# Multi-domain Collaborative Feature Representation for Robust Visual Object Tracking

**Jiqing Zhang**[1,†] · **Kai Zhao**[2,†] · **Bo Dong**[3] · **Yingkai Fu**[1] · **Yuxin Wang**[1] · **Xin Yang**[1,*] · **Baocai Yin**[1]

**Abstract** Jointly exploiting multiple different yet complementary domain information has been proven to be an effective way to perform robust object tracking. This paper focuses on effectively representing and utilizing complementary features from the frame domain and event domain for boosting object tracking performance in challenge scenarios. Specifically, we propose Common Features Extractor (CFE) to learn potential common representations from the RGB domain and event domain. For learning the unique features of the two domains, we utilize a Unique Extractor for Event (UEE) based on Spiking Neural Networks to extract edge cues in the event domain which may be missed in RGB in some challenging conditions, and a Unique Extractor for RGB (UER) based on Deep Convolutional Neural Networks to extract texture and semantic information in RGB domain. Extensive experiments on standard RGB benchmark and real event tracking dataset demonstrate the effectiveness of the proposed approach. We show our approach outperforms all compared state-of-the-art tracking algorithms and verify event-based data is a powerful cue for tracking in challenging scenes.

Jiqing Zhang, Yingkai Fu, Yuxin Wang, Xin Yang and Baocai Yin
Email: jqz@mail.dlut.edu.cn, yingkaifu@mail.dlut.edu.cn, wyx@dlut.edu.cn, xinyang@dlut.edu.cn and ybc@dlut.edu.cn

Kai Zhao
Email: kzhao@aiit.org.cn

Bo Dong
Email: bo.dong@sri.com

[1] Dalian University of Technology · [2] Advanced Institute of Information Technology, Peking University · [3] SRI International

† indicates equal contribution.
* Xin Yang is the Corresponding Author.

**Fig. 1** Visual examples of our tracker comparing with other three state-ot-the-art trackers including GradNet [32], MDNet [38], and SiamDW-RPN [72] on *Ironman* from OTB2013 [59]. *Ironman* is a challenging sequence with low light, motion blur, background clutters, and fast motion. Best viewed in zoom in.

**Keywords** Visual object tracking · Event-based camera · Multi-domain · Challenging conditions

## 1 Introduction

Visual object tracking is an important topic in computer vision, where the target object is identified in the first frame and tracked in all frames of a video. Due to the significant learning ability, deep convolutional neural networks (DCNNs) have been widely used to object detection [35, 34, 62], image matting [43, 42, 64], super-resolution [68, 67, 63], image enhancement [61, 65] and visual object tracking [2, 15, 19, 28, 32, 72, 38, 13, 70, 71, 22, 12, 11, 33, 58, 47]. However, RGB-based trackers suffer from bad environmental conditions, *e.g.*, low illumination, fast motion, and so on. Some works [52, 25, 24, 31, 30, 74] try to introduce additional information (*e.g.*, depth and thermal infrared) to improve tracking performance. However, when the tracking target is in a high-speed motion or an environment with a wide dynamic range, these sensors usually cannot provide satisfactory results.

Event-based cameras are bio-inspired vision sensors whose working principle is entirely different from traditional cameras. While conventional cameras obtain intensity frames at a fixed rate, event-based cameras measure light intensity changes and output events asynchronously. Compared with conventional cameras, event-based cameras have several advantages. First, with a high temporal resolution (around 1 $\mu s$), event-based cameras do not suffer from motion blur. Second, event-based cameras have a high dynamic range (*i.e.*, 120-140 dB). Thus they can work effectively even under over/under-exposure conditions.

We observe that events and RGB data are captured from different types of sensors, but they share some similar information like target boundaries. At the same time, stacked event images and RGB images have their own unique characteristics. In particular, RGB images contain rich low- and high-frequency texture information and provide abundant representations for describing target objects. Events can provide target edge cues that are not influenced by motion blur and bad illumination. Therefore, event-based data and RGB images are complementary, which calls for the development of novel algorithms capable of combining the specific advantages of both domains to perform computer vision in degraded conditions.

To the best of our knowledge, we are the first to jointly explore RGB and events for object tracking based on their similarities and differences in an end-to-end manner. This work is essentially object tracking with multi-modal data that includes RGB-D tracking [52, 25, 60, 24], RGB-T tracking [31, 27, 30, 74, 69], and so on. However, since the output of an event-based camera is an asynchronous stream of events, this makes event-based data fundamentally different from other sensors' data that have been addressed well by multi-model tracking methods. With the promise of increased computational ability and low power computation using neuromorphic hardware, Spiking Neural Networks (SNNs), a processing model aiming to improve the biological realism of artificial neural networks, show their potential as computational engines, especially for processing event-based data from neuromorphic sensors. Therefore, combining SNNs and DCNNs to process multi-domain data is worth exploring.

In this paper, focusing on the above two points, we propose Multi-domain Collaborative Feature Representation (MCFR) that can effectively extract the common features and unique features from both domains for robust visual object tracking in challenging conditions. Specifically, we employ the first three convolutional layers of VGGNet-M [51] as our Common Features Extractor (CFE) to learn similar potential rep-

resentations from the RGB domain and event domain. To model specific characteristics of each domain, the Unique Extractor for RGB (UER) is designed to extract unique texture and semantic features in the RGB domain. Furthermore, we leverage the Unique Extractor for Events (UEE) based on SNNs to efficiently extract edge cues in the event domain. Extensive experiments on the RGB benchmark and real event dataset suggest that the proposed tracker achieves outstanding performance. A visual example can be seen in Figure 1, which contains multiple challenging attributes. By analyzing quantitative results, we provide basic insights and identify the potentials of events in visual object tracking.

To sum up, our contributions are as follows:

• We propose a novel multi-domain feature representation network which can effectively extract and fuse the information from frame and event domains.

• We preliminarily explore combining SNNs and DCNNs for visual object tracking.

• The extensive experiments verify our approach outperforms other state-of-the-art methods. The ablation studies evidence the effectiveness of the designed components.

## 2 Related Work

### 2.1 Spiking Neural Networks

Spiking Neural Networks (SNNs) are bio-inspired models using spiking neurons as computational models. The inputs of spiking neurons are temporal events called spikes, and the outputs also are spikes. Spiking neurons have a one-dimensional internal state named potential, which is controlled by first-order dynamics. Whenever a spike arrives, if no other spikes are recorded in time, the potential will be excited but will decay again. When the potential reaches a certain threshold, the spiking neuron sends a spike to the connected neurons and resets its own potential. It has been shown that such networks are able to process asynchronous, without pre-processing events data [10, 16]. Since the spike generation mechanism cannot be differentiated and the spikes may introduce the problem of incorrect allocation of the time dimension, the traditional gradient backpropagation mechanism cannot be directly used in SNNs. Nonetheless, some researches [39, 66, 50, 48, 54, 17] on supervised learning for SNNs has taken inspiration from backpropagation to solve the error assignment problem. However, it is still unclear how to train multiple layers of SNNs, and combine them with DCNNs for tracking task.

## 2.2 Single-Domain Tracking

**RGB-based tracking.** Deep-learning-based methods have dominated the visual object tracking field, from the perspective of either one-shot learning [2,15,19,28, 32,72] or online learning [38,13,70,71,22,12,11,33]. Usually, the latter methods are more accurate (with less training data) but slower than the former ones. Among them, Nam *et al.* [38] proposed the Multi-Domain Network (MDNet), which used a CNN-based backbone pretrained offline to extract generic target representations, and the fully connected layers updated online to adapt temporal variations of target objects. In MDNet [38], each domain corresponds to one video sequence. Due to the effectiveness of this operation in visual tracking, we follow this idea to ensure the accuracy of tracking.

**Event-based tracking.** Compared with the frame-based object tracking methods, there are only a few works on event-based object tracking [41,36,73,44,1, 53,7]. Piatkowska *et al.* [41] presented a Gaussian mixture model to track the pedestrian motion. Barranco *et al.* [1] proposed a real-time clustering algorithm and used Kalman filters to smooth the trajectories. Zhu *et al.* [73] monitored the confidence of the velocity estimate and triggered a tracking command once the confidence reaches a certain threshold. Ramesh *et al.* [44] presented a long-term object tracking framework with a moving event camera under general tracking conditions. Mitrokhin *et al.* [36] proposed a motion compensation method for tracking objects by getting the possible areas that are not consistent with camera motion. Timo.S *et al.* [53] calculated the optical flow from the events at first, then warped the events' position to get the sharp edge event images according to the contrast principle. Besides, they gave each event a weight as its probability and fused them during the process of warping so that they can classify events into different objects or background. Chen *et al.* [7] proposed an end-to-end retinal motion regression network to regress 5-DoF motion features.

Although the above studies have achieved good performance in the RGB domain or the event domain, they ignore exploring the complementary information existing between the two domains. As a consequence, we investigate the similarities and differences between the event and RGB domain, and propose common features extractor and unique feature extractor to learn and fuse valuable complementary features.

## 2.3 Multi-Domain Tracking

The current popular visual object tracking based on multi-domain data mainly includes RGB-D (RGB +

depth) tracking [52,60,24,25] and RGB-T (RGB + thermal) tracking [31,27,30,74,69]. Depth cues are usually introduced to solve the occlusion problem in visual object tracking. Images from the thermal infrared sensors are not influenced by illumination variations and shadows, and thus can be combined with RGB to improve performance in bad environmental conditions. As the output of an event camera is an asynchronous stream of events, this makes raw event stream fundamentally different from other sensors data that have been addressed well by the above state-of-the-art multi-model tracking methods. Therefore, it is essential to design a tailored algorithm for leveraging RGB data and event data simultaneously.

## 3 Methodology

### 3.1 Backgroud: Event-based Camera

An event-based camera is a bio-inspired sensor. It asynchronously measures light intensity changes in scene-illumination at a pixel level. Therefore, it provides a very high-temporal resolution (*i.e.*, up to 1MHz). Due the light intensity changes are measured in the log scale, an event-based camera can offer a very high dynamic range (*i.e.*, up to 140 dB). An event was triggered when the change of a log-scale pixel intensity is higher or lower than a threshold, resulting in an "ON" and an "OFF" event, respectively. Mathematically, a set of events can be defined as:

$$\mathcal{E} = \{e_k\}_{k=1}^N = \{[x_k, y_k, t_k, p_k]\}_{k=1}^N, \tag{1}$$

where $e_k$ is the $k$th event; $[x_k, y_k]$ is the pixel location of event $e_k$; $t_k$ is the timestamp when the event is triggered; $p_k \in \{-1, 1\}$ is the polarity of an event, where $-1$ and $1$ represent OFF and ON events, respectively. In a constant lighting condition, events are normally triggered by moving edges (*e.g.*, object contour, texture and depth discontinuities), which makes an event-based camera be a natural edge extractor. Therefore, with these unique features, event-based cameras have been introduced to various tasks [53,55,4,26,9,40,6,37] in challenging scenes (*e.g.*, low-light, fast motion).

Even though event-based cameras are sensitive to edges, they cannot provide absolute intensity and texture information. Besides, since the asynchronous event stream differs significantly from the frames generated by conventional frame-based cameras, vision algorithms designed for frame-based cameras cannot be directly applied. To deal with it, events are typically aggregated into a grid-based representation first.
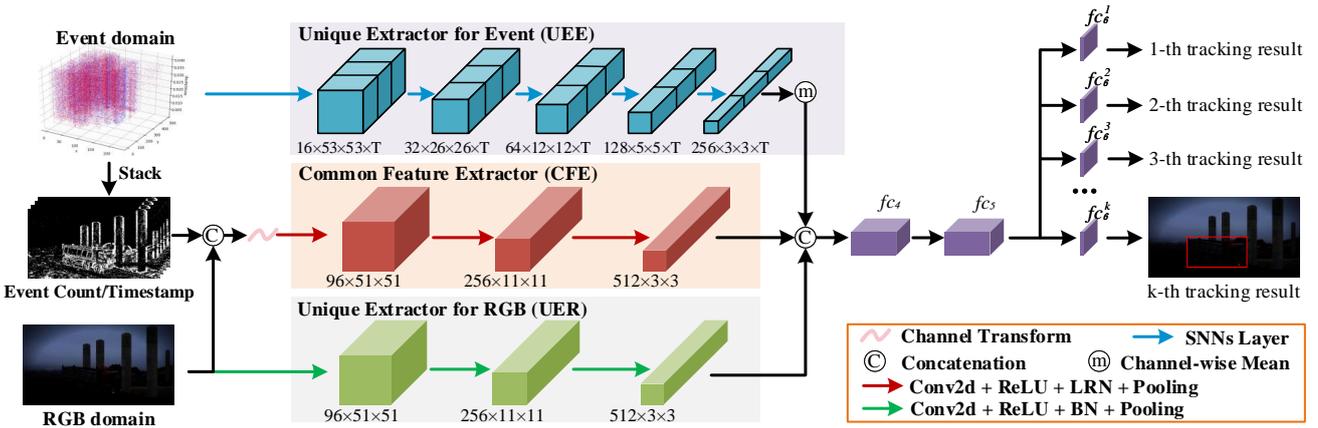
**Fig. 2** The overview of our proposed network. Our pipeline mainly consists of three parts, UEE for extracting special features from event domain, UER for extracting unique features from RGB domain, and CFE for extracting common shared features from both domains. The target is a moving truck in underexposure.

## 3.2 Network Overview

Our approach builds on two key observations. First, although events and RGB data are captured from different types of sensors, they share some similar information, such as target object boundaries. Similar features should be extracted using a consistent strategy. Second, rich textural and semantic cues can be easily captured by a conventional frame-based sensor. In contrast, an event-based camera can easily capture edge information which may be missed in RGB images under some challenging conditions. Therefore, fusing complementary advantages of each domain will enhance feature representation. Figure 2 illustrates the proposed Multi-domain Collaborative Feature Representation (MCFR) for robust visual object tracking. Specifically, for the first observation, we propose the Common Feature Extractor (CFE) which accepts stacked event images and RGB images as inputs to explore shared common features. For the second observation, we design a Unique Extractor for Event (UEE) based on SNNs to extract edge cues in the event domain which may be missed in the RGB domain under some challenging conditions, and a Unique Extractor for RGB (UER) based on DC-NNs to extract texture and semantic information in the RGB domain. The outputs of UEE, CFE, and UER are then concatenated, and a convolutional layer with $1 \times 1$ kernel size is used to adaptively select valuable combinative features. Finally, the combinative features are classified by three fully connected layers and softmax cross-entropy loss. Following [38], the network has $k$ branches, which are denoted by the last fully connected layers. In other words, training sequences $fc6^1 - fc6^k$.

## 3.3 Common Feature Extractor

To leverage a consistent scheme for extracting similar features of event and RGB domains, we first stack event stream according to the counts and latest timestamp of positive and negative polarities, which makes vision algorithms designed for frames can also be applied to asynchronous event streams. Mathematically,

$$C(x,y,p) = \sum_{k=1, t_k \in W}^{N} \delta\left(x_k, x\right) \delta\left(y_k, y\right) \delta\left(p_k, p\right)$$
$$T(x,y,p) = \max_{t_k \in W} t_k \delta\left(x_k, x\right) \delta\left(y_k, y\right) \delta\left(p_k, p\right) \tag{2}$$

where $\delta$ is the Kronecker delta function, $W$ is the time window (the interval between adjacent RGB frames), and $N$ is the number of events that occurred within $W$. The stacked event count image $C$ contains the number of events at each pixel, which implies the frequency and density information of targets. The stacked event timestamp image $T$ contains the temporal cues of the motion, which implies the direction and speed information of targets. An example of counts images and timestamp images is shown in Figure 3, we find that the stacked event images and RGB image indeed share some common features, such as the edge cues of targets.

We then employ a Common Feature Extractor (CFE) to extract shared object representations across different domains. To balance effectiveness and efficiency, we apply the first three layers from the VGGNet-M [51] as the main feature extraction structure of our CFE. Specifically, the convolution kernel sizes are $7 \times 7$, $5 \times 5$, and $3 \times 3$, respectively. The output channels are 96, 256, and 512, respectively. As shown in Figure 2, the whole process is formulated as $F_{cfe} = CFE(\tau([RGB, C, T]))$, where $RGB$ denotes RGB image, $[\cdot]$ is concatenation, $\tau$

indicates channel transformation, and $F_{cfe}$ is the output of CFE.

## 3.4 Unique Extractor for RGB

Since the raw event stream and RGB data storage methods and expressions are different, it is necessary to design an exclusive feature extraction method for each domain. For the RGB domain, we propose Unique Extractor for RGB (UER) to effectively extract unique texture and semantic features. Specifically, as shown in Figure 2, UER consists of three convolutional layers, and the size of the convolution kernel are $3 \times 3$, $1 \times 1$, and $1 \times 1$, respectively. It is noted that one major difference between UER and CFE is the size of the convolution kernel. CFE employs large-size convolution kernels to provide a larger receptive field so that the whole boundary from RGB and event domains can be better extracted, while UER can focus on the rich texture information in the RGB domain with small-size kernels. This process can be simply formulated as $F_{uer} = UER(RGB)$, where $F_{uer}$ is the output of UER.

## 3.5 Unique Extractor for Event

Compared with RGB images, the event-based data is not affected by HDR and motion blur. Besides, from Figure 3, we can see that events can provide clear cues about where object movement occurred, which will help the tracking process not be disturbed by the surrounding environment. Since SNNs can process raw event stream directly, we introduce it into our Unique Extractor for Events (UEE) (top branch in Figure 2) to effectively extract unique event features. There are different mathematical models to describe the dynamics of a spiking neuron, we use the Spike Response Model (SRM) [18] in this work. In the SRM [18], the net effect that firing has on the emitting and the receiving neuron is described by two response functions, $v(t)$ and $u(t)$. The refractory function $u(t)$ describes the response of the firing neuron to its own spike. The synaptic kernel $v(t)$ describes the effect of an incoming spike on the membrane potential at the soma of the postsynaptic neuron. Following [17,49], we define the feedforward SNNs with $n$ layers as:

$$v(t) = \frac{t}{\tau_s} e^{1 - \frac{t}{\tau_s}} H(t), \quad u(t) = -2\phi e^{-\frac{t}{\tau_r}} H(t) \quad (3)$$

$$\varepsilon_{i+1}(t) = W_i(v * s_i)(t) + (u * s_{i+1})(t) \quad (4)$$

$$s_i(t) = \sum \delta(t - t_i) \quad (t_i \in \{t | \varepsilon_i(t) = \phi\}) \quad (5)$$

$$F_{uee} = \mathcal{M}(W_n(v * s_n)(t)) \quad (6)$$

where $H$ is the Heaviside step function; $\tau_s$ and $\tau_r$ are the time constants of the synaptic kernel and refractory kernel, respectively. $s_i$ and $W_i$ are the input spikes and synaptic weights of the $i$th layer, respectively. $\phi$ denotes the neuron threshold, that means, when the sub-threshold membrane potential is strong enough to exceed $\phi$ the spiking neuron responds with a spike. To combine SNNs with DCNNs in the overall structure, we perform a mean operation $\mathcal{M}$ on the time dimension $T$ of SNNs output. $F_{uee}$ is the output of our UEE.

## 3.6 Discussion

After extracting common shared features and unique features from both domains, we fuse them with a concatenate operation. Considering different video sequences have different classes, movement styles, and challenging aspects, we further use three fully connected layers named as $fc_4$, $fc_5$, and $fc_6$ whose output channels are 512, 512, and 2, respectively, to further process fusion features. $fc_6$ is a domain-specific layer, that means each training has $k$ sequences, then there are $k$ $fc_6$ layers. Each of the $k$ sequences contains a binary classification layer with softmax cross-entropy loss, which is responsible for distinguishing target and background.

It should be noted that we did not use a very deep network or complex integration strategy because of the following reasons. First, compared with visual recognition problems, visual tracking requires much lower model complexity because it aims to distinguish only two categories of target and background. Second, since the target is usually small, it is desirable to reduce the input size, which will naturally reduce the depth of the network. Finally, due to the need for online training and testing, a smaller network will be more effective. Our main principle of network design is to make it simple yet work. To the best of our knowledge, this work is the first to explore and utilize the correlation between RGB images and event-based data for visual object tracking. We believe that more and more related works could be done to further improve such a compact network.

## 3.7 Training Details

For CFE, we initialize parameters of it using the pretrained model in VGGNet-M [51]. For UEE, by the public SLAYER [49], we can calculate the gradient of the loss function relative to the SNNs parameter based on the first-order optimization method. We initialize parameters of UEE using the pre-trained model in [17]
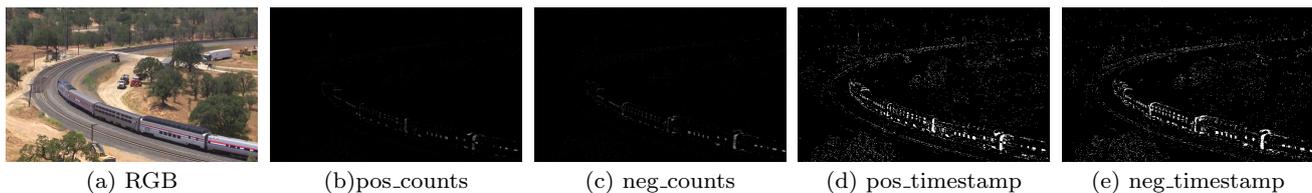
|  (a) RGB  |  (b) pos_counts  |  (c) neg_counts  |  (d) pos_timestamp  |  (e) neg_timestamp  |

**Fig. 3** Example of counts and timestamp images. Left to right: RGB image, positive counts image, negative counts image, positive timestamp image, and negative timestamp image. In timestamp images, each pixel represents the timestamp of the most recent event, and brighter is more recent.

and then fix them. We use the stochastic gradient descent algorithm (SGD) to train the network. The batch size is set to 8 frames which are randomly selected from training video sequences. We choose 32 positive samples (IoU overlap ratios with the ground truth bounding box are larger than 0.7) and 96 negative samples (IoU overlap ratios with the ground truth bounding box are less than 0.5) from each frame, which results in 256 positive and 768 negative samples altogether in a mini-batch. For multi-domain learning with $k$ training sequences, we train the network by softmax cross-entropy loss. The learning rates of all convolutional layers are set to 0.0001, the learning rates of *fc4* and *fc5* are set to 0.0001, and the learning rate of *fc6* is set to 0.001.

### 3.8 Tracker Details

During the tracking process, for each test video sequence, we replace $k$ branches of *fc6* with a single branch. To capture the context of a new sequence and learn video-specific information adaptively, we adopt online fine-tuning. Specifically, we fix all convolutional filters of UEE, CFE, and UER, and fine-tune the fully connected layers *fc4*, *fc5*, and a single branch *fc6*. The reason is that convolutional layers can extract the generic information about tracking, while the fully connected layers are able to learn video-specific information. For online updating, we collect 500 positive samples (IoU overlap ratios with the ground truth bounding box are greater than 0.7) and 5000 negative samples (IoU overlap ratios with the ground truth bounding box are less than 0.5) as the training samples in the first frame. For the $t$-th frame, we collect a set of candidate regions $z_t^i$ from previous tracking result $Z_{t-1}$ by Gaussian sampling. We then use these candidates as inputs to our network and obtain their classification scores. The positive and negative scores are computed using the trained network as $f^+(z_t^i)$ and $f^-(z_t^i)$, respectively. We select the candidate region with the highest score as the target location $Z_t^*$ of the current frame:

$$Z_t^* = \arg\max_{z_t^i} f^+(z_t^i), \quad i = 1, 2, ..., N \tag{7}$$

where $N$ is the number of candidate regions. We use the bounding box regression technique to improve the problem of target scale transformation in the tracking process and improve the accuracy of positioning.

## 4 Experiments

### 4.1 Training Dataset Generation

Supervised learning for visual object tracking requires a large quantity of data. In our case, we need a dataset that contains RGB data from a traditional APS camera (an APS (Active Pixel Sensor) is a conventional image sensor where each pixel sensor unit cell has a photodetector and one or more active transistors) and events from an event-based camera with ground truth bounding box. Our data set needs to meet the following needs: First, the RGB images and event-based data must be aimed at the same scene, and the data between different domains must be aligned. Second, we must have a large variety of scenes with ground truth bounding boxes to avoid overfitting to specific visual patterns. To our knowledge, such data sets do not yet exist. In order to meet the above requirements, we generate a synthetic dataset using event-camera simulator ESIM [45] on large-scale short-term generic object tracking database GOT-10k [21]. ESIM [45] has successfully been proven its effectiveness in previous works [57, 46, 53]. GOT-10k [21] is a large, high-diversity, and one-shot tracking database with a wide coverage of real-world moving objects. GOT-10k [21] collects over 10,000 videos of 563 object classes and annotates 1.5 million tight bounding boxes manually.

Actually, as we all know, traditional RGB frames suffer from motion blur under fast motion, and also have limited dynamic range resulting in the loss of details. Therefore, directly using the RGB and event pairs from ESIM [45] is not an ideal way for training the network, as our goal is to fully exploit the advantages of event cameras. Therefore, we randomly select 100 video sequences. For each RGB frame in the sequence, we randomly increase or decrease the exposure manually. In this way, we simulate the fact that event-based
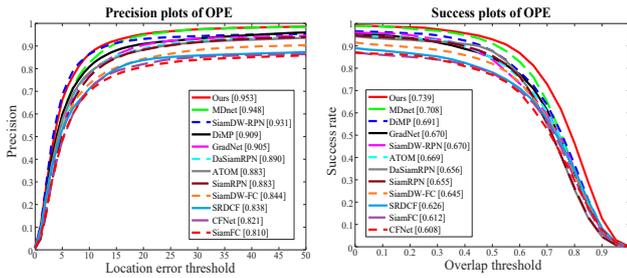
**Fig. 4** PR and SR curves of different tracking result on OTB2013 [59] dataset, where the representative PR and SR scores are presented in the legend.

cameras can provide valuable information that conventional cameras cannot capture in extreme exposure conditions. To verify the effectiveness of our proposed approach, we evaluate it on the standard RGB benchmark and the real event dataset, respectively.

### 4.2 Evaluation on Standard RGB Benchmark

To demonstrate the effectiveness of our MCFR, we first test it on the standard RGB benchmark OTB2013 [59]. The evaluation is based on two metrics: the Precision Rate (PR) and the Success Rate (SR). SR cares the frame of that overlap between ground truth and predicted bounding box is larger than a threshold; PR focuses on the frame of that the center distance between ground truth and predicted bounding box within a given threshold. The one-pass evaluation (OPE) is employed to compare our algorithm with the eleven state-of-the-art trackers including SiamDW-RPN [72], MDNet [38], SiamFC [2], CFNet [56], SiamRPN [29], SiamDW-FC [72], DaSiamRPN [75], SRDCF [14], GradNet [32], DiMP [3], and ATOM [12]. We also apply ESIM [45] to generate event-based data on OTB2013 [59].

The evaluation results are reported in Figure 4. From the results, we can see that our method outperforms the other trackers on OTB2013 [59]. In particular, our MCFR (95.3%/73.9% in PR/SR) outperforms 3.1% over the second-best tracker MDNet [38] in SR, and is superior to other trackers in PR. It demonstrates the effectiveness of our structure for extracting the common features and unique features from different domains. In addition, the remarkable superior performance over the state-of-the-art trackers like ATOM [12] and DiMP [3] suggests that our method is able to make the best use of event domain information to boost tracking performance.

In order to analyze what reliable information the event-based data provides, we report the results on various challenge attributes to show more detailed performance. As shown in Figure 5, our tracker can effec-

tively handle these challenging situations that traditional RGB trackers often lose targets. In particular, under the challenging scenes of fast motion and motion blur, our tracker greatly surpasses the other trackers. That's because the low latency and high temporal resolution of the event-based camera bring more information about the movement between adjacent RGB frames, which can effectively promote the performance of our tracker. From Figure 5, we can also find that our tracker has the best performance in illumination variation scenes. Moreover, in the *background_clutter* (the background near the target has the similar color or texture as the target), as event-based data pays more attention to moving objects rather than the color or texture of objects, our tracker has been significantly improved.

### 4.3 Evaluation on Real Event Dataset

To further prove the effectiveness of our method, we also evaluate it on the real event dataset EED [36]. The EED [36] was recorded using a DAVIS [5] event camera in real-world environments, which contains the events sequences and the corresponding RGB sequences for each video. The EED [36] also provides the ground truth for targets. The EED [36] contains five sequences: *fast_drone*, *light_variations*, *what_is_background*, *occlusions*, and *multiple_objects*. Since *multiple_objects* involves multiple targets, we use the first four video sequences here. Specifically, *fast_drone* describes a fast moving drone under a very low illumination condition, and in *light_variations*, a strobe light flashing at a stable frequency is placed in a dark room. A thrown ball with a dense net as foreground in *what_is_background*, and a thrown ball with a short occlusion under a dark environment in *occlusions*.

Following [8], we use two metrics: the Average Precision (AP) and the Average Robustness (AR) for evaluation. AP and AR describe the accuracy and robustness of the tracker, respectively. The AP can be formulated as follows:

$$AP = \frac{1}{N} \frac{1}{M} \sum_{a=1}^{N} \sum_{b=1}^{M} \frac{O_{a,b}^E \cap O_{a,b}^G}{O_{a,b}^E \cup O_{a,b}^G}, \tag{8}$$

where $N$ is the repeat times of the evaluation (here we set $N$ to 5), and $M$ is the number of objects in the current sequence. $O_{a,b}^E$ is the estimated bounding box in the $a$-th round of the evaluation for the $b$-th object, and $O_{a,b}^G$ is the corresponding ground truth. The AR can be formulated as follows:

$$AR = \frac{1}{N} \frac{1}{M} \sum_{a=1}^{N} \sum_{b=1}^{M} success_{a,b}, \tag{9}$$
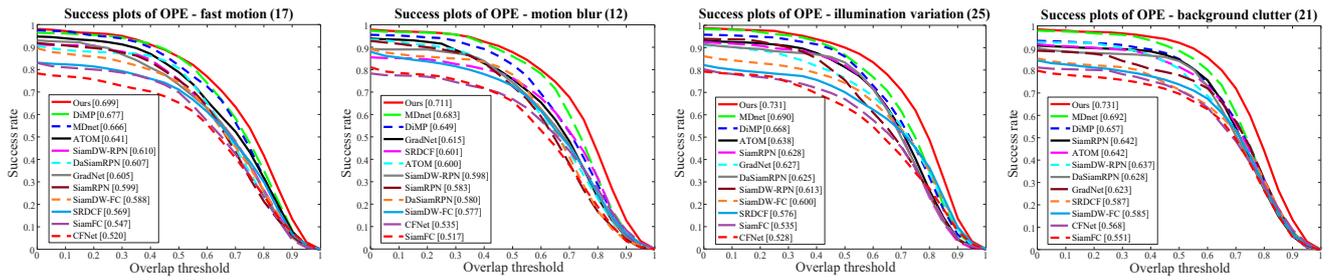
**Fig. 5** Evaluation results on various challenges comparing to the-state-of-the-art methods on OTB2013 [59]. Left to right: *fast_motion*, *motion_blur*, *illumination_variations* and *background_clutter*.



**Fig. 6** Qualitative evaluation of our method and other trackers including CFNet [56],GradNet [32], MDNet [38], SiamDW-RPN [72], SiamDW-FC [72], and SiamFC [2] on 8 challenging videos from OTB2013 [59]. From left to right and top to down are *Ironman*, *CarScale*, *Matrix*, *MotorRolling*, *Skating1*, *Skiing*, *Tiger2*, and *Trellis* respectively. Best viewed in zoom in.

where $success_{a,b}$ indicates that whether the tracking in the $a$-th round for the $b$-th object is successful or not. It will be considered a failure condition if the AP value is less than 0.5. We compare our algorithm with seven state-of-the-art methods including KCL [20], TLD [23], SiamFC [2], ECO [13], DaSiamRPN [75], E-MS [1], and ETD [8]. Herein, the first five algorithms are correlation filter-based or deep learning based traditional RGB object tracking methods, and the remaining are event-based tracking methods. The quantitative results are shown in Table 1, we can see that the traditional RGB trackers are severely affected by low light and fast motion. When there is too much noise in the events, due to lacking image texture information, the event-

based tracker cannot effectively obtain satisfactory performance. Instead, our proposed structure can simultaneously obtain texture information from RGB and target edge cues from events so that our method can effectively handle high dynamic range and fast motion conditions.

### 4.4 Ablation Study

To verify RGB images and event-based data can jointly promote the tracker performance, we implement three variants, including 1) $MCFR_{oE}$, that only applies CFE with events as inputs. 2) $MCFR_{oR}$, that only applies CFE with RGB images as inputs. 3) $MCFR_{ER}$, that

**Table 1** Results obtained by the competitors and our method on the EED [36] dataset. The best results are in red.

| Methods | fast_drone | | light variations | | what is background | | occlusions | |
|---|---|---|---|---|---|---|---|---|
| | AP ↑ | AR ↑ | AP ↑ | AR ↑ | AP ↑ | AR ↑ | AP ↑ | AR ↑ |
| KCL [20] | 0.169 | 0.176 | 0.107 | 0.066 | 0.028 | 0.000 | 0.004 | 0.000 |
| TLD [23] | 0.315 | 0.118 | 0.045 | 0.066 | 0.269 | 0.333 | 0.092 | 0.167 |
| SiamFC [2] | 0.559 | 0.667 | 0.599 | 0.675 | 0.307 | 0.308 | 0.148 | 0.000 |
| ECO [13] | 0.637 | 0.833 | 0.586 | 0.688 | 0.616 | 0.692 | 0.108 | 0.143 |
| DaSiamRPN [75] | 0.673 | 0.853 | 0.654 | 0.894 | 0.678 | 0.833 | 0.189 | 0.333 |
| E-MS [1] | 0.313 | 0.307 | 0.325 | 0.321 | 0.362 | 0.360 | 0.356 | 0.353 |
| ETD [8] | 0.738 | 0.897 | 0.842 | 0.933 | 0.653 | 0.807 | 0.431 | 0.647 |
| MCFR(Ours) | 0.802 | 0.931 | 0.853 | 0.933 | 0.734 | 0.871 | 0.437 | 0.644 |

**Table 2** Ablation analyses of MCFR and its variants.

| | UEE | CFE | UER | RGB | Event | C | T | PR(%) | SR(%) |
|---|---|---|---|---|---|---|---|---|---|
| $MCFR_{oE}$ | ✗ | ✓ | ✗ | ✗ | ✓ | ✓ | ✓ | 0.397 | 0.556 |
| $MCFR_{oR}$ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | 0.702 | 0.944 |
| $MCFR_{ER}$ | ✗ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | 0.719 | 0.949 |
| w/o UEE | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 0.729 | 0.950 |
| w/o CFE | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | 0.710 | 0.947 |
| w/o UER | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | 0.723 | 0.951 |
| $MCFR_C$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | 0.720 | 0.951 |
| $MCFR_T$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | 0.728 | 0.950 |
| MCFR | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 0.739 | 0.953 |



**Fig. 7** Failure cases. The target is stationary and a moving object similar to the target appears around. Red box is GT, green box is our result.

applies CFE with events and RGB data as inputs. The comparison results are shown in Table 2. The results illustrate that the collaborative use of multi-domain information is indeed superior to a single domain.

To validate our method can effectively extract common and unique features from RGB and event domains, we implement three variants based on MCFR, including 1) w/o *UEE*, that removes Unique Extractor for Event, 2) w/o *CFE*, that removes Common Feature Extractor, and 3) w/o *UER*, that removes Unique Extractor for RGB. From Table 2, we can see that our MCFR is superior over w/o *UEE*, which suggests the UEE with SNNs is helpful to take advantage of the event-based data, thereby improving the tracking performance. Besides, MCFR outperforms w/o *CFE* by a clear margin demonstrates that it is essential to extract common features of targets. The superior performance of MCFR over w/o *UER* suggests unique texture features from RGB are important for tracking.

We also explore the performance impact of different ways of stacking events. $MCFR_C$ and $MCFR_T$ represent stacking event streams according to counts and the latest timestamp, respectively. From Table 2, we can see that MCFR outperforms $MCFR_C$ and $MCFR_T$, which verifies that counts images $C$ can record all the events that occurred within a period, and timestamps images $T$ can encode features about the motion.

## 4.5 Failure Cases Analysis

Our method does have limitations. The failure examples are shown in Figure 7. Since the target is static, the event camera cannot effectively provide the edge cues of the target, resulting in the unavailability of information in the event domain. At the same time, an object similar to the target moves around the target, similar colors and textures will interfere with the target-related information provided by the RGB domain. In these cases, the event provides misleading information about moving object, which causes incorrect positioning.

## 5 Conclusion

In this paper, we propose Multi-domain Collaborative Feature Representation (MCFR) to effectively extract and fuse common features and unique features from the RGB and event domain for robust visual object tracking in some challenging conditions, such as fast motion and high dynamic range. Specifically, we apply CFE to extract common features and design UEE based on SNNs and UER based on DCNNs to present specific features of the RGB and event data. Extensive experiments on the RGB tracking benchmark and real event dataset suggest that the proposed tracker achieves outstanding performance. In future work, we will explore upgrading our event-based module so that it can be easily extended to existing RGB trackers for improving performance in challenging conditions.

## References

1. Barranco, F., Fermuller, C., Ros, E.: Real-time clustering and multi-target tracking using event-based sensors. In: IEEE/RSJ International Conference on Intelligent Robots and Systems (2018)

2. Bertinetto, L., Valmadre, J., Henriques, J.F., Vedaldi, A., Torr, P.H.: Fully-convolutional siamese networks for object tracking. In: Proceedings of the European Conference on Computer Vision. Springer (2016)

3. Bhat, G., Danelljan, M., Gool, L.V., Timofte, R.: Learning discriminative model prediction for tracking. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2019)

4. Bi, Y., Chadha, A., Abbas, A., Bourtsoulatze, E., Andreopoulos, Y.: Graph-based object classification for neuromorphic vision sensing. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2019)

5. Brändli, C., Berner, R., Yang, M., Liu, S.C., Delbrück, T.: A 240 × 180 130 db 3 $\mu s$ latency global shutter spatiotemporal vision sensor. IEEE Journal of Solid-state Circuits (2014)

6. Cadena, P.R.G., Qian, Y., Wang, C., Yang, M.: Spade-e2vid: Spatially-adaptive denormalization for event-based video reconstruction. IEEE Transactions on Image Processing (2021)

7. Chen, H., Suter, D., Wu, Q., Wang, H.: End-to-end learning of object motion estimation from retinal events for event-based object tracking. In: Proceedings of the AAAI Conference on Artificial Intelligence (2020)

8. Chen, H., Wu, Q., Liang, Y., Gao, X., Wang, H.: Asynchronous tracking-by-detection on adaptive time surfaces for event-based object tracking. In: Proceedings of the 27th ACM International Conference on Multimedia (2019)

9. Choi, J., Yoon, K.J., et al.: Learning to super resolve intensity images from events. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2020)

10. Cohen, G.K., Orchard, G., Leng, S.H., Tapson, J., Benosman, R.B., Van Schaik, A.: Skimming digits: neuromorphic classification of spike-encoded images. Frontiers in neuroscience (2016)

11. Dai, K., Wang, D., Lu, H., Sun, C., Li, J.: Visual tracking via adaptive spatially-regularized correlation filters. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2019)

12. Danelljan, M., Bhat, G., Khan, F.S., Felsberg, M.: Atom: Accurate tracking by overlap maximization. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2019)

13. Danelljan, M., Bhat, G., Shahbaz Khan, F., Felsberg, M.: Eco: Efficient convolution operators for tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2017)

14. Danelljan, M., Hager, G., Shahbaz Khan, F., Felsberg, M.: Learning spatially regularized correlation filters for visual tracking. In: Proceedings of the IEEE International Conference on Computer Vision (2015)

15. Fan, H., Ling, H.: Siamese cascaded region proposal networks for real-time visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2019)

16. Gehrig, D., Loquercio, A., Derpanis, K.G., Scaramuzza, D.: End-to-end learning of representations for asynchronous event-based data. In: Proceedings of the IEEE International Conference on Computer Vision (2019)

17. Gehrig, M., Shrestha, S.B., Mouritzen, D., Scaramuzza, D.: Event-based angular velocity regression with spiking networks. In: 2020 IEEE International Conference on Robotics and Automation (ICRA) (2020)

18. Gerstner, W.: Time structure of the activity in neural network models. Physical review E (1995)

19. He, A., Luo, C., Tian, X., Zeng, W.: A twofold siamese network for real-time object tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2018)

20. Henriques, J.F., Caseiro, R., Martins, P., Batista, J.: High-speed tracking with kernelized correlation filters. IEEE transactions on pattern analysis and machine intelligence (2014)

21. Huang, L., Zhao, X., Huang, K.: Got-10k: A large high-diversity benchmark for generic object tracking in the wild. IEEE Transactions on Pattern Analysis and Machine Intelligence (2019)

22. Jung, I., Son, J., Baek, M., Han, B.: Real-time mdnet. In: Proceedings of the European Conference on Computer Vision (2018)

23. Kalal, Z., Mikolajczyk, K., Matas, J.: Tracking-learning-detection. IEEE transactions on pattern analysis and machine intelligence (2011)

24. Kart, U., Kämäräinen, J.K., Matas, J., Fan, L., Cricri, F.: Depth masked discriminative correlation filter. In: 2018 24th International Conference on Pattern Recognition (2018)

25. Kart, U., Kämäräinen, J.K., Matas, J., Matas, J.: How to make an rgbd tracker? In: Proceedings of the European Conference on Computer Vision (2018)

26. Kepple, D.R., Lee, D., Prepsius, C., Isler, V., Park, I.M., Lee, D.D.: Jointly learning visual motion and confidence from local patches in event cameras. In: Proceedings of the European Conference on Computer Vision (2020)

27. Lan, X., Ye, M., Zhang, S., Yuen, P.C.: Robust collaborative discriminative learning for rgb-infrared tracking. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018)

28. Li, B., Wu, W., Wang, Q., Zhang, F., Xing, J., Yan, J.: Siamrpn++: Evolution of siamese visual tracking with very deep networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2019)

29. Li, B., Yan, J., Wu, W., Zhu, Z., Hu, X.: High performance visual tracking with siamese region proposal network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2018)

30. Li, C., Lu, A., Zheng, A., Tu, Z., Tang, J.: Multi-adapter rgbt tracking. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (2019)

31. Li, C., Zhu, C., Huang, Y., Tang, J., Wang, L.: Cross-modal ranking with soft consistency and noisy labels for robust rgb-t tracking. In: Proceedings of the European Conference on Computer Vision (2018)

32. Li, P., Chen, B., Ouyang, W., Wang, D., Yang, X., Lu, H.: Gradnet: Gradient-guided network for visual object tracking. In: Proceedings of the IEEE International Conference on Computer Vision (2019)

33. Li, W., Li, X., Bourahla, O.E., Huang, F., Wu, F., Liu, W., Wang, Z., Liu, H.: Progressive multistage learning for discriminative tracking. IEEE Transactions on Cybernetics (2020)

34. Mei, H., Liu, Y., Wei, Z., Zhou, D., Xiaopeng, X., Zhang, Q., Yang, X.: Exploring dense context for salient object detection. IEEE Transactions on Circuits and Systems for Video Technology (2021)

35. Mei, H., Yang, X., Wang, Y., Liu, Y., He, S., Zhang, Q., Wei, X., Lau, R.W.: Don't hit me! glass detection in real-world scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2020)

36. Mitrokhin, A., Fermuller, C., Parameshwara, C., Aloimonos, Y.: Event-based moving object detection and tracking. In: IEEE/RSJ International Conference on Intelligent Robots and Systems (2018)

37. Mostafavi, M., Wang, L., Yoon, K.J.: Learning to reconstruct hdr images from events, with applications to depth and flow prediction. International Journal of Computer Vision (2021)

38. Nam, H., Han, B.: Learning multi-domain convolutional neural networks for visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2016)

39. Neftci, E.O., Mostafa, H., Zenke, F.: Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks. IEEE Signal Processing Magazine (2019)

40. Pan, L., Liu, M., Hartley, R.: Single image optical flow estimation with an event camera. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2020)

41. Piatkowska, E., Belbachir, A.N., Schraml, S., Gelautz, M.: Spatiotemporal multiple persons tracking using dynamic vision sensor. In: IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (2012)

42. Qiao, Y., Liu, Y., Yang, X., Zhou, D., Xu, M., Zhang, Q., Wei, X.: Attention-guided hierarchical structure aggregation for image matting. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2020)

43. Qiao, Y., Liu, Y., Zhu, Q., Yang, X., Wang, Y., Zhang, Q., Wei, X.: Multi-scale information assembly for image matting. In: Computer Graphics Forum (2020)

44. Ramesh, B., Zhang, S., Yang, H., Ussa, A., Ong, M., Orchard, G., Xiang, C.: e-tld: Event-based framework for dynamic object tracking. IEEE Transactions on Circuits and Systems for Video Technology (2020)

45. Rebecq, H., Gehrig, D., Scaramuzza, D.: Esim: an open event camera simulator. In: Conference on Robot Learning (2018)

46. Rebecq, H., Ranftl, R., Koltun, V., Scaramuzza, D.: Events-to-video: Bringing modern computer vision to event cameras. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2019)

47. Ren, W., Wang, X., Tian, J., Tang, Y., Chan, A.B.: Tracking-by-counting: Using network flows on crowd density maps for tracking multiple targets. IEEE Transactions on Image Processing (2020)

48. Shrestha, S.B., Orchard, G.: Slayer: Spike layer error reassignment in time. In: Advances in Neural Information Processing Systems (2018)

49. Shrestha, S.B., Orchard, G.: SLAYER: Spike layer error reassignment in time. In: S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, R. Garnett (eds.) Advances in Neural Information Processing Systems. Curran Associates, Inc. (2018)

50. Shrestha, S.B., Song, Q.: Robustness to training disturbances in spikeprop learning. IEEE transactions on neural networks and learning systems (2017)

51. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)

52. Song, S., Xiao, J.: Tracking revisited using rgbd camera: Unified benchmark and baselines. In: Proceedings of the IEEE International Conference on Computer Vision (2013)

53. Stoffregen, T., Gallego, G., Drummond, T., Kleeman, L., Scaramuzza, D.: Event-based motion segmentation by motion compensation. In: Proceedings of the IEEE International Conference on Computer Vision (2019)

54. Tavanaei, A., Ghodrati, M., Kheradpisheh, S.R., Masquelier, T., Maida, A.: Deep learning in spiking neural networks. Neural Networks (2019)

55. Tulyakov, S., Fleuret, F., Kiefel, M., Gehler, P., Hirsch, M.: Learning an event sequence embedding for dense event-based deep stereo. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2019)

56. Valmadre, J., Bertinetto, L., Henriques, J., Vedaldi, A., Torr, P.H.: End-to-end representation learning for correlation filter based tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2017)

57. Wang, L., Ho, Y.S., Yoon, K.J., et al.: Event-based high dynamic range image and very high frame rate video generation using conditional generative adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2019)

58. Wang, X., Fan, B., Chang, S., Wang, Z., Liu, X., Tao, D., Huang, T.S.: Greedy batch-based minimum-cost flows for tracking multiple objects. IEEE Transactions on Image Processing (2017)

59. Wu, Y., Lim, J., Yang, M.H.: Object tracking benchmark. IEEE Transactions on Pattern Analysis and Machine Intelligence (2015)

60. Xiao, J., Stolkin, R., Gao, Y., Leonardis, A.: Robust fusion of color and depth data for rgb-d target tracking using adaptive range-invariant depth models and spatiotemporal consistency constraints. IEEE transactions on cybernetics (2017)

61. Xu, K., Yang, X., Yin, B., Lau, R.W.: Learning to restore low-light images via decomposition-and-enhancement. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2020)

62. Yang, X., Mei, H., Xu, K., Wei, X., Yin, B., Lau, R.W.: Where is my mirror? In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2019)

63. Yang, X., Mei, H., Zhang, J., Xu, K., Yin, B., Zhang, Q., Wei, X.: Drfn: Deep recurrent fusion network for single-image super-resolution with large factors. IEEE Transactions on Multimedia (2018)

64. Yang, X., Xu, K., Chen, S., He, S., Yin, B.Y., Lau, R.: Active matting. Advances in Neural Information Processing Systems (2018)
65. Yang, X., Xu, K., Song, Y., Zhang, Q., Wei, X., Lau, R.W.: Image correction via deep reciprocating hdr transformation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2018)
66. Zenke, F., Ganguli, S.: Superspike: Supervised learning in multilayer spiking neural networks. Neural computation (2018)
67. Zhang, J., Long, C., Wang, Y., Piao, H., Mei, H., Yang, X., Yin, B.: A two-stage attentive network for single image super-resolution. IEEE Transactions on Circuits and Systems for Video Technology (2021)
68. Zhang, J., Long, C., Wang, Y., Yang, X., Mei, H., Yin, B.: Multi-context and enhanced reconstruction network for single image super resolution. In: 2020 IEEE International Conference on Multimedia and Expo. IEEE (2020)
69. Zhang, L., Danelljan, M., Gonzalez-Garcia, A., van de Weijer, J., Shahbaz Khan, F.: Multi-modal fusion for end-to-end rgb-t tracking. In: Proceedings of the IEEE International Conference on Computer Vision Workshops (2019)
70. Zhang, T., Liu, S., Xu, C., Liu, B., Yang, M.H.: Correlation particle filter for visual tracking. IEEE Transactions on Image Processing (2017)
71. Zhang, T., Xu, C., Yang, M.H.: Learning multi-task correlation particle filters for visual tracking. IEEE transactions on pattern analysis and machine intelligence (2018)
72. Zhang, Z., Peng, H.: Deeper and wider siamese networks for real-time visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2019)
73. Zhu, Q., Triesch, J., Shi, B.E.: An event-by-event approach for velocity estimation and object tracking with an active event camera. IEEE Journal on Emerging and Selected Topics in Circuits and Systems (2020)
74. Zhu, Y., Li, C., Luo, B., Tang, J., Wang, X.: Dense feature aggregation and pruning for rgbt tracking. In: Proceedings of the 27th ACM International Conference on Multimedia (2019)
75. Zhu, Z., Wang, Q., Li, B., Wu, W., Yan, J., Hu, W.: Distractor-aware siamese networks for visual object tracking. In: Proceedings of the European Conference on Computer Vision (2018)