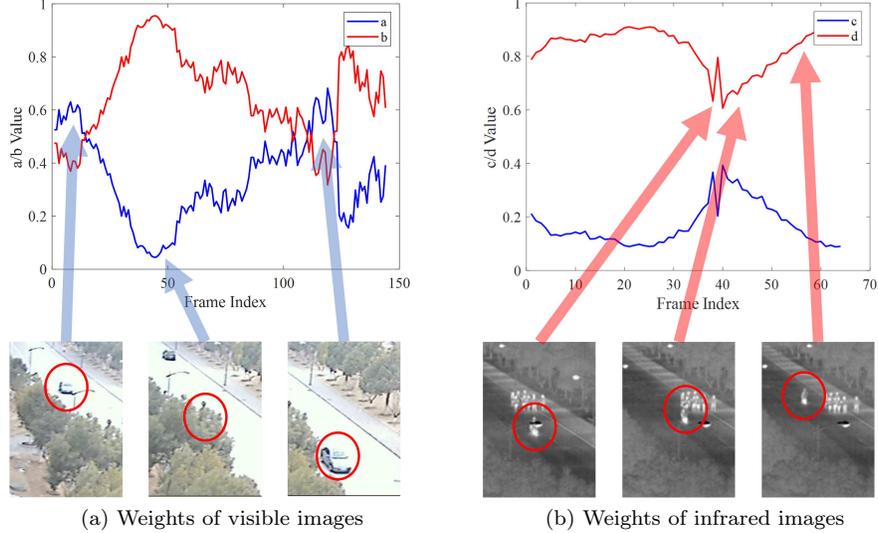


Figure 7: Dynamic fusion layer weights of different frames of (a) *OccCar-2* and (b) *FastMotorNig*

contributions of the individual features reduce, while the contributions of the common features increase. In the later stage of the video sequence, because the bicycle is close to the crowd, *thermal crossover* [47] occurs, so the contributions of the individual features of the infrared images are low. As the bicycle moves away from the crowd, the contributions of the individual features of the infrared images increase.

#### 5.4 Efficiency Analysis of different fusion methods

The speed of DFNet is 28.658 FPS. We compared the speed and performance of DFNet with other fusion tracking methods, the results are shown in Figure 8. The computational cost of DFNet is  $O(X) = 2(HWC_{in} + C_{in}C_{hidden} + 2C_{hidden} + HWC_{in}C_{out}k^2)$  Mult-Adds, where  $H$ ,  $W$  are the height and width of the input.  $C_{in}$ ,  $C_{out}$ , and  $C_{hidden}$  are the channel numbers of the input, output, and hidden layer, respectively.  $k$  is the size of the convolution kernel. Correspondingly, the computational cost of baseline (RTMDNet-feature) is  $O(X) = 2HWC_{in}C_{out}k^2$  Mult-Adds. Since the fusion of shared and non-shared convolution kernels is performed in convolution kernel space, compared with the non-shared-convolution-kernel-based fusion method, no additional calculations to increase. The increase of computational cost is mainly due to the attention which calculates weights according to the input. The computational cost caused by the attention is much smaller than convolution. In DFNet, it is less than 0.02%. While, compared with the baseline, MANet fuses the shared and non-

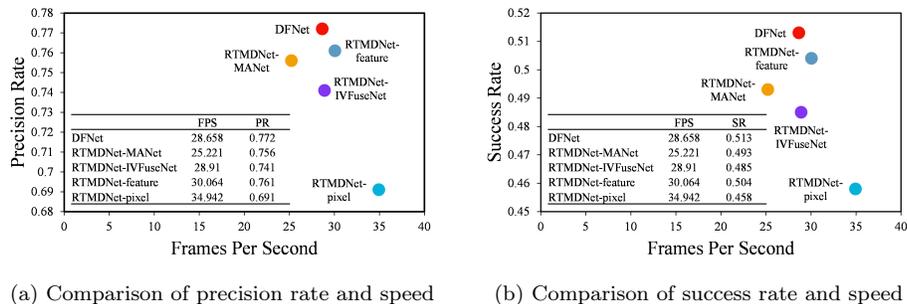


Figure 8: Comparison of speed and PR/SR.

Table 4: The computational cost of different fusion methods. C1, C2, C3, and total indicate the computational cost of the first, second, third, and all convolutional layers, respectively. Percent indicate computational cost expressed as a percentage of the non-shared-convolution-kernel-based fusion method.

Model	C1	C2	C3	total	percent
RTMDNet-feature (baseline)	323.14M	768.00M	285.47M	1376.61M	100.00%
RTMDNet-IVFuseNet	323.14M	768.00M	285.47M	1376.61M	100.00%
RTMDNet-MANet	382.49M	798.72M	317.19M	1498.40M	108.85%
DFNet (ours)	323.21M	768.12M	285.57M	1376.90M	100.02%

shared features in feature space, which causes calculations to increase by 8.85%. The specific computational cost is shown in the Table 4.

Based on all the experiments performed in this section, we conclude that:

1. Compared with the visible tracking method (MDNet, RT-MDNet, SiamFC, and SiamRPN) and the fusion tracking method (pixel-level fusion, feature-level fusion, MANet, and IVFuseNet), DFNet achieves the best PR and SR.
2. The performance of fusion method is better than that of methods based on single-modal images, which shows the advantage of fusion.
3. With consideration of the contributions of individual features and common features, DFNet can adaptively calculate the weights of shared and non-shared convolution kernels to cope with changes in modality reliability.
4. Compared with the fusion of shared and non-shared features in the feature space, the fusion of shared and non-shared convolution kernels in the convolution kernel space can effectively reduce the computational complexity and improve the tracking speed.

## 6 Conclusion

In this paper, we propose a novel RGBT tracking method, called dynamic fusion network (DFNet). DFNet is essentially a feature-level fusion method, which can use non-shared convolutions to respectively extract individual features according to the different characteristics of visible and infrared images. Furthermore, DFNet takes the advantage of shared convolution kernels to extract common features. In addition, because attention is used to adaptively calculate different convolution kernel weights according to inputs, DFNet can dynamically calculate the contributions of individual features and common features in the face of changes in modality reliability. The shared convolution kernels and non-shared convolution kernels are concatenated in convolution kernel space, so that, the computational cost is small. Extensive experiments on two RGBT datasets validate the effectiveness of DFNet. Future work will focus on adopting more advanced architectures, designing other adaptive weighting methods, and reducing the redundancy of features between different modalities.

## References

- [1] Sulan Zhai, Pengpeng Shao, Xinyan Liang, and Xin Wang. Fast RGB-T tracking via cross-modal correlation filters. *Neurocomputing*, 334:172–181, 2019.
- [2] Chaoqun Wang, Chunyan Xu, Zhen Cui, Ling Zhou, Tong Zhang, Xiaoya Zhang, and Jian Yang. Cross-modal pattern-propagation for RGB-T tracking. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7062–7071, June 2020.
- [3] Xingchen Zhang, Ping Ye, Henry Leung, Ke Gong, and Gang Xiao. Object fusion tracking based on visible and infrared images: A comprehensive review. *Information Fusion*, 63:166–187, 2020.
- [4] Amanda Berg, Jorgen Ahlberg, and Michael Felsberg. A thermal object tracking benchmark. In *2015 12th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–6, Aug 2015.
- [5] Xin Li, Qiao Liu, Nana Fan, Zhenyu He, and Hongzhi Wang. Hierarchical spatial-aware siamese network for thermal infrared object tracking. *Knowledge-Based Systems*, 166:71–81, 2019.
- [6] Qiao Liu, Xiaohuan Lu, Zhenyu He, Chunkai Zhang, and Wen-Sheng Chen. Deep convolutional neural networks for thermal infrared object tracking. *Knowledge-Based Systems*, 134:189–198, 2017.
- [7] Lichao Zhang, Abel Gonzalez-Garcia, Joost van de Weijer, Martin Danelljan, and Fahad Shahbaz Khan. Synthetic data generation for end-to-end thermal infrared tracking. *IEEE Transactions on Image Processing*, 28(4):1837–1850, April 2019.

- [8] Ningwen Xu, Gang Xiao, Xingchen Zhang, and Durga Prasad Bavirisetti. Relative object tracking algorithm based on convolutional neural network for visible and infrared video sequences. In *Proceedings of the 4th International Conference on Virtual Reality*, pages 44–49, 2018.
- [9] Ningwen Xu, Gang Xiao, Fang He, Xingchen Zhang, and Durga Prasad Bavirisetti. Object tracking via deep multi-view compressive model for visible and infrared sequences. In *2018 21st International Conference on Information Fusion (FUSION)*, pages 941–948, 2018.
- [10] Xingchen Zhang, Ping Ye, Dan Qiao, Junhao Zhao, Shengyun Peng, and Gang Xiao. Object fusion tracking based on visible and infrared images using fully convolutional siamese networks. In *2019 22th International Conference on Information Fusion (FUSION)*, pages 1–8, July 2019.
- [11] Xingming Zhang, Xuehan Zhang, Xuedan Du, Xiangming Zhou, and Jun Yin. Learning multi-domain convolutional network for RGB-T visual tracking. In *2018 11th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI)*, pages 1–6, Oct 2018.
- [12] Xingchen Zhang, Ping Ye, Shengyun Peng, Jun Liu, Ke Gong, and Gang Xiao. Siamft: An RGB-infrared fusion tracking method via fully convolutional siamese networks. *IEEE Access*, 7:122122–122133, 2019.
- [13] Chenglong Li, Xiaohao Wu, Nan Zhao, Xiaochun Cao, and Jin Tang. Fusing two-stream convolutional neural networks for RGB-T object tracking. *Neurocomputing*, 281:78–85, 2018.
- [14] Cheng Long Li, Andong Lu, Ai Hua Zheng, Zhengzheng Tu, and Jin Tang. Multi-adaptor RGBT tracking. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 2262–2270, Oct 2019.
- [15] Pengyu Zhang, Jie Zhao, Chunjuan Bo, Dong Wang, Huchuan Lu, and Xiaoyun Yang. Jointly modeling motion and appearance cues for robust RGB-T tracking. *IEEE Transactions on Image Processing*, 30:3335–3347, 2021.
- [16] Xiao Yun, Zhongliang Jing, Gang Xiao, Bo Jin, and Canlong Zhang. A compressive tracking based on time-space kalman fusion model. *Science China Information Sciences*, 59(1):1–15, Jan 2016.
- [17] N. Cvejic, S. G. Nikolov, H. D. Knowles, A. Loza, A. Achim, D. R. Bull, and C. N. Canagarajah. The effect of pixel-level fusion on object tracking in multi-sensor surveillance video. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–7, 2007.

- [18] F. Bunyak, K. Palaniappan, S. K. Nath, and G. Seetharaman. Geodesic active contour based fusion of visible and infrared video for persistent object tracking. In *2007 IEEE Workshop on Applications of Computer Vision (WACV '07)*, pages 35–35, 2007.
- [19] Supriya Mangale and Madhuri Khambete. Camouflaged target detection and tracking using thermal infrared and visible spectrum imaging. In *Intelligent Systems Technologies and Applications 2016*, pages 193–207, 2016.
- [20] Xue Mei and Haibin Ling. Robust visual tracking using l1 minimization. In *2009 IEEE 12th International Conference on Computer Vision*, pages 1436–1443, 2009.
- [21] Liying Li, J.P. McGinnis, and Kausik Si. Translational control by prion-like proteins. *Trends in Cell Biology*, 28(6):494–505, 2018.
- [22] Xiangyuan Lan, Mang Ye, Shengping Zhang, Huiyu Zhou, and Pong C. Yuen. Modality-correlation-aware sparse representation for RGB-infrared object tracking. *Pattern Recognition Letters*, 130:12–20, 2020.
- [23] Chenglong Li, Nan Zhao, Yijuan Lu, Chengli Zhu, and Jin Tang. Weighted sparse representation regularized graph learning for RGB-T object tracking. In *Proceedings of the 25th ACM International Conference on Multimedia*, pages 1856–1864, 2017.
- [24] Chenglong Li, Chengli Zhu, Yan Huang, Jin Tang, and Liang Wang. Cross-modal ranking with soft consistency and noisy labels for robust RGB-T tracking. In *Computer Vision – ECCV 2018*, pages 831–847, 2018.
- [25] Juxin Li, Weizhi Liu, Giulia Pedrielli, Loo Hay Lee, and Ek Peng Chew. Optimal computing budget allocation to select the nondominated systems—a large deviations perspective. *IEEE Transactions on Automatic Control*, 63(9):2913–2927, 2018.
- [26] Chenglong Li, Chengli Zhu, Shaofei Zheng, Bin Luo, and Jing Tang. Two-stage modality-graphs regularized manifold ranking for RGB-T tracking. *Signal Processing: Image Communication*, 68:207–217, 2018.
- [27] Ying Tai, Jian Yang, and Xiaoming Liu. Image super-resolution via deep recursive residual network. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2790–2798, 2017.
- [28] Guang Chen, Haitao Wang, Kai Chen, Zhijun Li, Zida Song, Yinlong Liu, Wenkai Chen, and Alois Knoll. A survey of the four pillars for small object detection: Multiscale representation, contextual information, super-resolution, and region proposal. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, pages 1–18, 2020.
- [29] Jingchao Peng, Haitao Zhao, Zhengwei Hu, Yi Zhuang, and Bofan Wang. Siamese infrared and visible light fusion network for RGB-T tracking, 2021.

- [30] Chenglong Li, Lei Liu, Andong Lu, Qing Ji, and Jin Tang. Challenge-aware rgbt tracking. In *Computer Vision – ECCV 2020*, pages 222–237, 2020.
- [31] Xingchen Zhang, Ping Ye, Shengyun Peng, Jun Liu, and Gang Xiao. DSiamMFT: An RGB-T fusion tracking method via dynamic siamese networks using multi-layer feature fusion. *Signal Processing: Image Communication*, 84:115756, 2020.
- [32] Yabin Zhu, Chenglong Li, Jin Tang, and Bin Luo. Quality-aware feature aggregation network for robust RGBT tracking. *IEEE Transactions on Intelligent Vehicles*, 6(1):121–130, March 2021.
- [33] Zhangyang Wang, Shiyu Chang, Yingzhen Yang, Ding Liu, and Thomas S. Huang. Studying very low resolution recognition using deep networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4792–4800, June 2016.
- [34] Yuqi Li, Haitao Zhao, Zhengwei Hu, Qianqian Wang, and Yuru Chen. IVFuseNet: Fusion of infrared and visible light images for depth prediction. *Information Fusion*, 58:1–12, 2020.
- [35] Jihong Pei, Mi Zou, and Lixia Wang. A local adaptive threshold noise detection linear interpolation filter (lalif) for stripe noise removal in infrared images. In *2016 IEEE 13th International Conference on Signal Processing (ICSP)*, pages 681–686, 2016.
- [36] Wei Hu, Yiqi Zhuang, Junlin Bao, and Qifeng Zhao. Effects of radiation-induced changes in low-frequency noise of infrared detectors. In *2015 International Conference on Noise and Fluctuations (ICNF)*, pages 1–3, 2015.
- [37] Brandon Yang, Gabriel Bender, Quoc V Le, and Jiquan Ngiam. CondConv: Conditionally parameterized convolutions for efficient inference. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [38] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11027–11036, June 2020.
- [39] Changqian Yu, Yuanjie Shao, Changxin Gao, and Nong Sang. CondNet: Conditional classifier for scene segmentation. *IEEE Signal Processing Letters*, 28:758–762, 2021.
- [40] Hao Jiang, Siqi Wang, Huikun Bi, Xiaolei Lv, Binqiang Zhao, Zheng Wang, and Zhaoqi Wang. Synthesizing indoor scene layouts in complicated architecture using dynamic convolution networks. *Proceedings of the ACM on Computer Graphics and Interactive Techniques*, 4(1):16, april 2021.

- [41] Xin Ma, Xiaoqiang Zhou, Huaibo Huang, Zhenhua Chai, Xiaolin Wei, and Ran He. Free-form image inpainting via contrastive attention network. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 9242–9249, 2021.
- [42] Yan Li, Guitao Cao, and Wenming Cao. A dynamic group equivariant convolutional networks for medical image analysis. In *2020 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 1056–1062, 2020.
- [43] Borui Jiang, Ruixuan Luo, Jiayuan Mao, Tete Xiao, and Yuning Jiang. Acquisition of localization confidence for accurate object detection. In *Computer Vision – ECCV 2018*, pages 816–832, 2018.
- [44] Ilchae Jung, Jeany Son, Mooyeol Baek, and Bohyung Han. Real-time MD-Net. In *Computer Vision – ECCV 2018*, pages 89–104, 2018.
- [45] Ken Chatfield, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *Proceedings of the British Machine Vision Conference*, 2014.
- [46] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, June 2009.
- [47] Chenglong Li, Hui Cheng, Shiyi Hu, Xiaobai Liu, Jin Tang, and Liang Lin. Learning collaborative sparse representation for grayscale-thermal tracking. *IEEE Transactions on Image Processing*, 25(12):5743–5756, Dec 2016.
- [48] Chenglong Li, Xinyan Liang, Yijuan Lu, Nan Zhao, and Jin Tang. RGB-T object tracking: Benchmark and baseline. *Pattern Recognition*, 96:106977, 2019.
- [49] Hyeonseob Nam and Bohyung Han. Learning multi-domain convolutional neural networks for visual tracking. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4293–4302, June 2016.
- [50] Luca Bertinetto, Jack Valmadre, João F. Henriques, Andrea Vedaldi, and Philip H. S. Torr. Fully-convolutional siamese networks for object tracking. In *Computer Vision – ECCV 2016 Workshops*, pages 850–865, 2016.
- [51] Bo Li, Junjie Yan, Wei Wu, Zheng Zhu, and Xiaolin Hu. High performance visual tracking with siamese region proposal network. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8971–8980, 2018.