# BTMTrack: Robust RGB-T Tracking via Dual-template Bridging and Temporal-Modal Candidate Elimination

Zhongxuan Zhang[a], Bi Zeng[a,*], Xinyu Ni[a], Yimin Du[a]

[a]*Guangdong University of Technology, No. 100 Waihuan Xi Road, Guangzhou Higher Education Mega Center, Guangzhou, 510006, Guangdong, China*

**Abstract**

RGB-T tracking leverages the complementary strengths of RGB and thermal infrared (TIR) modalities to address challenging scenarios such as low illumination and adverse weather. However, existing methods often fail to effectively integrate temporal information and perform efficient cross-modal interactions, which constrain their adaptability to dynamic targets. In this paper, we propose BTMTrack, a novel framework for RGB-T tracking. The core of our approach lies in the dual-template backbone network and the Temporal-Modal Candidate Elimination (TMCE) strategy. The dual-template backbone effectively integrates temporal information, while the TMCE strategy focuses the model on target-relevant tokens by evaluating temporal and modal correlations, reducing computational overhead and avoiding irrelevant background noise. Building upon this foundation, we propose the Temporal Dual Template Bridging (TDTB) module, which facilitates precise cross-modal fusion through dynamically filtered tokens. This approach further strengthens the interaction between templates and the search region. Extensive experiments conducted on three benchmark datasets demonstrate the effectiveness of BTMTrack. Our method achieves state-of-the-art performance, with a 72.3% precision rate on the LasHeR test set and competitive results on RGBT210 and RGBT234 datasets.

*Keywords:* Object Tracking, RGB-T Tracking, Cross-Modal Interaction

## 1. Introduction

Visual Object Tracking (VOT) is a key task in computer vision that aims to localize a target in consecutive video frames. While many approaches (Gao et al., 2023; Chen et al., 2022; Zeng et al., 2024; Ye et al., 2022) have achieved impressive performance under standard conditions, RGB-based trackers often struggle in complex scenarios such as low illumination, cluttered backgrounds, and adverse weather conditions. These challenges, stemming from the reliance of visible-light imaging on environmental conditions, significantly constrain their robustness in real-world applications. Thermal infrared (TIR) imaging, which captures heat emissions, serves as a complementary modality to RGB

---

data, enabling reliable performance in low-light or obstructed scenarios. By leveraging the strengths of both RGB and TIR modalities, RGB-T tracking has shown great potential in diverse applications such as surveillance (Alldieck et al., 2016) and autonomous driving (Dai et al., 2021).
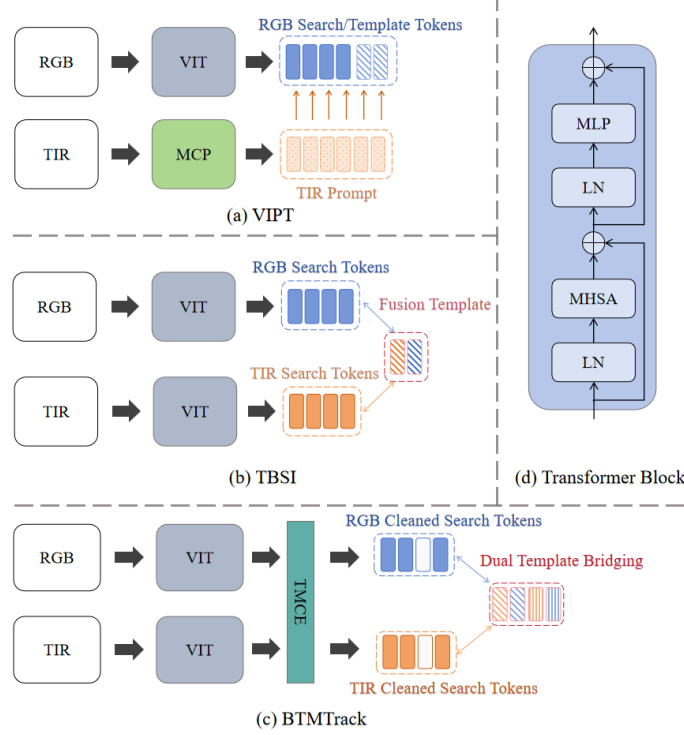


Figure 1: Comparison of our cross-modal fusion approach with previous methods. (a) VIPT, injects TIR modality information as prompt-based auxiliary input into the RGB modality network. (b) TBSI, uses template tokens as a bridge to mediate interactions between the search regions of the two modalities. (c) Our model, filters target-relevant search region tokens before performing dual-temporal template bridging. (d) Example of a Transformer block.

The primary goal of RGB-T tracking is to effectively fuse the complementary information from RGB and TIR modalities, ensuring robust tracking in diverse environments. Recent advancements, including ViPT (Zhu et al., 2023) and TBSI (Hui et al., 2023), have demonstrated notable success by proposing innovative strategies for multi-modal feature fusion. As shown in Fig.1 (a) and (b), ViPT (Zhu et al., 2023) utilizes TIR as prompts for RGB to enhance feature interaction, while TBSI (Hui et al., 2023) employs a Template-Bridged Search region Interaction structure for cross-modal fusion. However, these approaches often overlook the importance of leveraging temporal information to address target variations, and they fail to effectively mitigate the impact of background noise, which reduces the model's discriminative capability. Therefore, two significant challenges still remain. The first challenge lies in handling target variations during tracking, such as changes in aspect ratio, deformation, or fast motion. Without considering temporal information, trackers face difficulties in adapting to these dynamic

changes, resulting in issues such as tracking drift or inaccurate localization. The second challenge is the impact of background noise, which introduces redundancy and misguides the model's judgment. Conventional fusion strategies often fail to effectively suppress irrelevant information, further weakening the tracker's ability to accurately distinguish the target from the background.

To address the challenges of dynamic target variations, some recent works, such as TATrack (Wang et al., 2024), have incorporated temporal information by introducing online template updates. These methods enable the tracker to adapt to target variations by leveraging historical information during tracking process. However, they primarily focus on limited template interactions, failing to fully exploring the synergy between static and dynamic templates. As a result, these approaches often fail to strike a balance between capturing target variations and maintaining stable representation, which can lead to tracking drift in complex scenarios. To tackle the aforementioned challenges, we propose BTMTrack, a novel RGB-T tracking framework that integrates dynamic templates as a core component of its design. Static and dynamic templates interact within a unified framework, where self-attention mechanisms across transformer blocks. This design enables the tracker to effectively balance object appearance and recent changes, thereby enhancing its adaptability to dynamic target variations. By balancing these two perspectives, our framework achieves robust performance in tracking scenarios involving significant target variations. While dynamic templates enhance the tracker's adaptability to target variations, they are insufficient on their own to effectively address the challenge of background noise. This prompts a critical question: can we design an approach that simultaneously leverages temporal information, models cross-modal associations, and mitigates background interference without incurring excessive computational overhead?

Inspired by recent advancements such as BAT (Cao et al., 2024) and USTrack (Xia et al., 2023), we observe that the dominant modality should adapt dynamically to varying scenarios rather than relying exclusively on either RGB or TIR modalities. For instance, RGB images excel in well-lit environments with rich color and texture, while TIR provides higher confidence in low-visibility conditions. To address this, we propose the Temporal-Modal Candidate Elimination (TMCE) strategy, which dynamically prunes noisy tokens in the search region by leveraging both temporal and modal features. By jointly evaluating the contributions of RGB and TIR modalities, TMCE mitigates single-modality reliance and selectively retains the most effective tokens. Furthermore, it incorporates scores from both static and dynamic templates to handle target variations, leveraging temporal and multi-modal features for precise token selection. TMCE reduces the computational cost of dynamic templates while suppressing background noise, thereby enhancing robustness in complex environments. To complement TMCE, we introduce the Temporal Dual Template Bridging (TDTB) module, which builds upon and extends the template-bridged search region interaction inspired by TBSI (Hui et al., 2023). Through bidirectional attention mechanisms, TDTB enables static and dynamic templates to collaboratively interact with the search region, deeply fusing RGB and TIR information while enriching templates with contextual relevance. Together, TMCE and TDTB collaboratively adapt to temporal and modal variations, ensuring robust and efficient performance in RGB-T tracking. The primary contributions of this work are summarized as follows:

- We propose BTMTrack, an innovative RGB-T tracking framework that leverages

both static and dynamic templates to enable adaptive temporal modeling and robust object tracking.

- We develop the TMCE strategy, which conducts temporal and modality-aware token pruning, effectively suppressing background noise while ensuring computational efficiency.

- We propose the TDTB module, which facilitates comprehensive cross-modal interactions between templates and search regions, thereby improving feature fusion and enhancing overall robustness.

- Our tracker achieves state-of-the-art performance across three benchmarks, especially attaining a precision of 72.3% on the LasHeR test set, while also delivering competitive results on RGBT210 and RGBT234.

## 2. Related Works

### 2.1. RGB-T Tracking

RGB-based trackers often struggle in challenging scenarios such as low illumination, thermal crossover, and adverse weather due to their reliance on visible light. By leveraging the complementary strengths of visible (RGB) and thermal infrared (TIR) modalities, RGB-T tracking achieves robust performance across diverse environments. mfDiMP (Zhang et al., 2019a) incorporates multimodal fusion at multiple stages to align and combine features from RGB and TIR modalities. APFNet (Xiao et al., 2022) adopts an attribute-based progressive fusion network, aggregating features under specific conditions such as occlusion or illumination variation. SiamCDA (Zhang et al., 2021) reduces modality differences before fusion, enhancing the quality of the fused features. While effective, these methods often struggle with background noise and irrelevant regions during fusion, limiting the complementary benefits of RGB and TIR modalities in complex scenarios.

In contrast to traditional approaches, Transformer architectures have significantly advanced RGB-T tracking by leveraging self-attention mechanisms for cross-modal interaction. DRGCNet (Mei et al., 2023) and MIRNet (Hou et al., 2022) use cross-attention to transfer discriminative features across modalities while filtering redundant information through gating mechanisms. TBSI (Hui et al., 2023) introduces a stacked extraction-fusion module to improve interaction between RGB and TIR features, whereas ViPT (Zhu et al., 2023) employs lightweight fusion modules to simplify the fusion process and reduce computational costs. However, these methods typically process all tokens in the search region indiscriminately, failing to prioritize target-relevant features, which introduces noise and inefficiencies. In this paper, we propose the Temporal Dual Template Bridging (TDTB) module, which selectively processes tokens that are dynamically pruned through a temporal-modal evaluation process. Unlike existing methods, TDTB focuses on the most target-relevant tokens, jointly determined by RGB and TIR modalities. By unifying static and dynamic templates with these filtered tokens, TDTB facilitates precise and efficient cross-modal interaction, enabling the model to focus on critical features from both modalities. This design enhances the model's ability to capture complementary information between modalities while dynamically adapting to temporal variations, ensuring robust and efficient RGB-T tracking across diverse scenarios.

## 2.2. Temporal Information Exploitation

Temporal information is critical for robust visual tracking, enabling trackers to adapt to target variations over time. Existing works can be categorized into three main approaches: template update methods, temporal propagation strategies, and multi-modal temporal modeling. Template update methods aim to maintain dynamic templates that adapt to target appearance changes. For example, UpdateNet (Zhang et al., 2019b) uses additional networks to optimize template updates, STARK (Yan et al., 2021) introduces a dual-template design with scoring-based dynamic updates, and MixFormer (Cui et al., 2022) selects the highest-scoring frame as its dynamic template. However, these methods often treat temporal updates as isolated processes, neglecting the interplay between temporal and spatial information. Temporal propagation strategies emphasize utilizing historical information to guide future predictions. ARTrack (Wei et al., 2023) employs autoregressive frameworks to propagate motion dynamics but suffers from high complexity and computational overhead. TCTrack (Cao et al., 2022) leverages online temporally adaptive convolution to refine spatial features with historical data but remains limited to single-modality tracking, missing cross-modal complementarities.

Multi-modal temporal modeling integrates temporal and cross-modal information. TATrack (Wang et al., 2024) introduces an online template for dual-branch tracking, leveraging temporal and modal complementarity. However, it primarily uses TIR to enhance RGB tracking without fully involving TIR in the tracking process. Additionally, its complex spatio-temporal interaction modules introduce significant computational costs. In this paper, we propose the Temporal-Modal Candidate Elimination (TMCE) strategy. TMCE dynamically prunes noisy tokens in the search region by synchronously evaluating temporal, cross-modal, and spatial features. By filtering background noise, it enhances target focus while leveraging temporal and cross-modal correlations. This approach ensures robust and efficient RGB-T tracking across diverse scenarios.

## 3. Method

In this section, we introduce BTMTrack, an RGB-T tracking method that fully leverages temporal and modal information to address background noise challenges. The overall framework of BTMTrack is illustrated in Fig.2

### 3.1. RGB-T Tracking Process Based on ViT

Similar to other state-of-the-art RGB-T tracking methods like TBSI (Hui et al., 2023) and BAT (Cao et al., 2024), we follow the ViT-based framework introduced in VOT (Ye et al., 2022; Cui et al., 2022; Lin et al., 2022) and use the target states and consecutive video frames from both modalities as inputs. Before feeding the data into the model, we also preprocess the dynamic templates specifically based on the current tracking state and input them into the framework accordingly.

Before the tracking process begins, the RGB-T tracker first extracts the appropriate templates $z_{\text{static}}^{\text{rgb}}, z_{\text{static}}^{\text{tir}} \in \mathbb{R}^{H_z \times W_z \times 3}$ from the target state of the first frame and simultaneously initializes the dynamic template $z_{\text{dynamic}}^{\text{rgb}}, z_{\text{dynamic}}^{\text{tir}}$, which is identical to the static template. These templates and their corresponding search regions, denoted
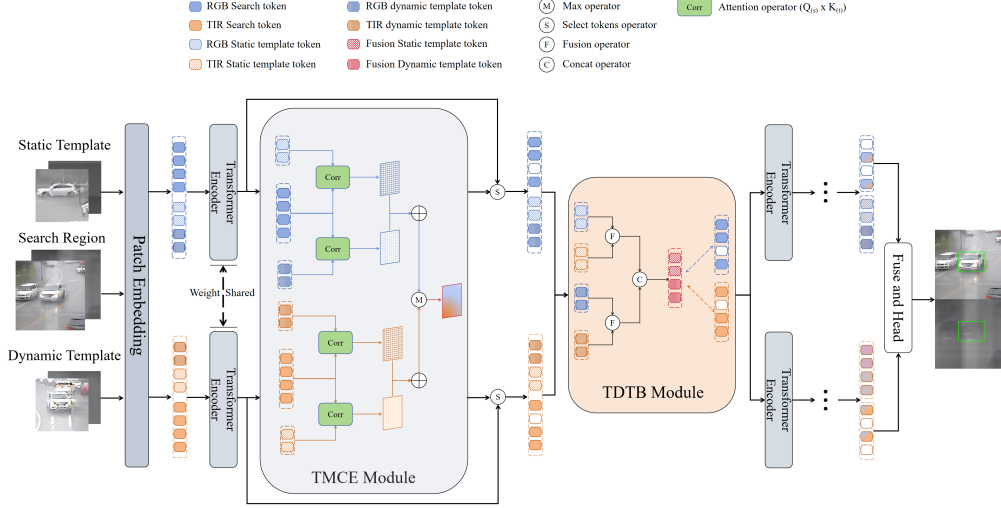
Figure 2: The overall framework of our method. It integrates static and dynamic templates with search regions from RGB and TIR patches. These patches are tokenized and processed by a ViT backbone for feature extraction. The proposed TMCE strategy filters tokens based on temporal and modal relevance to reduce background noise. The TDTB module enables interactions between dual-temporal templates and search regions of both modalities. Finally, fused RGB and TIR features are passed to the tracking head to predict the target's location.

as $x^{\mathrm{rgb}}, x^{\mathrm{tir}} \in \mathbb{R}^{H_x \times W_x \times 3}$, are then passed through the same Patch Embedding layer for tokenization as:

$$x^{\mathrm{rgb}}, x^{\mathrm{tir}} \in \mathbb{R}^{N_x \times (3P^2)}$$
$$z_{\mathrm{static}}^{\mathrm{rgb}}, z_{\mathrm{static}}^{\mathrm{tir}}, z_{\mathrm{dynamic}}^{\mathrm{rgb}}, z_{\mathrm{dynamic}}^{\mathrm{tir}} \in \mathbb{R}^{N_z \times (3P^2)} \quad (1)$$

Here, $P \times P$ represents the resolution of each patch, and $N_z$ and $N_x$ denote the number of patches in the template and search images, respectively. Subsequently, they are enhanced with positional encodings, and the tokens from both the RGB and TIR modalities are then concatenated as:

$$I^{\mathrm{rgb}} = \left( z_{\mathrm{static}}^{\mathrm{rgb}}, z_{\mathrm{dynamic}}^{\mathrm{rgb}}, x^{\mathrm{rgb}} \right), I^{\mathrm{tir}} = \left( z_{\mathrm{static}}^{\mathrm{tir}}, z_{\mathrm{dynamic}}^{\mathrm{tir}}, x^{\mathrm{tir}} \right) \quad (2)$$

These are then input into a series of transformer blocks, where multi-modal joint feature extraction and selective feature fusion take place.

After feature extraction through the backbone network, the token information from both modalities is converted into feature maps and undergoes a simple fusion through convolutional layers. The fused features are then fed into the prediction head to obtain the final output bounding box. The process is as follows:

$$B = H \left( \mathrm{Conv} \left( F \left( I_{\mathrm{rgb}}, I_{\mathrm{tir}} \right) \right) \right) \quad (3)$$

Where B represents the final predicted bounding box, H denotes the prediction head, Conv refers to the 3×3 convolutional layer, and F represents our ViT backbone network.

*3.2. Temporal-Modal Candidate Elimination*

Transformer-based object tracking models possess strong global modeling capabilities, capturing long-range dependencies between the target and various regions in the scene, enabling accurate target localization even in complex scenarios. However, retaining all tokens is not necessary. Since the contributions of different tokens to target localization are uneven, a large number of tokens originating from the background or irrelevant regions often lack useful information, introduce noise and redundancy, and ultimately degrade model performance while increasing computational cost. Consequently, selecting key tokens and filtering irrelevant information is an effective strategy for improving both model efficiency and accuracy. Addressing the problem of selective token interaction, previous work (Ye et al., 2022) proposed the CE mechanism, which leverages the correlation map from the attention computation to score and filter search region tokens, retaining only the highest-scoring tokens for interaction with template tokens. However, the CE method is limited to single-modality information and does not account for temporal dependencies.

To address this issue, we propose TMCE (as illustrated in Fig.2), which dynamically filters effective target candidate tokens based on temporal information and the synchronized evaluation of RGB and TIR modality features, thereby enhancing tracking accuracy and efficiency.

**Balanced Contribution of Dual Templates.** The static template provides global and stable features of the target, helping to anchor its original appearance, while the dynamic template captures temporal variations and adapts to changes in complex scenarios. Therefore, we believe that the static and dynamic templates contribute equally to feature recognition, and we directly add the two correlation maps, which are individually generated by the static template and the dynamic template with the search region, respectively. The process is formally defined as follows:

$$Corr\_Map_{\text{static}} = \text{Softmax}\left(Q_{\text{static template}} K_{\text{search}}^{\top}\right)$$

$$Corr\_Map_{\text{dynamic}} = \text{Softmax}\left(Q_{\text{dynamic template}} K_{\text{search}}^{\top}\right)$$

$$Corr\_Map = (Corr\_Map_{\text{static}} + Corr\_Map_{\text{dynamic}}) \tag{4}$$

**Spatial-level Adaptive Modality Selection.** Inspired by outstanding RGB-T tracking works such as BAT (Cao et al., 2024) and USTrack (Xia et al., 2023), we observe that the confidence levels of RGB and TIR modalities vary significantly across different environments. Specifically, in well-lit and clear weather conditions, the RGB modality typically provides more detailed appearance features, such as color, texture, and target boundaries. In contrast, under nighttime, foggy, or high-glare conditions, the confidence of the RGB modality decreases, while the TIR modality becomes more reliable by capturing thermal radiation characteristics of the target.

Unlike other works that emphasize dominant-subordinate modality selection or modality weight allocation, we adopt a spatial perspective to select the modality with the higher score at each specific patch. We believe that this approach retains the most valuable spatial information from the perspectives of both modalities, enabling more accurate classification and localization tasks. By comparing the corresponding positions on the correlation maps of RGB and TIR ($Corr\_Map_{\text{RGB}}, Corr\_Map_{\text{TIR}}$), which are respectively derived through Eq.4, we obtain a new correlation map. The process is formally defined as follows:

$$Corr\_Map_{\text{Final}} = \max\left(Corr\_Map_{\text{RGB}} + Corr\_Map_{\text{TIR}}\right) \tag{5}$$

Subsequently, we sort all search tokens by their scores on $Corr\_Map_{\text{Final}}$, retaining the top k tokens with the highest scores. Our TMCE strategy effectively integrates temporal, spatial, and cross-modality information, fully exploiting the strengths of the dual-template baseline. It enhances the model's robustness while significantly reducing computational complexity.

### 3.3. Temporal Dual Template Bridging Module

To complement the effects of TMCE, we propose the Temporal Dual Template Bridging (TDTB) module, which further optimizes the interaction between the static and dynamic templates. This module is inserted after the first feature selection to ensure effective information fusion on the basis of retaining the most representative features of both modalities. Inspired by the TBSI (Hui et al., 2023), we adopt a similar template bridging search region interaction method for modality fusion, while effectively connecting the static and dynamic templates along the temporal dimension. By fully utilizing temporal information, TDTB helps the model better adapt to target variations and enhances its performance in dynamic scenarios. The data flow of this module is shown in Fig.3.

**Temporal Template Fusion.** To facilitate dynamic multimodal interaction across different temporal dimensions, we separately fuse the static and dynamic templates of the RGB and TIR modalities, resulting in two multimodal bridging intermediates: $Z_{static} \in \mathbb{R}^{N_z \times C}$ and $Z_{dynamic} \in \mathbb{R}^{N_z \times C}$, which capture temporal and modality-specific features. The formulas for this process are as follows:

$$Z_{static} = [Z_{static}^{rgb}; Z_{static}^{tir}]W_m$$

$$Z_{dynamic} = [Z_{dynamic}^{rgb}; Z_{dynamic}^{tir}]W_m \tag{6}$$

Where $W_m \in \mathbb{R}^{2C \times C}$ is the parameter of the linear layer used in the fusion operator. Subsequently, $Z_{static}$ and $Z_{dynamic}$ influence the fusion process between the two modalities from different temporal perspectives.

**Dual Template-Bridged Interaction with Search Regions.** To enable effective cross-modal interaction, we employ multi-head cross-attention between the template intermediates and their corresponding search regions. This process ensures that information from the static and dynamic templates is seamlessly integrated with the search regions of both RGB and TIR modalities. In the following equations, we will define multi-head cross-attention as $\text{MHCA}(Q, K, V)$, where Q, K, and V represent the query, key, and value matrices, respectively. The details of the Transformer block are shown in Fig.1(d). We define the MHCA operation as follows, omitting the multi-head operation for clarity:

$$\text{MHCA}(Q, K, V) = \text{Softmax}\left(\frac{(QW_q)(KW_k)^T}{\sqrt{C}}\right)(VW_v) \tag{7}$$

Where $W_q$, $W_k$, and $W_v$ are learnable parameter matrices that project the input features
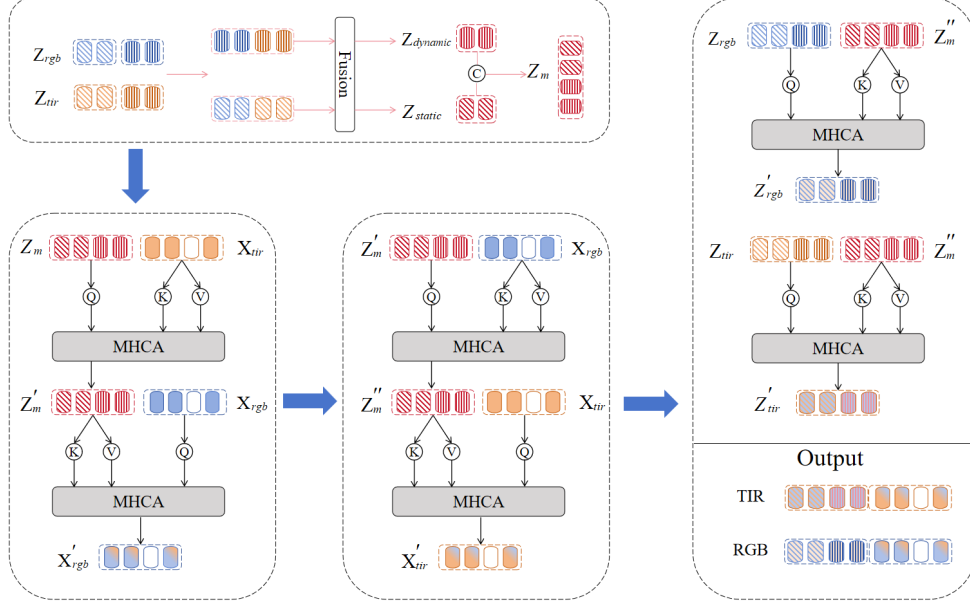
Figure 3: The diagram demonstrates the process of dual-temporal template fusion and six MHCA operations within the TDTB module. For clearer presentation, we omit details such as LN, MLP, and the residual connections typically performed in each Transformer block.

into the query, key, and value spaces, respectively. In the cross-attention mechanism of this module, K is equal to V, so the equation can be simplified as follows:

$$\text{MHCA}(X, Y) = \text{Softmax}\left(\frac{(XW_q)(YW_k)^T}{\sqrt{C}}\right)(YW_v) \qquad (8)$$

First, we concatenate the $Z_{static}$ and $Z_{dynamic}$ directly to obtain $Z_m = [Z_{static}; Z_{dynamic}]$, which is then used in the subsequent interactions. Then, we update the template $Z_m$ to a new template $Z_m^{'}$ using the target-relevant contextual information from the TIR search region $X_{tir}$. The process is as follows:

$$Z_{temp}^{'} = \text{LN}(Z_m + \text{MHCA}(Z_m, X_{tir}))$$

$$Z_m^{'} = \text{LN}(Z_{temp}^{'} + \text{MLP}(Z_{temp}^{'})) \qquad (9)$$

After obtaining the target-relevant contextual information from the TIR search region $X_{tir}$, the updated temporal dual-template $Z_m^{'}$ will propagate the information to the RGB search region $X_{rgb}$ via MHCA, LN, and MLP. The detailed process is as follows:

$$X_{rgb\_temp}^{'} = \text{LN}(X_{rgb} + \text{MHCA}(X_{rgb}, Z_m^{'})$$

$$X_{rgb}^{'} = \text{LN}(X_{rgb\_temp}^{'} + \text{MLP}(X_{rgb\_temp}^{'})) \qquad (10)$$

Similarly, we use the updated temporal dual-template $Z_m^{'}$ as the new dual-template input, and by reversing the interaction order between the search region modalities (from

TIR → RGB to RGB → TIR), we perform the same operations to obtain a new fused template $Z_m^{''}$, which has interacted with both search regions, as well as the updated TIR search region tokens $X_{tir}^{'}$ that have interacted with the RGB search region's contextual information, as illustrated in the flow of Fig.3.

**Original Templates Interaction with Medium.** Finally, we leverage the fused template $Z_m^{''}$, enriched with cross-modal information, to update the original RGB and TIR templates($Z_{rgb} = [Z_{static}^{rgb}; Z_{dynamic}^{rgb}]$ and $Z_{tir} = [Z_{static}^{tir}; Z_{dynamic}^{tir}]$). The core of this process lies in utilizing the comprehensive information provided by Z to fully explore the potential correlations between the RGB and TIR modalities, thereby enhancing the adaptability and accuracy of the templates in object tracking. By effectively extracting and interacting the multimodal information within the fused template, we can dynamically adjust each modality's template, further improving the model's robustness in complex environments. The process for obtaining the updated RGB template $Z_{rgb}^{'}$ is as follows:

$$Z_{rgb\_temp}^{'} = \text{LN}(Z_{rgb} + \text{MHCA}(Z_{rgb}, Z_m^{''})$$
$$Z_{rgb}^{'} = \text{LN}(Z_{rgb\_temp}^{'} + \text{MLP}(Z_{rgb\_temp}^{'})) \tag{11}$$

Similarly, we can obtain the updated TIR template $Z_{tir}^{'}$, which will then be combined with $X_{tir}^{'}$, $Z_{rgb}^{'}$, and $X_{rgb}^{'}$ to serve as the output of the TDTB module.

### 3.4. Prediction Head and Loss Function

In this work, we utilize a unified prediction head design, consistent with prior methods (Ye et al., 2022). The prediction head transforms token sequences into a 2D spatial feature map via fully convolutional networks (FCN). This process outputs three components: the target classification score map (indicating the target location), the offset, and the normalized bounding box dimensions.

The total loss function is defined as:

$$L_{\text{total}} = L_{\text{cls}} + \lambda_1 L_{\text{iou}} + \lambda_2 L_1 \tag{12}$$

where $L_{\text{cls}}$ represents the classification loss computed with weighted Focal Loss (Ross and Dollár, 2017), and $L_{\text{iou}}$ and $L_1$ are the bounding box regression losses, optimized using Generalized IoU Loss (Rezatofighi et al., 2019) and L1 Loss, respectively. The trade-off parameters $\lambda_1$ and $\lambda_2$ balance the contribution of these components.

## 4. Experiments

### 4.1. Implementation Details

Our model is implemented using PyTorch (Paszke et al., 2019), with all experiments conducted on two RTX 3090 GPUs. Additionally, speed tests were performed on a single RTX 2080Ti GPU.

**Architectures**. For the backbone network, we utilize the pre-trained Vision Transformer (ViT) model, which has been trained on the SOT task, as adopted in several advanced approaches (Hui et al., 2023; Zhu et al., 2023; Cao et al., 2024) in the RGB-T tracking domain. The TMCE strategy is applied after the 4th, 7th, and 10th Transformer

| Method | Source | Backbone | LasHeR | | | RGBT234 | | RGBT210 | |
|---|---|---|---|---|---|---|---|---|---|
| | | | SR | PR | NPR | MSR | MPR | MSR | MPR |
| MANet++ (Lu et al., 2021) | TIP 2021 | VGG-M | 31.7 | 46.7 | 40.8 | 55.4 | 80.0 | 55.3 | 78.5 |
| DMCNet (Lu et al., 2022) | TNNLS 2022 | VGG-M | 35.5 | 49.0 | 43.1 | 59.3 | 83.9 | 55.9 | 79.7 |
| APFNet (Xiao et al., 2022) | AAAI 2022 | VGG-M | 36.2 | 50.0 | - | 57.9 | 82.7 | 57.1 | 82.1 |
| CAT++ (Liu et al., 2024) | TIP 2024 | VGG-M | 35.6 | 50.9 | 44.4 | 59.2 | 84.0 | 56.1 | 82.2 |
| HMFT (Zhang et al., 2022a) | CVPR 2022 | ResNet-50 | - | - | - | 56.8 | 78.8 | 53.5 | 78.6 |
| CMD (Zhang et al., 2023) | CVPR 2023 | ResNet-18 | 46.4 | 59.0 | 54.6 | 58.4 | 82.4 | 59.3 | 83.4 |
| MFNet (Zhang et al., 2022b) | IVC 2022 | ResNet-50 | 46.7 | 59.7 | 55.4 | 60.1 | 84.4 | - | - |
| mfDiMP (Zhang et al., 2019a) | CVPRW 2019 | ResNet-50 | 46.7 | 59.9 | - | 59.1 | 84.2 | 59.3 | 84.9 |
| VIPT (Zhu et al., 2023) | CVPR 2023 | ViT-Base | 52.5 | 65.1 | 61.7 | 61.7 | 83.5 | 60.3 | 82.1 |
| QueryTrack (Fan et al., 2024) | TIP 2024 | ViT-Base | 52.0 | 66.0 | - | 60.0 | 84.1 | - | - |
| OneTracker (Hong et al., 2024) | CVPR 2024 | ViT-Base | 53.8 | 67.2 | - | 64.2 | 85.7 | - | - |
| TBSI (Hui et al., 2023) | CVPR 2023 | ViT-Base | 56.5 | 69.2 | 66.5 | 63.7 | 87.1 | 62.5 | 85.3 |
| BAT (Cao et al., 2024) | AAAI 2024 | ViT-Base | 56.3 | 70.2 | 66.4 | 64.1 | 86.8 | - | - |
| TATrack (Wang et al., 2024) | AAAI 2024 | ViT-Base | 56.1 | 70.2 | 66.7 | 64.4 | 87.2 | 61.8 | 85.3 |
| GMMT (Tang et al., 2024) | AAAI 2024 | ViT-Base | 56.6 | 70.7 | 67.0 | 64.7 | 87.9 | - | - |
| BTMTrack | Ours | ViT-Base | 58.0 | 72.3 | 68.5 | 65.4 | 88.3 | 63.7 | 87.5 |

Table 1: Comparison of our model with other state-of-the-art methods on the LasHeR, RGBT234, and RGBT210 benchmarks. The top three results are highlighted in red, blue, and green, respectively.

blocks, while the TDTB module is inserted only after the first application of the TMCE strategy. The resolution of the template is set to $128 \times 128$ pixels, while the search region has a resolution of $256 \times 256$ pixels.

**Training**. Our model is trained on the LasHeR (Li et al., 2021) dataset, and evaluated on the LasHeR (Li et al., 2021), RGBT210 (Li et al., 2017), and RGBT234 (Li et al., 2019) datasets. We set the training to 15 epochs, with 60,000 image pairs per epoch. Each GPU handles 32 image pairs, resulting in a total batch size of 64. The learning rate for the ViT backbone is set to 1e-5, while the learning rate for the other parameters is set to 1e-4. After 10 epochs, the learning rate is decayed by a factor of 10.

**Inference**. During inference, the initial static template, dynamic template, and search region are used as inputs to the model. We determine whether to update the dynamic template based on the target classification score predicted by the prediction head.

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

To evaluate the performance and robustness of our model, we compare it with recent state-of-the-art RGB-T tracking models on three different RGB-T tracking benchmarks.

**LasHeR**. LasHeR (Li et al., 2021) is a large-scale and diverse RGB-T tracking dataset, consisting of 1224 RGB and thermal infrared video sequences, with a total of over 730K frames. The dataset's test set includes 245 challenging video sequences, which capture RGB and thermal infrared image pairs under various scenes and conditions. The goal is to evaluate the performance of RGB-T tracking models in complex environments. We evaluate our model using three metrics: Success Rate(SR), Precision Rate(PR), and Normalized Precision Rate(NPR), and compare it with 15 other advanced trackers on the LasHeR dataset. The RGB-T trackers included in the comparison are: TATrack (Wang et al., 2024), BAT (Cao et al., 2024), TBSI (Hui et al., 2023), GMMT (Tang et al., 2024),

OneTracker (Hong et al., 2024), QueryTrack (Fan et al., 2024), VIPT (Zhu et al., 2023), mfDiMP (Zhang et al., 2019a), MFNet (Zhang et al., 2022b), CMD (Zhang et al., 2023), CAT++ (Liu et al., 2024), APFNet (Xiao et al., 2022), DMCNet (Lu et al., 2022), and MANet++ (Lu et al., 2021) . The results, as shown in Tab.1, indicate that our model outperforms previous state-of-the-art methods on this dataset. Specifically, BTMTrack surpasses GMMT (Tang et al., 2024) by 1.4% in SR and 1.6% in PR. The effectiveness of our model on the large-scale RGB-T dataset LasHeR (Li et al., 2021) clearly demonstrates that combining spatiotemporal information while avoiding irrelevant background noise can significantly improve performance.

**RGBT210**. RGBT210 (Li et al., 2017) is a classic RGB-T tracking benchmark consisting of 210 video sequences, with a total of approximately 210K frames. Each video pair in the dataset contains up to 8K frames, offering long-term tracking challenges for RGB-T trackers. As shown in Tab.1, BTMTrack outperforms TBSI (Hui et al., 2023) by 1.2% and 2.2% in terms of Maximum Success Rate (MSR) and Maximum Precision Rate (MPR), respectively.

**RGBT234**. RGBT234 (Li et al., 2019) is an extension of the RGBT210 (Li et al., 2017) dataset, featuring additional environmental challenges. It consists of 234 video sequences, totaling approximately 234K frames. The dataset includes 12 attributes, such as Low Illumination (LI), Occlusion, Deformation (DEF), and Movement, providing a more diverse benchmark for evaluating RGB-T tracking models. Ground truth labels are provided for both RGB and thermal infrared (TIR) modalities, enabling performance evaluations across multiple modalities. As shown in Tab.1, BTMTrack achieves the best performance on the RGBT234 (Li et al., 2019) dataset compared to 15 other advanced RGB-T trackers. It surpasses GMMT (Tang et al., 2024) by 0.7% in MSR and by 0.4% in MPR, further validating the effectiveness of our model.

### 4.3. Ablation Studies

**Component analysis**. In Tab.2, we conducted ablation experiments on the individual components of our model using the LasHeR (Li et al., 2021) dataset, all under the same parameter settings. Additionally, to validate the effectiveness of the TDTB module, we included the TBSI module as a component for comparison in this ablation study.

We denote **#1** as the baseline, which consists of the ViT backbone pre-trained on the SOT task dataset without any additional components. **#2** adds the TBSI (Hui et al., 2023) module to the baseline. **#3** further introduces a dynamic template enriched with temporal information into the backbone network, which is then concatenated and fed into the TBSI (Hui et al., 2023) module. **#4** replaces the TBSI (Hui et al., 2023) module in **#3** with our TDTB module. The results from **#1** to **#4** demonstrate that the introduction of dual templates can improve the model's performance to some extent; however, it also introduces a larger resource overhead, which slows down the model. Moreover, the TBSI (Hui et al., 2023) module cannot fully leverage the temporal information within the dual templates. In contrast, by replacing it with our more lightweight TDTB module, which is inserted only after the fourth layer of the backbone network and more effectively utilizes temporal information, the model's performance is significantly improved.

We denote **#5** as the backbone network that applies the TMCE strategy with the dual-template design, but without the fusion module. **#6** represents the case where the TBSI (Hui et al., 2023) module is added to **#5**, while **#7** represents our complete

12

| | Components | | | | LasHeR | | | |
|---|---|---|---|---|---|---|---|---|
| | TBSI | Dual-Template | TMCE | TDTB | SR | PR | NPR | FPS (2080Ti) |
| #1 | - | - | - | - | 53.2 | 65.9 | 62.6 | 47.04 |
| #2 | ✓ | - | - | - | 55.6 | 69.2 | 65.8 | 33.03 |
| #3 | ✓ | ✓ | - | - | 55.9 | 69.8 | 66.1 | 29.71 |
| #4 | - | ✓ | - | ✓ | 57.4 | 71.5 | 67.8 | 36.82 |
| #5 | - | ✓ | ✓ | - | 56.1 | 70.4 | 66.6 | **47.48** |
| #6 | ✓ | ✓ | ✓ | - | 57.0 | 71.0 | 67.3 | 32.22 |
| #7 | - | ✓ | ✓ | ✓ | **58.0** | **72.3** | **68.5** | 40.65 |

Table 2: Effectiveness comparison of the proposed BTMTrack components and the TBSI (Hui et al., 2023) module on the LasHeR dataset. Speed tests are conducted on a single RTX 2080Ti GPU. The best results are highlighted in **bold**.

BTMTrack, where the TDTB module is incorporated into **#5**. As shown by the metrics from **#5** to **#7**, introducing the TMCE strategy into the dual-template backbone significantly improves model speed while also enhancing performance to some extent. Adding the TBSI (Hui et al., 2023) module further improves performance. However, the TDTB module, which is more synergistic with TMCE, achieves the best results. Compared to the baseline network, our complete BTMTrack improves the Success Rate (SR) by 4.8% and the Precision Rate (PR) by 6.4%. Additionally, the introduction of the TMCE strategy reduces computational cost, compensating for the overhead introduced by the dual-template, thus striking an excellent balance between speed and performance.

**Comparison of Elimination Strategies**. To evaluate the effectiveness of our proposed TMCE strategy, we conducted experiments on the LasHeR (Li et al., 2021) dataset under the same conditions, testing the changes in model performance when using different Elimination Strategies. *w/o any Strategy* refers to removing the TMCE strategy entirely from BTMTrack, while *with CE* replaces the TMCE strategy with the traditional CE (Ye et al., 2022) strategy. Furthermore, to verify the contribution allocation discussed in Sec.3.2, which incorporates both multi-modal and temporal information into the Elimination Strategy, we designed two variants of the CE strategy as baselines for the ablation study: *with Add-CE* directly adds the correlation maps between the templates and the search region from the two modalities and uses the resulting score to filter tokens; *with Max-CE* directly takes the maximum value of the correlation maps between the templates and the search region from the two modalities and uses the resulting score for token filtering. *with TMCE* represents the full BTMTrack model equipped with the TMCE strategy.

As shown in Tab.3, although the CE (Ye et al., 2022) strategy does not improve performance, it slightly enhances the model's speed. On the other hand, while the Add-CE and Max-CE strategies are faster, their incomplete filtering mechanisms introduce more background noise, leading to degraded performance. These two variants fail to fully utilize the dual-template information within the backbone network, which results in reduced model effectiveness. In contrast, the TMCE strategy significantly improves model performance while also enhancing speed to some extent, further demonstrating

| Elimination Strategy | SR | PR | NPR | FPS (2080Ti) |
|---|---|---|---|---|
| w/o any Strategy | 57.4 | 71.5 | 67.8 | 36.82 |
| with CE | 57.3 | 71.5 | 67.9 | 40.58 |
| with Add-CE | 56.7 | 70.6 | 67.0 | 41.65 |
| with Max-CE | 56.7 | 70.7 | 67.1 | **42.58** |
| with TMCE | **58.0** | **72.3** | **68.5** | 40.65 |

Table 3: Performance comparison of the proposed TMCE with other elimination strategies on the LasHeR dataset.

| Layers | | | Success | Precision | NormPrec |
|---|---|---|---|---|---|
| 4 | 7 | 10 | | | |
| - | - | - | 56.1 | 70.4 | 66.6 |
| ✓ | - | - | **58.0** | **72.3** | **68.5** |
| ✓ | ✓ | - | 56.2 | 70.2 | 66.5 |
| ✓ | - | ✓ | 57.0 | 71.5 | 67.4 |
| - | ✓ | ✓ | 54.9 | 68.9 | 65.0 |
| ✓ | ✓ | ✓ | 57.8 | 72.1 | 68.3 |

Table 4: Ablation study on the insertion layers of the proposed TDTB module on the LasHeR dataset.

the superiority of the TMCE method.

**Inserting Layers of TDTB module**. To explore the optimal insertion layer for the TDTB module within the model, we evaluated its effectiveness when inserted after different layers of the backbone network using the LasHeR dataset. The results, shown in Tab.4, indicate that inserting the TDTB module only after the 4th layer, compared to inserting it after the 4th, 7th, and 10th layers, not only maintains the model's performance but even slightly improves it. Furthermore, using a single insertion significantly reduces the number of parameters and improves the model's speed.

We analyzed the reasons behind this phenomenon: The TDTB module integrates four templates (two modalities and two temporal states) and performs six cross-modal attention operations. When inserted after the 4th layer, it better facilitates the integration of information from different modalities, enhancing the representation of mid-level features and improving contextual consistency. These mid-level features are already sufficient for the task. However, when TDTB is inserted after higher layers, such as the 7th or 10th layers, the model has already captured highly complex semantic information. Introducing a complex fusion module at this stage may disrupt the learned patterns, ultimately leading to a performance drop. Additionally, due to the TMCE strategy employed in BTMTrack, performing cross-modal fusion after token filtering in the earlier and middle stages of the network allows the model to fully utilize the most relevant tokens. At later stages, as the number of tokens decreases significantly, the remaining information is already well-refined and focused. Introducing complex cross-modal fusion operations at this point could lead to redundancy or conflicts in the information. Specifically, cross-

|      | VIPT      | BAT       | TBSI              | GMMT              | BTMTrack          |
|------|-----------|-----------|-------------------|-------------------|-------------------|
| NO   | 84.1/68.4 | 90.2/73.3 | **91.4**/74.1     | 90.9/**74.3**     | 90.3/73.6         |
| PO   | 62.5/50.4 | 67.5/54.0 | 67.8/54.0         | 68.0/54.4         | **70.1/55.6**     |
| TO   | 57.7/46.2 | 64.1/**51.1** | **64.3**/51.0 | 64.1/50.9         | **64.3**/50.7     |
| HO   | 46.8/43.4 | 56.5/51.0 | 60.6/**53.4**     | 57.6/51.2         | **61.0/53.4**     |
| MB   | 57.6/46.0 | 62.4/49.6 | 63.1/49.5         | 64.0/50.5         | **64.8/50.6**     |
| LI   | 49.9/41.2 | 60.4/48.2 | 61.3/49.3         | 61.9/49.5         | **62.7/49.9**     |
| HI   | 67.8/54.2 | 75.3/59.6 | 73.8/58.2         | 75.2/59.3         | **79.1/62.2**     |
| AIV  | 36.3/34.2 | 51.4/45.3 | **58.2/49.8**     | 53.0/46.4         | 55.6/48.6         |
| LR   | 56.8/41.8 | 62.7/46.2 | 63.9/**47.3**     | 63.7/**47.3**     | **65.0**/47.3     |
| DEF  | 67.6/55.8 | **72.2**/58.5 | 71.6/58.7     | 70.8/58.0         | 72.1/**58.8**     |
| BC   | 65.0/51.9 | 67.7/53.9 | **69.9/55.7**     | 68.8/54.7         | 69.8/55.2         |
| SA   | 57.4/46.5 | 61.3/49.2 | 62.2/50.2         | 61.9/50.2         | **64.6/51.5**     |
| CM   | 62.0/50.0 | 68.0/54.4 | 69.5/55.0         | 69.0/54.7         | **71.1/56.4**     |
| TC   | 57.4/46.0 | 62.7/50.1 | 62.6/50.1         | 63.1/50.5         | **64.1/50.8**     |
| FL   | 59.7/46.9 | 62.0/49.0 | 60.9/47.5         | **63.6/49.8**     | 62.4/49.3         |
| OV   | 76.2/65.0 | **76.8/66.0** | 64.6/55.9     | 67.5/58.8         | 72.3/62.1         |
| FM   | 63.2/51.5 | 68.5/55.1 | 69.4/55.7         | 69.2/55.6         | **71.4/57.1**     |
| SV   | 65.0/52.5 | 69.7/55.9 | 70.2/56.2         | 70.6/56.7         | **73.0/58.1**     |
| ARC  | 59.4/49.5 | 63.1/51.9 | 64.3/52.5         | 64.7/52.8         | **67.1/54.5**     |

Table 5: Attribute-based Precision/Success scores on the LasHeR dataset. All metrics are evaluated based on the original results published by the authors of the compared methods.

modal fusion may involve information from different sources, and if this information is imprecise or conflicting, it could reduce the model's focus on core features.

Therefore, we adopt the design of inserting the TDTB module only after the 4th layer, which achieves optimal performance while reducing the model's parameter count and striking an excellent balance between speed and performance.

### 4.4. Analysis and Visualization

**Attribute-Based Performance**. To evaluate the performance of RGB-T tracking models under diverse challenges, the LasHeR dataset defines 19 attributes that represent real-world tracking difficulties. These attributes include no occlusion (NO), partial occlusion (PO), total occlusion (TO), hyaline occlusion (HO), motion blur (MB), low illumination (LI), high illumination (HI), abrupt illumination variation (AIV), low resolution (LR), deformation (DEF), background clutter (BC), similar appearance (SA), camera moving (CM), thermal crossover (TC), frame lost (FL), out-of-view (OV), fast motion (FM), scale variation (SV), and aspect ratio change (ARC). We compared our model with four of the most advanced models in recent years, including VIPT, BAT, TBSI, and GMMT, based on these attributes. The results, as shown in Tab.5, demonstrate that our model achieves the best performance across most scenarios. Notably, it exhibits significant advantages in conditions such as fast motion, scale variation, aspect ratio change, thermal crossover, high illumination, and low illumination. This further

highlights the importance of integrating temporal information and cross-modal information for feature fusion and selection.
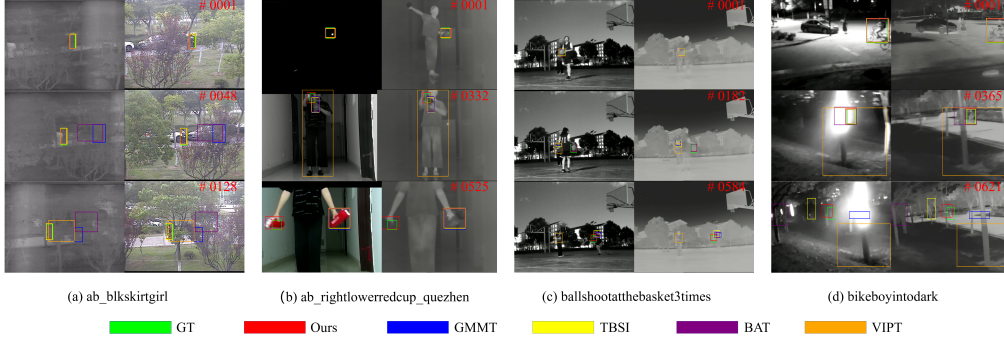


(a) ab_blkskirtgirl   (b) ab_rightlowerredcup_quezhen   (c) ballshootatthebasket3times   (d) bikeboyintodark

GT   Ours   GMMT   TBSI   BAT   VIPT

Figure 4: Qualitative comparison of our method with other RGB-T trackers on four representative sequences from the LasHeR dataset.

**Visualization**. As shown in Fig.4, we visually compared our model with four of the most advanced RGB-T trackers in recent years. To better demonstrate the robustness of our model, we selected four representative video sequences from the LasHeR dataset, which include challenges such as occlusion, low illumination, high illumination, and fast motion. For example, in the first video sequence, a girl walking on the road is partially occluded by trees; the second sequence presents the challenge of a low-light environment; in the third sequence, a basketball moves at high speed; and in the fourth sequence, high illumination creates significant difficulty in determining the position of a boy riding a bicycle. The results show that our model performs exceptionally well under a variety of challenging conditions.

## 5. Conclusion

In this paper, we propose an efficient cross-modal interaction method for RGB-T tracking that integrates dynamic templates into the backbone network. Our approach aims to fully exploit temporal and cross-modal information. Unlike previous methods, which either overlook or fail to effectively utilize temporal information, BTMTrack leverages the TMCE strategy to focus the model on target-related features, reducing computational overhead while avoiding irrelevant background noise. Additionally, we introduce the TDTB module to further enhance the interaction between templates and the search region by extracting more robust target-related contextual information. Evaluations on three widely used RGB-T benchmark datasets demonstrate that our method accurately locates targets and achieves state-of-the-art performance.

## Acknowledgments

# References

Alldieck, T., Bahnsen, C.H., Moeslund, T.B., 2016. Context-aware fusion of rgb and thermal imagery for traffic monitoring. Sensors 16, 1947.

Cao, B., Guo, J., Zhu, P., Hu, Q., 2024. Bi-directional adapter for multimodal tracking, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 927–935.

Cao, Z., Huang, Z., Pan, L., Zhang, S., Liu, Z., Fu, C., 2022. Tctrack: Temporal contexts for aerial tracking, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 14798–14808.

Chen, B., Li, P., Bai, L., Qiao, L., Shen, Q., Li, B., Gan, W., Wu, W., Ouyang, W., 2022. Backbone is all your need: A simplified architecture for visual object tracking, in: European Conference on Computer Vision, Springer. pp. 375–392.

Cui, Y., Jiang, C., Wang, L., Wu, G., 2022. Mixformer: End-to-end tracking with iterative mixed attention, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 13608–13618.

Dai, X., Yuan, X., Wei, X., 2021. Tirnet: Object detection in thermal infrared images for autonomous driving. Applied Intelligence 51, 1244–1261.

Fan, H., Yu, Z., Wang, Q., Fan, B., Tang, Y., 2024. Querytrack: Joint-modality query fusion network for rgbt tracking. IEEE Transactions on Image Processing .

Gao, S., Zhou, C., Zhang, J., 2023. Generalized relation modeling for transformer tracking, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 18686–18695.

Hong, L., Yan, S., Zhang, R., Li, W., Zhou, X., Guo, P., Jiang, K., Chen, Y., Li, J., Chen, Z., et al., 2024. Onetracker: Unifying visual object tracking with foundation models and efficient tuning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 19079–19091.

Hou, R., Ren, T., Wu, G., 2022. Mirnet: A robust rgbt tracking jointly with multi-modal interaction and refinement, in: 2022 IEEE International Conference on Multimedia and Expo (ICME), IEEE. pp. 1–6.

Hui, T., Xun, Z., Peng, F., Huang, J., Wei, X., Wei, X., Dai, J., Han, J., Liu, S., 2023. Bridging search region interaction with template for rgb-t tracking, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 13630–13639.

Li, C., Liang, X., Lu, Y., Zhao, N., Tang, J., 2019. Rgb-t object tracking: Benchmark and baseline. Pattern Recognition 96, 106977.

Li, C., Xue, W., Jia, Y., Qu, Z., Luo, B., Tang, J., Sun, D., 2021. Lasher: A large-scale high-diversity benchmark for rgbt tracking. IEEE Transactions on Image Processing 31, 392–404.

Li, C., Zhao, N., Lu, Y., Zhu, C., Tang, J., 2017. Weighted sparse representation regularized graph learning for rgb-t object tracking, in: Proceedings of the 25th ACM international conference on Multimedia, pp. 1856–1864.

Lin, L., Fan, H., Zhang, Z., Xu, Y., Ling, H., 2022. Swintrack: A simple and strong baseline for transformer tracking. Advances in Neural Information Processing Systems 35, 16743–16754.

Liu, L., Li, C., Xiao, Y., Ruan, R., Fan, M., 2024. Rgbt tracking via challenge-based appearance disentanglement and interaction. IEEE Transactions on Image Processing .

Lu, A., Li, C., Yan, Y., Tang, J., Luo, B., 2021. Rgbt tracking via multi-adapter network with hierarchical divergence loss. IEEE Transactions on Image Processing 30, 5613–5625.

Lu, A., Qian, C., Li, C., Tang, J., Wang, L., 2022. Duality-gated mutual condition network for rgbt tracking. IEEE Transactions on Neural Networks and Learning Systems .

Mei, J., Zhou, D., Cao, J., Nie, R., He, K., 2023. Differential reinforcement and global collaboration network for rgbt tracking. IEEE Sensors Journal 23, 7301–7311.

Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al., 2019. Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems 32.

Rezatofighi, H., Tsoi, N., Gwak, J., Sadeghian, A., Reid, I., Savarese, S., 2019. Generalized intersection over union: A metric and a loss for bounding box regression, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 658–666.

Ross, T.Y., Dollár, G., 2017. Focal loss for dense object detection, in: proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2980–2988.

Tang, Z., Xu, T., Wu, X., Zhu, X.F., Kittler, J., 2024. Generative-based fusion mechanism for multi-modal tracking, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 5189–5197.

Wang, H., Liu, X., Li, Y., Sun, M., Yuan, D., Liu, J., 2024. Temporal adaptive rgbt tracking with modality prompt, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 5436–5444.

Wei, X., Bai, Y., Zheng, Y., Shi, D., Gong, Y., 2023. Autoregressive visual tracking, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9697–9706.

Xia, J., Shi, D., Song, K., Song, L., Wang, X., Jin, S., Zhou, L., Cheng, Y., Jin, L., Zhu, Z., et al., 2023. Unified single-stage transformer network for efficient rgb-t tracking. arXiv preprint arXiv:2308.13764 .

Xiao, Y., Yang, M., Li, C., Liu, L., Tang, J., 2022. Attribute-based progressive fusion network for rgbt tracking, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 2831–2838.

Yan, B., Peng, H., Fu, J., Wang, D., Lu, H., 2021. Learning spatio-temporal transformer for visual tracking, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 10448–10457.

Ye, B., Chang, H., Ma, B., Shan, S., Chen, X., 2022. Joint feature learning and relation modeling for tracking: A one-stream framework, in: European Conference on Computer Vision, Springer. pp. 341–357.

Zeng, G., Zeng, B., Wei, Q., Hu, H., Zhang, H., 2024. Visual object tracking with mutual affinity aligned to human intuition. IEEE Transactions on Multimedia .

Zhang, L., Danelljan, M., Gonzalez-Garcia, A., Van De Weijer, J., Shahbaz Khan, F., 2019a. Multi-modal fusion for end-to-end rgb-t tracking, in: Proceedings of the IEEE/CVF International conference on computer vision workshops, pp. 0–0.

Zhang, L., Gonzalez-Garcia, A., Weijer, J.V.D., Danelljan, M., Khan, F.S., 2019b. Learning the model update for siamese trackers, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 4010–4019.

Zhang, P., Zhao, J., Wang, D., Lu, H., Ruan, X., 2022a. Visible-thermal uav tracking: A large-scale benchmark and new baseline, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8886–8895.

Zhang, Q., Liu, X., Zhang, T., 2022b. Rgb-t tracking by modality difference reduction and feature re-selection. Image and Vision Computing 127, 104547.

Zhang, T., Guo, H., Jiao, Q., Zhang, Q., Han, J., 2023. Efficient rgb-t tracking via cross-modality distillation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5404–5413.

Zhang, T., Liu, X., Zhang, Q., Han, J., 2021. Siamcda: Complementarity-and distractor-aware rgb-t tracking based on siamese network. IEEE Transactions on Circuits and Systems for Video Technology 32, 1403–1417.

Zhu, J., Lai, S., Chen, X., Wang, D., Lu, H., 2023. Visual prompt multi-modal tracking, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 9516–9526.