

SSTtrack: A Unified Hyperspectral Video Tracking Framework via Modeling Spectral-Spatial-Temporal Conditions

Yuzeng Chen^a, Qiangqiang Yuan^{a, b, *}, Yuqi Tang^c, Yi Xiao^a, Jiang He^a, Te Han^c, Zhenqi Liu^d, and Liangpei Zhang^e

^a School of Geodesy and Geomatics, Wuhan University, Wuhan, Hubei 430079, China

^b Hubei Luojia Laboratory, Wuhan, Hubei 430079, China

^c School of Geosciences and Info-Physics, Central South University, Changsha, China

^d College of Artificial Intelligence, Southwest University, Chongqing, 400715, China

^e State Key Laboratory of Information Engineering, Survey Mapping and Remote Sensing, Wuhan University, Wuhan, China

* Corresponding author

ARTICLE INFO

Keywords:

Hyperspectral video

Spectral awareness

Temporal awareness

Prompt learning

Multi-modal tracking

ABSTRACT

Hyperspectral video contains rich spectral, spatial, and temporal conditions that are crucial for capturing complex object variations and overcoming the inherent limitations (e.g., multi-device imaging, modality alignment, and finite spectral bands) of regular RGB and multi-modal video tracking. However, existing hyperspectral tracking methods frequently encounter issues including data anxiety, band gap, huge volume, and weakness of the temporal condition embedded in video sequences, which result in unsatisfactory tracking performance. To address these dilemmas, we propose a unified hyperspectral video tracking framework via modeling spectral-spatial-temporal conditions in an end-to-end fashion, dubbed SSTtrack. First, we design a multi-modal generation adapter (MGA) to explore the interpretability benefits of combining physical and machine models for learning the multi-modal generation and bridging the band gap. Then, we construct a spectral-spatial adapter (SSA) to dynamically transfer and interact with multiple modalities. Finally, we design a temporal condition adapter (TCA) for injecting the temporal condition to guide spectral and spatial feature representations to capture static and instantaneous object properties. SSTtrack follows the prompt learning paradigm with the addition of fewer trainable parameters (0.575M), resulting in superior performance in extensive comparisons. The code and model will be available at <https://github.com/YZCU/SSTtrack>.

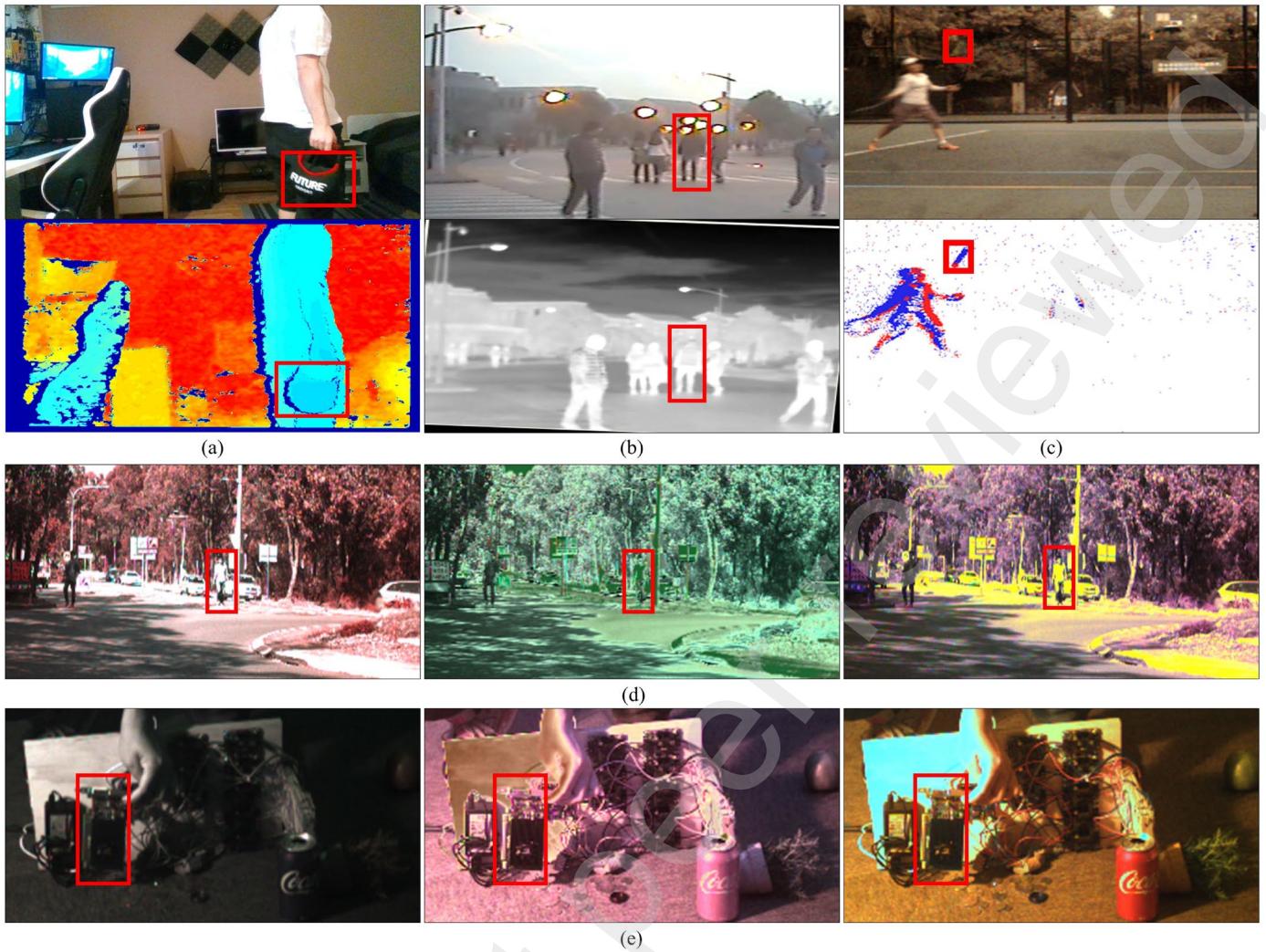
11

12 1. Introduction

Visual tracking, essential for establishing the association of the object in videos, has diverse applications including augmented reality, autonomous vehicles, and industrial automation (Javed et al., 2022). While progress using RGB modality is notable, challenges persist in scenarios with similar appearances, low light, and background clutter due to limited spectral information. To conquer these challenges, multi-modal tracking has emerged via combining RGB with event (RGB-E) (Zhu et al., 2023a), thermal infrared (RGB-T) (Cao et al., 2024), and depth (RGB-D) (Xuefeng Zhu et al., 2024) data, as shown in Fig. 1. However, regular multi-modal tracking requires multiple imaging devices, posing challenges in capturing the same scene accurately, especially for distant or small objects (Chen et al., 2024a). Consequently, aligning multiple modalities has become a common procedure (Li et al., 2022), albeit potentially causing image distortion, as illustrated in Figs. 1(a) and (b). Moreover, regular multi-modal tracking typically couples only two modalities in an RGB+X pattern (e.g., RGB-E, RGB-T, and RGB-D) (Zhu et al., 2023a), yielding suboptimal results in complex scenarios. Indeed, multi-modal data comprises information from various radiated bands (e.g., RGB and thermal infrared) or different mechanisms (e.g., event and depth). Hyperspectral (HS) cameras capture material-radiated signals across distinct bands, enabling trackers to identify object materials (Xiong et al., 2020). The abundant HS bands, captured from the same viewpoint, can generate multiple modalities. Thus, extending HS video tracking to multi-modal tracking is expected to address modality alignment and spectral limitations (Chen et al., 2024a). However, effectively utilizing multi-modal information,

40 including spectral, spatial, and temporal conditions, remains a challenge in HS video tracking.

41 Currently, HS video tracking encounters several dilemmas. First, data anxiety refers to the scarcity of HS videos, which hinders the direct training of robust tracking models (Liu et al., 2022). Second, the band gap between RGB and HS modalities makes it difficult to dynamically utilize rich spectral and spatial conditions (Chen et al., 2024b). Third, the huge volume caused by dense bands slows down tracking speed (Sun et al., 2023). At last, the temporal condition is far from being mined, albeit invaluable for robust tracking (Chen et al., 2024a). To address the data anxiety, one conceivable idea is to leverage large-scale RGB tracking datasets like TrackingNet (Muller et al., 2018), LaSOT (Fan et al., 2019), and COCO (Lin et al., 2014) for initial training of foundation models followed by performing full fine-tuning on HS training sets. However, fine-tuning poses risks of overfitting due to the scarcity of HS video datasets. Additionally, full fine-tuning is time-consuming and resource-intensive compared to promising prompt-tuning (Zhu et al., 2023a). To bridge the band gap, certain methods convert HS images into a three-channel representation, such as SiamHT (Tang et al., 2023), BAHT (Tang et al., 2022b), BS-SiamRPN (Wang et al., 2022), (Su et al., 2022), and (Zhang et al., 2022b). Nonetheless, this conversion inevitably results in the loss and distortion of the vital spectral condition (Chen et al., 2024b). Others focus on leveraging rich spectral conditions, yielding competitive results, such as DSP-Net (Zhu et al., 2023b), MHT (Xiong et al., 2020), SMT (Chen et al., 2023a), SiamBAG (Li et al., 2023a), CBFF-Net (Gao et al., 2023b), Trans-HST (Gao et al., 2023a), and TFTN (Zhao et al., 2022). Regarding the large volume, decision-level fusion methods are common but can limit efficiency



¹ Fig. 1. Sample of multi-modal data. (a) RGB and depth modalities. (b) RGB and thermal infrared modalities. (c) RGB and event modalities. (d) and (e) Two
² samples of multi-modal data generated by the hyperspectral modality.
³

4 by producing multiple weak results in each inference session, such as
5 SST-Net (Li et al., 2021), BRRF-Net (Ouyang et al., 2022), SEE-Net
6 (Li et al., 2023b), and BAE-Net (Li et al., 2020b). In addition,
7 previous studies have mainly focused on spatial-spectral conditions
8 and struggled to dynamically fuse visual feature prompts across
9 modalities. Few methods have integrated the temporal condition
10 contained in videos for HS tracking, such as SPIRIT (Chen et al.,
11 2024b), SENSE (Chen et al., 2024a), and SST-Net (Li et al., 2021).
12 Despite their success, these methods often necessitate the manual
13 design of update strategies and the introduction of hyper-parameters
14 (e.g., intervals and thresholds), which fall short of adequately
15 exploiting the temporal condition.

16 Motivated by the above analysis, we propose SSTtrack, a unified
17 HS video tracking framework via modeling spectral-spatial-temporal
18 conditions in an end-to-end fashion. Adhering to the prompt learning
19 paradigm, SSTtrack harnesses the potential of foundation models to
20 address data anxiety. The framework comprises three main
21 components: the multi-modal generation adapter (MGA), spectral-
22 spatial adapter (SSA), and temporal condition adapter (TCA). First,
23 we design the MGA to bridge the band gap and stimulate prior
24 knowledge. Drawing on the interpretability advantages of combining
25 physical and machine models, MGA achieves multi-modal
26 generation adaptively. Next, we construct the SSA to fully exploit
27 spectral and spatial conditions. SSA comprises a quintuple-stream
28 architecture with an interactive modal perception module to learn

29 prompts from other modalities. Finally, TCA is designed to inject the
30 temporal condition for guiding spectral and spatial feature
31 representations. Notably, the feature-level fusion strategy permeates
32 SSTtrack, alleviating the huge volume.

33 Major contributions are outlined as follows:

- 34 ● To fully model spectral-spatial-temporal conditions, we propose
35 SSTtrack, a unified hyperspectral video tracking framework via
36 modeling spectral-spatial-temporal conditions in an end-to-end
37 fashion. SSTtrack inherits the prompt learning paradigm by
38 adding fewer trainable parameters (0.575M), empowering the
39 foundation model to effectively tackle data anxiety.
- 40 ● A multi-modal generation adapter (MGA) is proposed to learn the
41 multi-modal generation by capitalizing on the interpretability
42 benefits of integrating physical and machine models. Benefiting
43 from MGA, SSTtrack can bridge the band gap to stimulate prior
44 knowledge of foundation models.
- 45 ● A spectral-spatial adapter (SSA) is proposed to mine the spectral-
46 spatial conditions. Within SSA, an interactive modal perception
47 module (IMPM) facilitates the transfer and interaction of
48 complementary features across modalities, enhancing spectral-
49 spatial feature representations.
- 50 ● A temporal condition adapter (TCA) is proposed to inject the
51 temporal condition to guide spectral-spatial conditions. Within
52 TCA, a spectral-spatial-temporal filtering module (SSTFM)
53 filters spectral-spatial-temporal conditions, facilitating the

1 integration of static and instantaneous object properties.
 2 We conduct extensive experiments to validate the SSTtrack
 3 framework. The remainder is organized as follows: Section II reviews
 4 related research. In Section III, we detail the proposed method.
 5 Section IV presents the results and analysis. Finally, Section V
 6 formulates the conclusions and highlights the main contributions.

7 2. Review on related research

8 2.1. Tracking via HS video

9 Generative and discriminative paradigms are prevalent in HS
 10 video tracking (Chen et al., 2024a). Initially, researchers concentrate
 11 on the generative paradigm, developing models to represent objects
 12 and identify similar regions. Recently, HS trackers have shifted
 13 towards the discriminative paradigm, which includes correlation
 14 filters and Siamese networks. Certain HS trackers aim to leverage full
 15 band information by correlation filter models, such as TASSCF
 16 (Tang et al., 2022a), MFI (Zhang et al., 2021), MHT (Xiong et al.,
 17 2020), and TSCFW (Hou et al., 2022). Despite the integration of
 18 hand-crafted and/or deep features, they have encountered limited
 19 success, largely due to the relative simplicity of the model (Li et al.,
 20 2023b). Siamese networks have received significant attention in the
 21 HS tracking domain. HS trackers like MMF-Net (Li et al., 2024),
 22 SiamOHOT (Sun et al., 2023), SEE-Net (Li et al., 2023b), SiamBAG
 23 (Li et al., 2023a), SiamHYPER (Liu et al., 2022), SiamCAT (Jiang et
 24 al., 2024), CBFF-Net (Gao et al., 2023b), DSP-Net (Zhu et al., 2023b),
 25 SiamHT (Tang et al., 2023), BRRF-Net (Ouyang et al., 2022), and
 26 Trans-HST (Gao et al., 2023a) have integrated Siamese networks to
 27 achieve outstanding performance, laying a solid foundation for our
 28 research. For instance, MMF-Net (Li et al., 2024) combines HS, false
 29 color, and material information to mitigate tracking drift. CBFF-Net
 30 (Gao et al., 2023b) introduces a bidirectional multiple deep feature
 31 fusion module and a cross-band group attention module to interact
 32 with multiple bands for stable performance. SiamCAT (Jiang et al.,
 33 2024) introduces a channel adaptive module and a guided learning
 34 attention module to improve effectiveness. However, previous
 35 studies have mainly focused on spatial-spectral conditions and
 36 struggled to dynamically fuse visual feature prompts across
 37 modalities. Few methods have integrated the temporal condition
 38 contained in videos for HS tracking, such as SPIRIT (Chen et al.,
 39 2024b), SST-Net (Li et al., 2021), and SENSE (Chen et al., 2024a).
 40 SPIRIT (Chen et al., 2024b) integrates an update module to assess the
 41 tracking confidence for adapting to object changes and mitigating
 42 tracking drift. SST-Net (Li et al., 2021) inherits the temporal attention
 43 with RNN-like structure to select valuable bands for deep ensemble
 44 tracking. While SENSE (Chen et al., 2024a) designs a motion-aware
 45 module comprising an awareness selector to determine the reliability
 46 of material and motion cues and a motion prediction scheme to
 47 manage abnormal states, facilitating continuous tracking. Despite
 48 their success, these methods often require manual design of update
 49 strategies and the introduction of hyper-parameters (e.g., intervals
 50 and thresholds), which falls short in adequately exploiting the
 51 temporal condition. Moreover, the aforementioned challenges of data
 52 anxiety, band gap, and huge volume inherently impede the robustness
 53 of HS video tracking. To alleviate these dilemmas, we propose
 54 SSTtrack, a unified HS video tracking framework via modeling
 55 spectral-spatial-temporal conditions in an end-to-end fashion,
 56 attaining comparable tracking results.

57 2.2. Multi-modal tracking via prompt learning

58 Recently, prompt learning has emerged as a technique that
 59 significantly enhances the performance of various language-
 60 processing tasks (Liu et al., 2023). Interestingly, this approach is
 61 gaining traction in the visual tracking community (Han et al., 2024;
 62 Zhu et al., 2023a). For instance, ProTrack (Yang et al., 2022) converts
 63 multi-modal inputs into a single modality using the prompt paradigm,
 64 achieving effective tracking capabilities learned from foundation
 65 models. ViPT (Zhu et al., 2023a) introduces prompt learning to multi-
 66 modal tracking by learning prompts relevant to each modality.
 67 OneTrack (Han et al., 2024) employs prompt-tuning techniques for
 68 RGB+X tracking tasks by pre-training the foundation model on the
 69 RGB dataset. Benefiting from the prompt learning, BAT (Cao et al.,
 70 2024) introduces a bi-directional adapter to improve multi-modal
 71 tracking based on the RGB foundation model. Considerable research
 72 has demonstrated the effectiveness and generalization of prompt
 73 learning in the RGB+X tracking community. In this work, we focus
 74 on prompt learning for the HS tracking community.

75 2.3. Tracking with temporal condition

76 Temporal condition is crucial in capturing object state changes and
 77 motion patterns (Cao et al., 2023). Consequently, mainstream studies
 78 have delved into temporal conditions in the visual tracking
 79 community (Lei et al., 2024; Wei et al., 2023; Xie et al., 2024). One
 80 common strategy involves updating the appearance representation.
 81 For instance, certain studies like Stark (Yan et al., 2021a) and
 82 SeqTrack (Chen et al., 2023b) use dynamic templates to update the
 83 object appearance for capturing changes, enhancing the match
 84 between the template and the search images. Others concentrate on
 85 learning features that encode information about the previous state or
 86 motion, such as KYS (Bhat et al., 2020) and SwinTrack (Lin et al.,
 87 2022). Additionally, integrating historical appearances is another
 88 effective approach for leveraging temporal conditions, such as
 89 UpDateNet (Zhang et al., 2019a). Recent trackers like JTBP (Lei et
 90 al., 2024), TCTrack++ (Cao et al., 2023), EVPTrack (Shi et al.,
 91 2024), TATrack (Wang et al., 2024), AQATrack (Xie et al., 2024),
 92 Seqtrack (Chen et al., 2023b), and ARTTrack (Wei et al., 2023) have
 93 thoroughly demonstrated the significance of temporal condition in
 94 tracking community from various perspectives. Despite the
 95 remarkable achievements of existing RGB/RGB+X trackers, they
 96 predominantly focus on the spatial and temporal aspects, limiting
 97 their performance in complex environments where the spectral
 98 condition is crucial. Considering the substantial spectral advantages
 99 offered by HS video, we propose a unified HS video tracking
 100 framework SSTtrack that models spectral-spatial-temporal
 101 conditions in an end-to-end fashion, thereby enhancing tracking
 102 robustness in challenging scenarios.

103 3. Method

104 This section starts by overviewing our SSTtrack framework. We
 105 then delve into its key components including the multi-modal
 106 generation adapter, spectral-spatial adapter, and temporal condition
 107 adapter. Finally, we introduce the head and loss function.

108 3.1. Overview

109 As shown in Fig. 2, the proposed SSTtrack framework primarily
 110 consists of MGA, SSA, and TCA. The process begins with feeding a
 111 pair of HS patches (i.e., the search patch $X \in \mathbb{R}^{H_x \times W_x \times 16}$ and

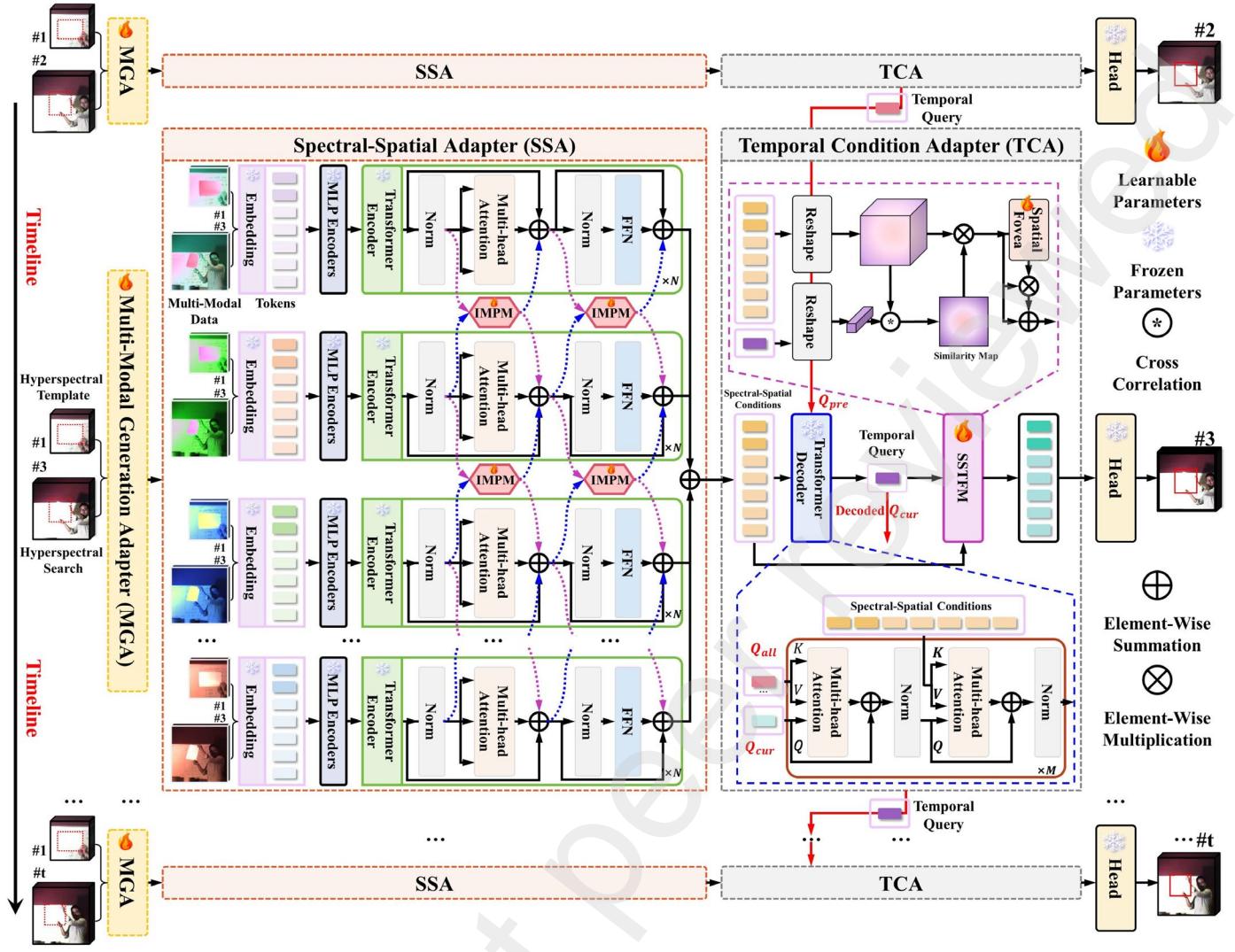


Fig. 2. Overview of the proposed SSTtrack: a unified hyperspectral video tracking framework via modeling spectral-spatial-temporal conditions. SSTtrack comprises three main components: the multi-modal generation adapter (MGA), spectral-spatial adapter (SSA), and temporal condition adapter (TCA). First, MGA achieves the multi-modal generation adaptively for bridging the band gap and stimulating the prior knowledge. Then, SSA follows a quintuple-stream architecture with an interactive modal perception module (IMPM) to learn prompts from other modalities for full exploitation of spectral and spatial conditions. While the temporal condition adapter (TCA) aims to inject the temporal condition for guiding spectral and spatial feature representations. At last, the output is fed into the head network for predicting results.

8 template patch $Z \in \mathbb{R}^{H_z \times W_z \times 16}$) into MGA to adaptively generate 9 multi-modal data for capturing the rich spectral condition. Right after 10 that, SSA dynamically extracts multi-modal complementary 11 information using a quintuple-stream architecture, equipped with 12 IMPM to learn feature prompts across modalities. For TCA, it has 13 two inputs. The first one is spectral-spatial conditions derived from 14 SSA, while the second one is several autoregressive and learnable 15 temporal queries. It aims to inject the temporal condition to guide 16 spectral-spatial conditions. Within TCA, the SSTFM filters spectral- 17 spatial-temporal conditions, facilitating the integration of static and 18 instantaneous object properties. Finally, the output features of TCA 19 will be forwarded to the head network for result prediction.

3.2. Multi-modal generation adapter (MGA)

HS modality compiles the richer spectral condition, potentially 22 overcoming the shortcomings of RGB modality (Li et al., 2023b). 23 However, the scarcity of large-scale HS datasets poses a challenge 24 for training generalized HS tracking models, especially given the 25 data-driven models, which increases the risk of overfitting. To

address this issue, we inherit the prompt learning to leverage the pre- 27 learned knowledge embedded in RGB foundation models. 28 Specifically, we introduce the MGA to learn the multi-modal 29 generation by capitalizing on the interpretability benefits of 30 integrating physical and machine models. MGA mainly consists of 31 the inception prompt enhancing module (IPEM) and the spectral 32 contribution prompting module (SCPM).

3.2.1. Inception prompt enhancing module (IPEM)

The self-expression model selectively identifies informative bands 34 from C bands using the coefficient matrix $G \in \mathbb{R}^{C \times C}$ (Cai et al., 35 2020). The matrix reveals the relationships among bands and their 36 contributions to downstream tracking tasks. In the HS video, each 37 frame T is denoted as $T = [t_1, t_2, \dots, t_C] \in \mathbb{R}^{D \times C}$, where $D = m \times n$ 38 signifies the pixel count, and $t_i \in \mathbb{R}^{m \times n}$ represents the vector for the 39 i -th band. The HS self-expression model is thus expressed as:

$$\text{argmin} \|G\|_{1,2}, \text{s.t. } T = TG + E, \text{diag}(G) = 0, G \geq 0$$

$$41 \quad \|G\|_{1,2} = \sum_{i=1}^C \|g^i\|_2, \sum_{i=1}^C g_{ij} = 1, \forall j. \quad (1)$$

1 where $E \in \mathbb{R}^{D \times C}$ represents the residual item. $\text{diag}(G) = 0$ is
 2 enforced to prevent trivial solutions. $G \geq 0$ ensures that each element
 3 of g_j signifies the probability of representing t_j . $\|G\|_{1,2}$ denotes the
 4 sum of l_2 -norm of all row vectors g^i . To harness the interpretability
 5 benefits of combining physical and machine models, we introduce the
 6 IPEM (Fig. 3). This module innovatively learns the solution for the
 7 matrix G of the self-expression physical model in a learning-to-
 8 optimize manner, circumventing the laborious iterative optimization
 9 process (Eq. 1). The IPEM extracts visual prompts with diverse fine-
 10 grained features, adapting to changes in object scale. Specifically, the
 11 HS image is divided into patches of size (P, P) , each of which is then
 12 mapped to a 2-D tensor through a projection layer. With P values set
 13 to 18, 16, and 14, we capture features at different scales, denoted as
 14 fp_{18} , fp_{16} , and fp_{14} . Spatial max pooling Max and spatial average
 15 pooling Avg are applied to extract salient information, followed by
 16 summation, reshaping, and concatenation to obtain the token ft_{cat} .
 17 Finally, the token is fed into an MLP layer for solving solve G . This
 18 process is formulated as follows:

$$19 \quad \begin{cases} ft_{18} = R(\text{Max}(fp_{18}) \oplus \text{Avg}(fp_{18})), \\ ft_{16} = R(\text{Max}(fp_{16}) \oplus \text{Avg}(fp_{16})), \\ ft_{14} = R(\text{Max}(fp_{14}) \oplus \text{Avg}(fp_{14})), \end{cases} \quad (2)$$

$$20 \quad ft_{cat} = \text{Cat}(ft_{18}, ft_{16}, ft_{14}), \quad (3)$$

$$21 \quad Y = MLP(MLP(ft_{cat}) \oplus ft_{cat}), \quad (4)$$

$$22 \quad G = Y^T Y, \quad (5)$$

23 where R denotes the Reshape. Cat represents concatenation.

24 3.2.2. Spectral contribution prompting module (SCPM)

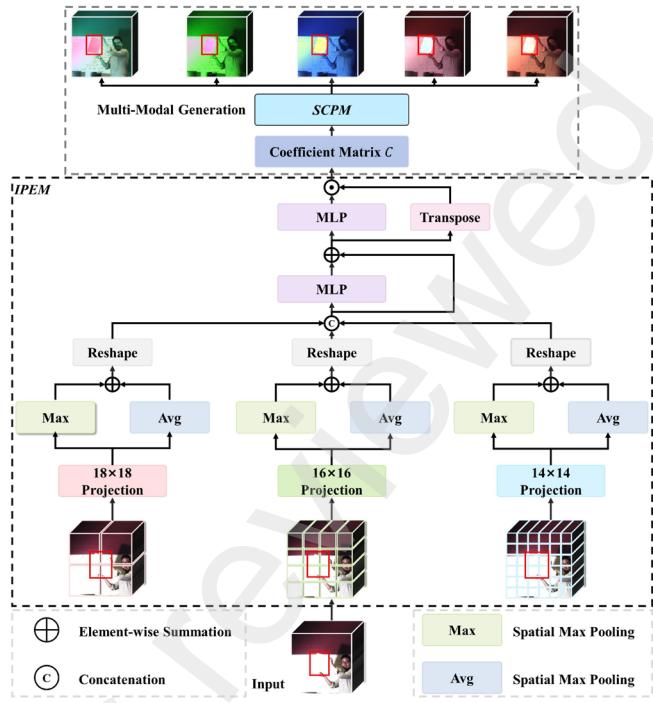
25 Upon IPEM, we derive the matrix G . Within G , the j -th column, i -
 26 th row, and (i, j) -th element are indicated by g_j , g^i , and g_{ij} ,
 27 respectively. g_j denotes the coefficient of the j -th band represented
 28 by all the remaining bands, while g^i denotes the contribution of the
 29 i -th band to the reconstruction. The importance of a band is reflected
 30 by its contribution, which allows us to rank them accordingly. Thus,
 31 G unveils the relationship between bands and facilitates multi-modal
 32 generation (Chen et al., 2024a).

33 Specifically, G undergoes column-wise normalization to yield
 34 $\hat{g}_j = |g_{ij}| / \|g_j\|_2$, for all i , followed by $z_i = \|\hat{g}^i\|_1$, where \hat{g}_j is the
 35 normalization output of the j -th column, and \hat{g}^i stands for the i -th
 36 row of the normalized G . $Z \in \mathbb{R}^{C \times 1}$ indicates the desired cumulative
 37 contributions for each HS band. Subsequently, we sort all HS bands
 38 in descending order and group them to obtain multi-modal
 39 representations, denoted as $[q_1, q_2, \dots, q_k]$, where $k = \text{int}(C/3)$ and
 40 $q_i \in \mathbb{R}^{m \times n \times 3}$.

41 3.3. Spectral-spatial adapter (SSA)

42 As depicted in Fig. 2, we propose the SSA to fully exploit spectral
 43 and spatial conditions. SSA comprises a quintuple-stream
 44 architecture with IMPM to cross-prompt multi-modal features. Each
 45 stream shares the parameters.

46 First, the multi-modal templates $\mathcal{L}_{tem} = [\mathcal{L}_{tem}^1, \mathcal{L}_{tem}^2, \dots, \mathcal{L}_{tem}^k]$
 47 and searches $\mathcal{L}_{sea} = [\mathcal{L}_{sea}^1, \mathcal{L}_{sea}^2, \dots, \mathcal{L}_{sea}^k]$ undergo the patch
 48 embedding and position embedding to yield tokens $x_i =$
 49 $[x_0^1, x_0^2, \dots, x_0^k]$. Then, the IMPM is embedded within the i -th layer
 50 of the Transformer encoder, penetrating the encoders of multiple
 51 modalities. At each encoder's $(i+1)$ layer, it learns to fuse the
 52 modality-specific features with complementary information from
 53



55 56 **Fig. 3.** Overview of the proposed multi-modal generation adapter (MGA). It
 57 consists of IPEM for solving the coefficient matrix and SCPM for generating
 58 multi-modal data.

59 other modalities, progressively refining feature prompts in a layer-
 60 by-layer fashion, as follows:

$$61 \quad (x_{i+1}^1, x_{i+1}^2, \dots, x_{i+1}^k) = \mathcal{F}_i^A(x_i^1, x_i^2, \dots, x_i^k), i = 1, 2, \dots, N, \quad (6)$$

62 where \mathcal{F}_i^A means the quintuple-stream encoder layer paralleled with
 63 IMPM. The N -layer transformer encoder progressively and
 64 dynamically extracts multi-modal features of the object. Finally, the
 65 features of the multi-branch are summed.

66 The structure of IMPM (Fig. 4) is essentially a bottleneck design
 67 for reducing parameter burden. It comprises the down-/mid-/up-
 68 projection layers. Rectified Linear Unit (ReLU) activation serves to
 69 introduce non-linearity, while a shortcut structure is implemented to
 70 mitigate potential interference from low-level feature maps.
 71 Furthermore, a scale factor is integrated to constrain feature
 72 contributions. The output of IMPM is fed to the encoder layer of
 73 another modality as feature prompts, facilitating efficient cross-
 74 modal tracking. Taking the transfer of $x_i^2 \rightarrow x_i^1$ for example, the
 75 formula is present as follows:

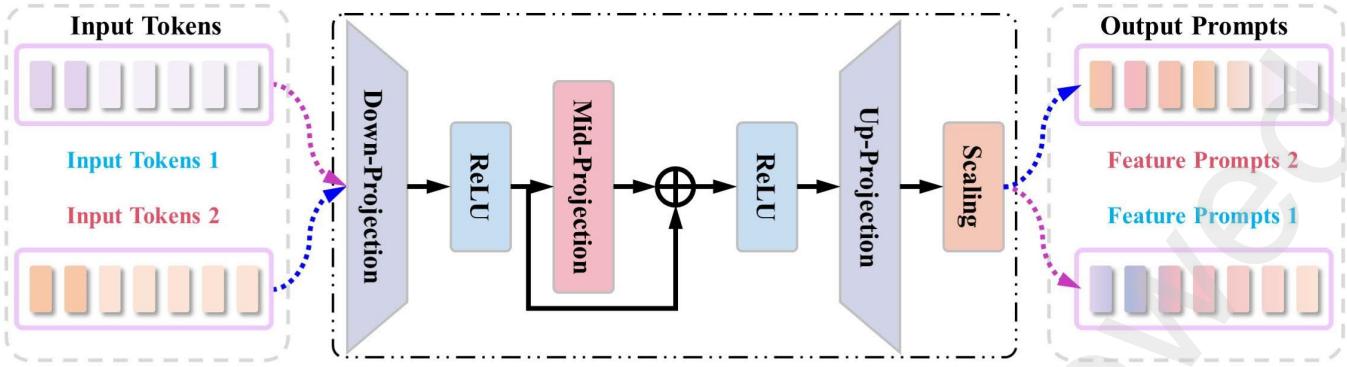
$$76 \quad \begin{cases} x_i^{1'} = x_i^1 \oplus \mathcal{F}^{Att}(x_i^1) \oplus P_i^2 \\ P_i^2 = \mathcal{F}^{IMPM}(x_i^2), i = 1, 2, \dots, N \end{cases} \quad (7)$$

$$77 \quad \begin{cases} x_{i+1}^{1'} = x_i^{1'} \oplus \mathcal{F}^{MLP}(x_i^{1'}) \oplus P_i^{2'} \\ P_i^{2'} = \mathcal{F}^{IMPM}(x_i^{2'}), i = 1, 2, \dots, N \end{cases} \quad (8)$$

78 where $\mathcal{F}^{Att}(\cdot)$, $\mathcal{F}^{IMPM}(\cdot)$, and $\mathcal{F}^{MLP}(\cdot)$ mean the multi-head attention
 79 (MHA), IMPM, and multi-layer perceptron (MLP) operations,
 80 respectively. P_i^2 is the feature prompt extracted derived from the x_i^2
 81 modality. The formula for $\mathcal{F}^{IMPM}(x_i^2)$ is denoted as follows:

$$82 \quad \mathcal{F}^{IMPM}(x_i^2) = UP \left(\delta \left(MP \left(\delta \left(DP(x_i^2) \right) \right) \oplus \delta \left(DP(x_i^2) \right) \right) \right) \cdot s, \quad (9)$$

83 where DP , MP , and UP denote the down-projection, mid-projection,
 84 and up-projection layers, respectively. δ denotes the ReLU
 85 activation, while s is the scale factor.



2 Fig. 4. Overview of the proposed interactive modal perception module (IMPM).

3.3.4. Temporal condition adapter(TCA)

4 Tracking challenges such as motion and similar objects cannot be
 5 effectively tackled when using solely spectral-spatial conditions.
 6 Therefore, we design the TCA (Fig. 2) to inject the temporal
 7 condition for guiding spectral-spatial feature representations. TCA
 8 comprises the Transformer decoder for temporal condition transfer
 9 and SSTFM for filtering spectral-spatial-temporal conditions.

10 The Transformer decoder has two inputs: temporal queries and
 11 spectral-spatial conditions from the SSA. There are two types of
 12 queries: Q_{pre} and Q_{cur} . Q_{pre} is passed down from previous frames.
 13 Q_{cur} is a query learned from the current frame, implying the object
 14 state while considering spectral-spatial-temporal conditions. As
 15 illustrated in Fig. 2, the Transformer decoder's queries comprise Q_{all} ,
 16 Q_{cur} , and the middle parameter Q_{pre} , defined as follows:

$$17 \quad Q_{all} = \text{Cat}(Q_{pre}, Q_{cur}), \quad (10)$$

$$18 \quad Q_{pre} = \begin{cases} \text{Cat}(Q_1, \dots, Q_{i-1}), & t < l, \\ \text{Cat}(Q_{i-l+1}, \dots, Q_{i-1}), & t \geq l, \end{cases} \quad (11)$$

19 where l denotes the length of spectral-spatial-temporal conditions
 20 (tokens). The decoded Q_{cur} is propagated to the next frame of the
 21 video clip as one of the Q_{pre} in a sliding window mode. The decoder
 22 (Fig. 2) comprises M layers, each primarily consisting of two MHAs
 23 and a feed-forward neural network (FFN). The Transformer decoder
 24 interacts with historical queries to generate the decoded query,
 25 capturing temporal states and motion properties. The formula for
 26 MHA is as follows:

$$27 \quad \left\{ \begin{array}{l} \text{MultiHead}(Q, \mathcal{K}, \mathcal{V}) = \text{Cat}(\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_{n_h}) \mathcal{W}^O, \\ \mathcal{H}_i = \text{Attention}(Q \mathcal{W}_i^Q, \mathcal{K} \mathcal{W}_i^K, \mathcal{V} \mathcal{W}_i^V), \\ \text{Attention}(q, k, v) = \text{Softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)v, \end{array} \right. \quad (12)$$

28 where $\mathcal{W}^O \in \mathbb{R}^{n_h \times d_m}$, $\mathcal{W}_i^Q \in \mathbb{R}^{d_m \times d_k}$, $\mathcal{W}_i^K \in \mathbb{R}^{d_m \times d_k}$, and $\mathcal{W}_i^V \in$
 29 $\mathbb{R}^{d_m \times d_v}$ are learnable parameters. $n_h = 8$ is the number of heads.
 30 $d_m = 512$ signifies the hidden size, and $d_k = d_v = d_m/8 = 64$.

31 The SSTFM (Fig. 2) is designed to filter spectral, spatial, and
 32 temporal conditions with almost no learnable parameters. We first
 33 obtain the search region token f_{sea} from the encoder output. Next, we
 34 calculate the similarity between f_{sea} and the decoded token f_{tem} . The
 35 similarity map undergoes element-wisely multiplication with the
 36 search region feature followed by spatial fovea and shortcut
 37 operations to enhance the significant regions while filtering out less
 38 discriminative ones.

39 Subsequently, the output feature is forwarded to the head network.
 40 In summary, the process is described as follows:

$$41 \quad \left\{ \begin{array}{l} f_{sea'} = (R(f_{sea}) \odot R(f_{dec})) \otimes f_{sea}, \\ f_{sea''} = f_{sea'} \otimes f_{fovea}, \\ f_{fovea} = \left\{ \begin{array}{l} e^{f_{sea'}^{[i,j]}} \\ \sum e^{f_{sea'}^{[i,j]}} \lambda \end{array} \right\}, \\ f_{sea} = f_{sea''} \oplus f_{sea'}, \end{array} \right. \quad (13)$$

42 where $i = 1, 2, \dots, H_s$ and $j = 1, 2, \dots, W_s$. H_s and W_s are the height
 43 and width of the reshaped feature of f_{sea} . \odot denotes the cross-
 44 correlation operation, and λ is a learnable weight.

45 3.5. Head and loss

46 The center-based head network is used for prediction by producing
 47 three outputs: the classification score map with $\mathbb{R}^{\frac{H_x \times W_x}{P}}$, the bounding
 48 box size with $\mathbb{R}^{2 \times \frac{H_x \times W_x}{P}}$, and the offset size with $\mathbb{R}^{2 \times \frac{H_x \times W_x}{P}}$. The
 49 object's location is determined by the highest classification score, and
 50 the final result considers both the bounding box size and offset size.
 51 We apply the weighted focal loss (Lin et al., 2017) L_{cls} for
 52 classification, L1 loss L_1 and GIoU loss (Rezatofighi et al., 2019)
 53 L_{iou} for bounding box regression, and another L1 loss L_{rec} for
 54 spectral reconstruction. The overall loss L_{tot} is obtained by:

$$55 \quad L_{tot} = L_{cls} + \varphi_{L_1} L_1 + \varphi_{iou} L_{iou} + \varphi_{rec} L_{rec} \quad (14)$$

56 where φ_{L_1} , φ_{iou} , and φ_{rec} are regularization parameters and set to 5,
 57 2, and 2, respectively.

58 4. Experimental results and analysis

59 4.1. Datasets

60 SSTtrack is trained and evaluated on HS datasets acquired from
 61 the Hyperspectral Object Tracking Competition (HOTC), with 35
 62 sets for testing and 40 sets for training (Xiong et al., 2020). Each
 63 dataset has frame-level annotations and three types of videos: HS,
 64 false color, and RGB. The false color video is created from the
 65 corresponding HS video, while the RGB video is captured from a
 66 similar perspective. The HOTC dataset facilitates tracker
 67 performance assessment, featuring 11 attributes: Scale Variation
 68 (SV), Motion Blur (MB), Fast Motion (FM), Low Resolution (LR),
 69 Out-of-View (OV), Occlusion (OCC), Background Clutter (BC),
 70 Illumination Variation (IV), Deformation (DEF), Out-of-Plane
 71 Rotation (OPR), and In-Plane Rotation (IPR).

72 4.2. Implementation details

73 The proposed SSTtrack employs PyTorch 2.0.0 in Python 3.8 and

1 **Table 1**

2 Comparative results and properties of state-of-the-art RGB tracking methods. Trackers are listed in chronological order.

No.	Tracker	Venue	Feature/Backbone/Tag	RGB		FAC/HS		PreD	SucD
				Pre	Suc	Pre	Suc		
1	CSK (Henriques et al., 2012)	ECCV 2012	I	0.575	0.331	0.615	0.343	-4.0%	-1.2%
2	CN (Danelljan et al., 2014)	CVPR 2014	CN+I	0.646	0.380	0.643	0.379	0.3%	0.1%
3	SAMF (Li et al., 2015)	ECCV 2015	HOG+CN+I	0.693	0.418	0.660	0.388	3.3%	3.0%
4	DAT (Possegger et al., 2015)	CVPR 2015	CH	0.647	0.394	0.542	0.327	10.5%	6.7%
5	KCF (Henriques et al., 2015)	TPAMI 2015	HOG	0.613	0.377	0.586	0.358	2.7%	1.9%
6	SRDCF (Danelljan et al., 2015)	ICCV 2015	HOG	0.846	0.568	0.830	0.554	1.6%	1.4%
7	Staple (Bertinetto et al., 2016a)	CVPR 2016	HOG+CN	0.801	0.518	0.770	0.507	3.1%	1.1%
8	DSST (Danelljan et al., 2017)	TPAMI 2017	HOG+I	0.776	0.504	0.731	0.480	4.5%	2.4%
9	BACF (Galoogahi et al., 2017)	ICCV 2017	HOG	0.793	0.533	0.819	0.544	-2.6%	-1.1%
10	CSRDCF (Lukezic et al., 2018)	IJCV 2018	HOG+CN+CH	0.870	0.563	0.833	0.533	3.7%	3.0%
11	STRCF (Li et al., 2018b)	CVPR 2018	HOG+CN	0.829	0.569	0.838	0.568	-0.9%	0.1%
12	ARCF (Huang et al., 2019)	ICCV 2019	HOG+CN+I	0.818	0.564	0.798	0.542	2.0%	2.2%
13	AutoTrack (Li et al., 2020a)	CVPR 2020	HOG+CN+I	0.831	0.534	0.818	0.534	1.3%	0.0%
14	SiamRPN (Li et al., 2018a)	CVPR 2018	AlexNet	0.902	0.592	0.757	0.486	14.5%	10.6%
15	DaSiamRPN (Zhu et al., 2018)	ECCV 2018	AlexNet	0.878	0.622	0.850	0.575	2.8%	4.7%
16	ATOM (Danelljan et al., 2019)	CVPR 2019	ResNet-18	0.917	0.614	0.867	0.556	5.0%	5.8%
17	DiMP (Bhat et al., 2019)	ICCV 2019	ResNet-50	0.944	0.641	0.836	0.556	10.8%	8.5%
18	SiamRPN++ (Li et al., 2019)	CVPR 2019	ResNet-50	0.912	0.653	0.847	0.591	6.5%	6.2%
19	UpdateNet (Zhang et al., 2019a)	ICCV 2019	AlexNet/DaSiamRPN	0.863	0.595	0.833	0.551	3.0%	4.4%
20	SiamDW (Zhang et al., 2019b)	CVPR 2019	CIRNext22/SiamFC	0.872	0.565	0.812	0.529	6.0%	3.6%
21	SiamMask (Wang et al., 2019)	CVPR 2019	ResNet-50	0.877	0.611	0.813	0.554	6.4%	5.7%
22	PrDiMP (Danelljan et al., 2020)	CVPR 2020	ResNet-50	0.917	0.634	0.829	0.565	8.8%	6.9%
23	SiamBAN (Chen et al., 2020)	CVPR 2020	ResNet-50	0.853	0.610	0.863	0.587	-1.0%	2.3%
24	SiamFC++ (Xu et al., 2020)	AAAI 2020	GoogLeNet	0.865	0.635	0.820	0.578	4.5%	5.7%
25	KeepTrack (Mayer et al., 2021)	ICCV 2021	ResNet-50	0.951	0.656	0.900	0.617	5.1%	3.9%
26	SiamGAT (Guo et al., 2021)	CVPR 2021	GoogLeNet	0.889	0.649	0.820	0.576	6.9%	7.3%
27	LightTrack (Yan et al., 2021b)	CVPR 2021	Custom	0.814	0.593	0.761	0.530	5.3%	6.3%
28	Stark (Yan et al., 2021a)	ICCV 2021	ResNet-50/ST	0.900	0.637	0.814	0.579	8.6%	5.8%
29	SiamCAR (Cui et al., 2022)	IJCV 2022	ResNet-50	0.882	0.636	0.846	0.586	3.6%	5.0%
30	OSTrack (Ye et al., 2022)	ECCV 2022	ViT-Base/256	0.910	0.662	0.869	0.620	4.1%	4.2%
31	SimTrack (Chen et al., 2022)	ECCV 2022	ViT-Base/16	0.925	0.664	0.852	0.602	7.3%	6.2%
32	SBT (Xie et al., 2022)	CVPR 2022	SBT-Base	0.938	0.677	0.866	0.626	7.2%	5.1%
33	GRM (Gao et al., 2023c)	CVPR 2023	ViT-Large/320	0.921	0.671	0.892	0.648	2.9%	2.3%
34	SeqTrack (Chen et al., 2023b)	CVPR 2023	ViT-Base/256GOT	0.917	0.643	0.856	0.594	6.1%	4.9%
35	ARTTrack (Wei et al., 2023)	CVPR 2023	ViT-Base/256	0.929	0.673	0.899	0.640	3.0%	3.3%
36	SMAT (Xie et al., 2024)	WACV 2024	MobileViTv2	0.894	0.637	0.831	0.581	6.3%	5.6%
37	AQATrack (Xie et al., 2024)	CVPR 2024	HiViT-Base/256	0.926	0.670	0.880	0.640	4.6%	3.0%
38	SSTtrack	Ours	HiViT-Base	n/a	n/a	0.956	0.713	n/a	n/a

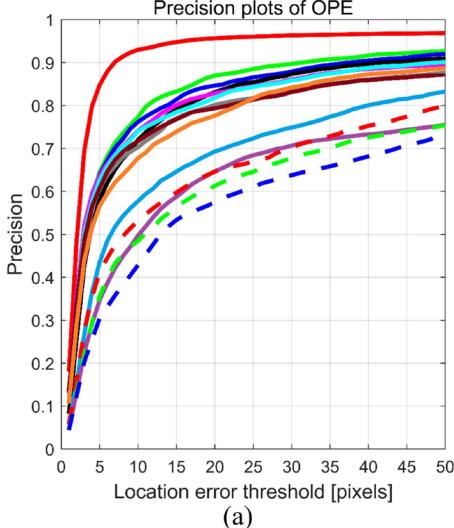
3 Benchmark scores are color-coded for clarity: highest (**red**), second (**green**), and third (**blue**). Video categories are represented as RGB (red-green-blue), FAC (false color),
4 and HS (hyperspectral). Pre-degradation and Suc-degradation transformations from RGB to false color are abbreviated as PreD and SucD, respectively. The dashed line
5 separates trackers employing hand-crafted and deep features. n/a indicates not applicable.

6 is trained on an NVIDIA RTX 3090 GPU. The foundation model
7 undergoes training for 150 epochs on large-scale RGB (three bands)
8 data, including GOT-10k (Huang et al., 2021), TrackingNet (Muller
9 et al., 2018), LaSOT (Fan et al., 2019), and COCO (Lin et al., 2014).
10 The optimizer is AdamW (Loshchilov et al., 2018) with a weight
11 decay of 10^{-4} , an initial learning rate of 4×10^{-5} for the encoder with
12 HiViT (Zhang et al., 2022a) as the backbone, and other parameters of
13 4×10^{-4} . Subsequently, prompt-tuning is conducted on HS data with
14 16 bands for 30 epochs using the same AdamW optimizer. The initial
15 learning rate is set to 4×10^{-4} , decayed to 4×10^{-5} at the 25-th epoch.

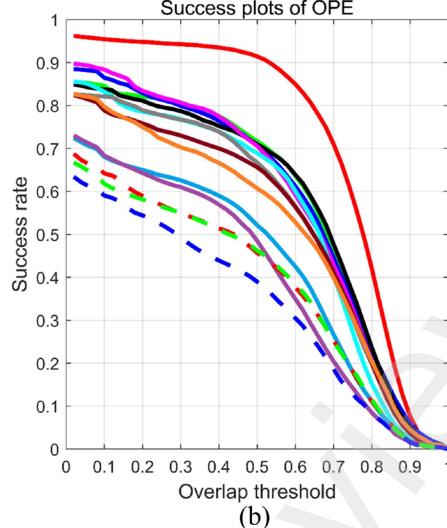
16 Template and search region sizes are fixed at 128 and 256 pixels,
17 respectively, with a batch size of 16 using video-level sampling. The
18 tracker is tested at speed on an NVIDIA RTX 4060 GPU.

19 4.3. Assessment metrics

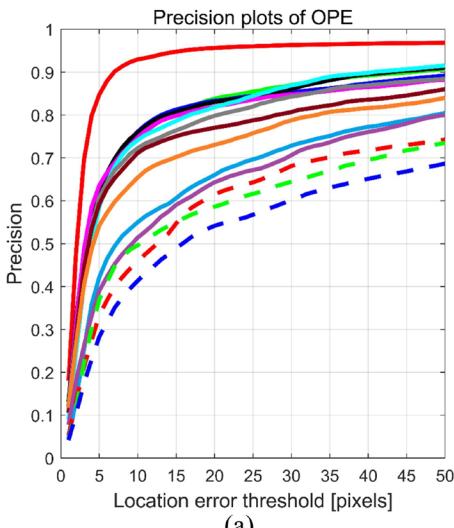
20 Trackers are subjected to comparative analysis through one-pass
21 evaluations utilizing precision and success plots (Wu et al., 2015).
22 Precision plots measure the percentage of frames where the center
23 location error (ν) falls below predefined pixel thresholds ranging
24 from 1 to 50 pixels. ν is defined as $\nu = \sqrt{(x - X)^2 + (y - Y)^2}$,



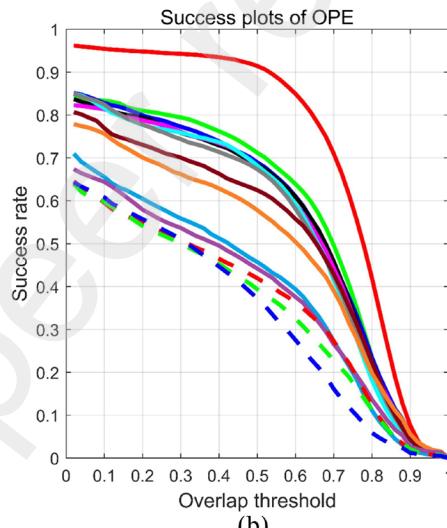
(a)



(b)

Fig. 5. Comparative results with hand-crafted feature-based trackers on RGB videos. (a) Precision plot. (b) Success plot.

(a)



(b)

Fig. 6. Comparative results with hand-crafted feature-based trackers on false color videos. (a) Precision plot. (b) Success plot.

6 where (x, y) and (X, Y) are the center of the prediction r_t and the
7 ground truth r_g . Success plots assess the percentage of successful
8 frames where the overlap score (s) exceeds thresholds ranging from
9 0 to 1. s is computed using $s = |r_t \cap r_g| / |r_t \cup r_g|$, where union \cup
10 and intersection \cap are mathematical set operations, and $|\cdot|$ is the pixel
11 count in a given region. Trackers are ranked according to their
12 precision at 20 pixels of the precision plot and the area under the
13 curve of the success plot, i.e., Pre and Suc. The inference speed is
14 assessed in frames per second (FPS). Moreover, the Params (M) is
15 also present for detailed ablation comparisons.

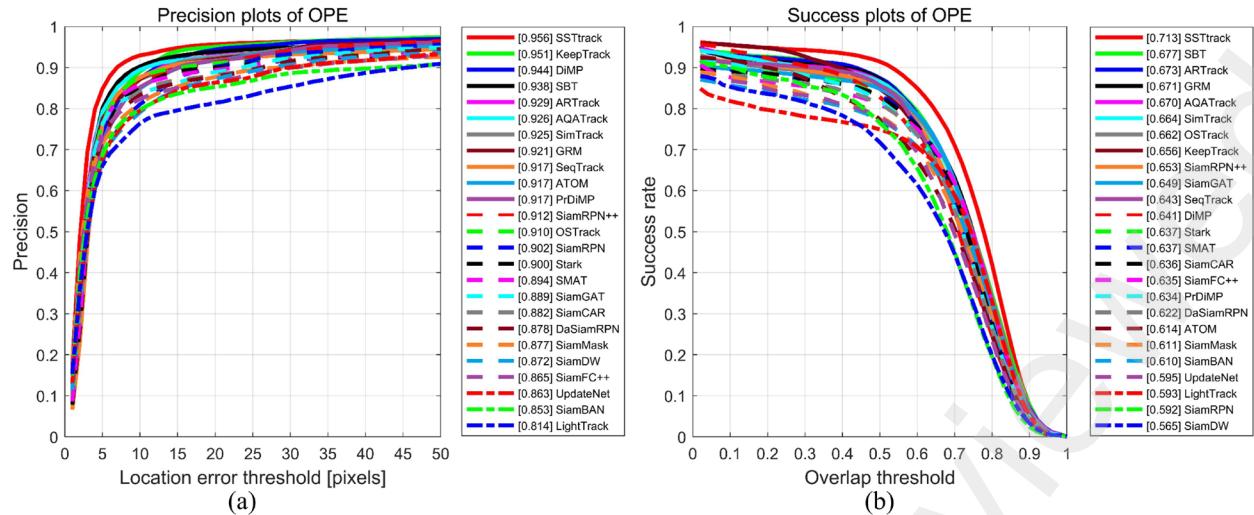
16 4.4. Evaluations against leading RGB trackers

17 This section evaluates SSTtrack with 34 state-of-the-art (SOTA)
18 RGB trackers, categorized into two groups: (i) trackers using hand-
19 crafted features and (ii) those utilizing deep features.

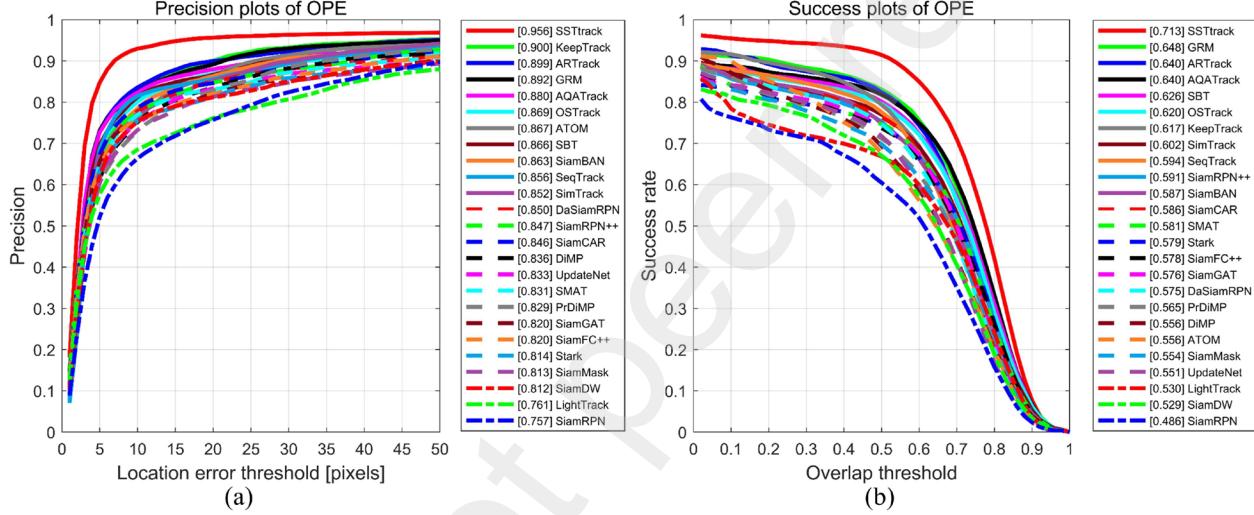
20 4.4.1. Trackers employing hand-crafted features

21 In the benchmark analysis, we compare SSTtrack with 13
22 representative RGB trackers that employ hand-crafted features,
23 namely AutoTrack (Li et al., 2020a), ARCF (Huang et al., 2019),
24 STRCF (Li et al., 2018b), CSRDCF (Lukezic et al., 2018), BACF
47

25 (Galoogahi et al., 2017), DSST (Danelljan et al., 2017), Staple
26 (Bertinetto et al., 2016a), SRDCF (Danelljan et al., 2015), KCF
27 (Henriques et al., 2015), DAT (Possegger et al., 2015), SAMF (Li et
28 al., 2015), CN (Danelljan et al., 2014), and CSK (Henriques et al.,
29 2012). Compared trackers are tested on both RGB and false-color
30 videos, whereas SSTtrack is assessed using HS videos. Table 1
31 provides an overview of the characteristics and results of trackers.
32 Precision and success plots for RGB video trials are illustrated in Fig.
33 5, showcasing SSTtrack's remarkable performance with a Pre of
34 0.956 and a Suc of 0.713. Markedly outperforming STRCF, SSTtrack
35 achieves a 12.7% enhancement in Pre and a 14.4% lift in Suc.
36 Furthermore, it outperforms SRDCF and ARCF by 14.5% and 14.9%
37 in Suc, respectively, underscoring the potential of harnessing the rich
38 spectral-spatial-temporal conditions inherent in HS data.
39 Transforming the HS video into a false color representation facilitates
40 the feasible implementation of the RGB tracker. Subsequently, we
41 proceed to perform experiments utilizing the false color video, as
42 depicted in Fig. 6 and Table 1. It is observed that STRCF
43 demonstrates notable results, with SRDCF, BACF, and ARCF
44 slightly trailing behind. However, SSTtrack stands out as the
45 frontrunner, surpassing them by margins of 14.5%, 15.9%, 16.9%,
46 and 17.1%, respectively.



3 Fig. 7. Comparative results with deep feature-based trackers on RGB videos. (a) Precision plot. (b) Success plot.



5 Fig. 8. Comparative results with deep feature-based trackers on false color videos. (a) Precision plot. (b) Success plot.

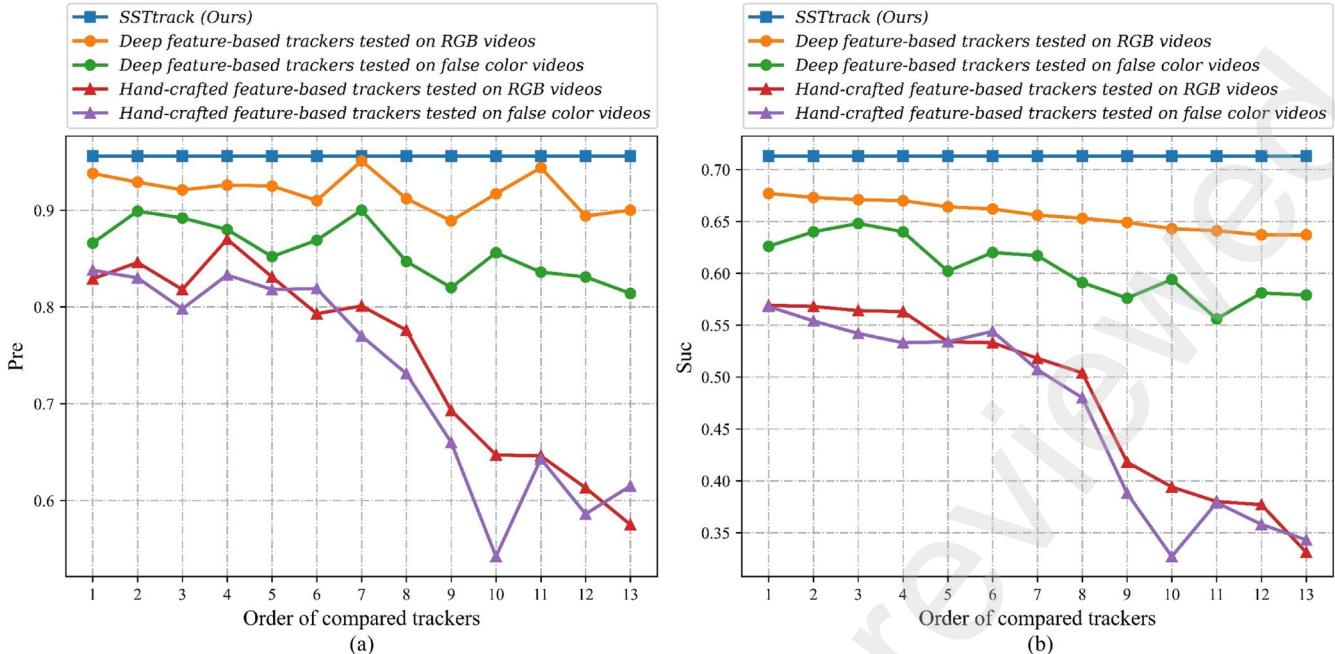
6 4.4.2. Trackers employing deep features

In the realm of RGB tracking, deep features, renowned for their discriminative capabilities, have garnered considerable attention (Javed et al., 2022). Here, SSTtrack undergoes evaluation against 24 SOTA deep feature-based trackers, involving AQATrack (Xie et al., 2024), SMAT (Yelluru Gopal et al., 2024), ARTrack (Wei et al., 2023), SeqTrack (Chen et al., 2023b), GRM (Gao et al., 2023c), SBT (Xie et al., 2022), SimTrack (Chen et al., 2022), OSTRack (Ye et al., 2022), SiamCAR (Cui et al., 2022), Stark (Yan et al., 2021a), LightTrack (Yan et al., 2021b), SiamGAT (Guo et al., 2021), KeepTrack (Mayer et al., 2021), SiamFC++ (Xu et al., 2020), SiamBAN (Chen et al., 2020), PrDiMP (Danelljan et al., 2020), SiamMask (Wang et al., 2019), SiamDW (Zhang et al., 2019b), UpdateNet (Zhang et al., 2019a), SiamRPN++ (Li et al., 2019), DiMP (Bhat et al., 2019), ATOM (Danelljan et al., 2019), DaSiamRPN (Zhu et al., 2018), and SiamRPN (Li et al., 2018a). Table 1 showcases a detailed overview of the results and properties. Figs. 7 and 8 depict the precision plot and success plot for experiments conducted on RGB and false color videos, respectively. Our tracker consistently maintains the leading performance across 24 SOTA trackers. Compared to the baseline AQATrack tested on false color videos, our SSTtrack exhibits gains of 7.6% in Pre and 7.3% in Suc, respectively.

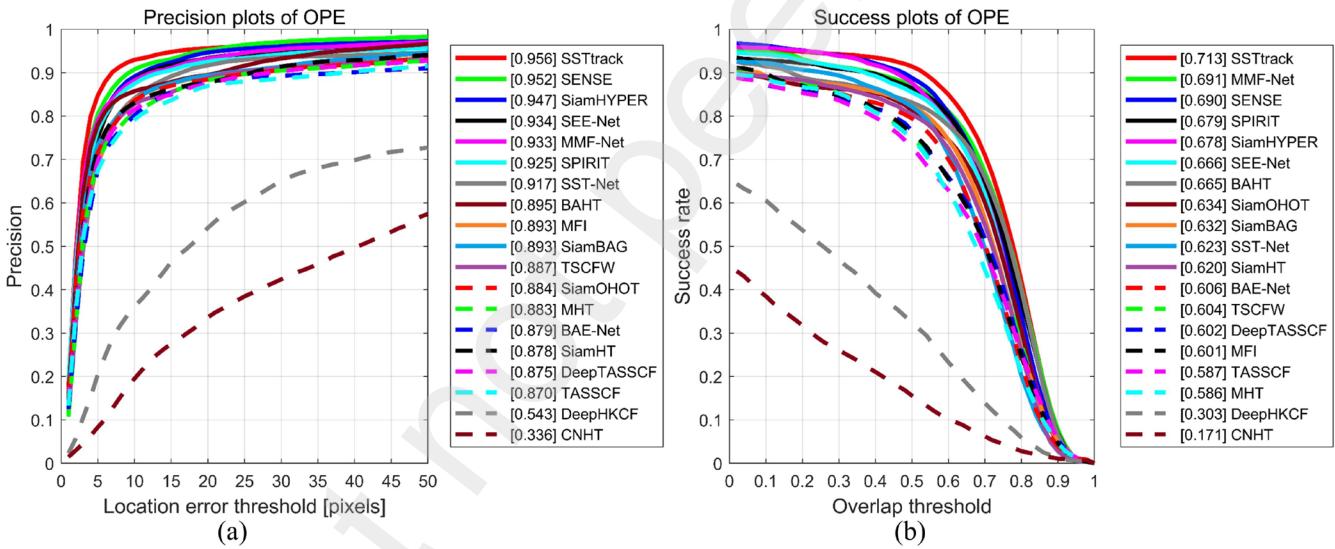
Especially, SSTtrack outperforms top-ranked KeepTrack and SBT by 5.7% and 3.6%, respectively, in terms of Suc on RGB videos, and by 9.6% and 8.7% on false color videos. To summarize, the exhaustive results validate the impressive performance of our SSTtrack.

32 4.5. Parallel analysis with RGB trackers

Fig. 9 depicts the parallel analysis of the top 13 trackers utilizing both hand-crafted and deep features across RGB and false color videos. Through the analysis of Fig. 9, we have gleaned some valuable insights that will enrich our comprehension in the realm of HS video tracking. First, it is noteworthy that trackers employing deep features often outperform those relying on hand-crafted features, regardless of RGB or false color video. This superiority arises from the capability of deep models to acquire generalized representations from training sets. Additionally, the performance of RGB trackers generally diminishes on false color videos compared to RGB videos, albeit with relatively consistent rankings. Third, trackers utilizing deep features demonstrate more pronounced performance degradation, measured in terms of PreD and SucD, compared to their counterparts relying on hand-crafted features. This decline can be attributed to the heavy dependence on properties present in RGB format training datasets. Remarkably, our tracker achieves superior performance by effectively leveraging the spectral-



1
2 **Fig. 9.** Parallel analysis alongside RGB trackers. (a) Pre score. (b) Suc score. The order of compared trackers is determined by their Suc scores on the RGB video.
3 Trackers employing hand-crafted features are listed as follows: STRCF, SRDCF, ARCF, CSRDCF, AutoTrack, BACF, Staple, DSST, SAMF, DAT, CN, KCF,
4 and CSK. Trackers employing deep features are listed as follows: SBT, ARTTrack, GRM, AQATrack, SimTrack, OSTRack, KeepTrack, SiamRPN++, SiamGAT,
5 SeqTrack, DiMP, SMAT, and Stark.



6
7 **Fig. 10.** Comparative results with HS trackers on HS videos. (a) Precision plot. (b) Success plot.

8 spatial-temporal conditions inherent in HS videos. In addition, the
9 transformation of HS video to a false color format unavoidably leads
10 to the loss and distortion of critical spectral-spatial conditions,
11 thereby impacting the stability of trackers.

12 4.6. Comparison with HS trackers

13 Here, we analyze SSTtrack in comparison to 18 representative HS
14 trackers, i.e., MMF-Net (Li et al., 2024), SENSE (Chen et al., 2024a),
15 SPIRiT (Chen et al., 2024b), SiamHT (Tang et al., 2023), SiamBAG
16 (Li et al., 2023a), SiamOHOT (Sun et al., 2023), SEE-Net (Li et al.,
17 2023b), DeepTASSCF (Tang et al., 2022a), TASSCF (Tang et al.,
18 2022a), BAHT (Tang et al., 2022b), TSCFW (Hou et al., 2022),
19 SiamHYPER (Liu et al., 2022), SST-Net (Li et al., 2021), MFI
20 (Zhang et al., 2021), BAE-Net (Li et al., 2020b), MHT (Xiong et al.,
21 2020), DeepHKCF (Uzkent et al., 2019), and CNHT (Qian et al.,

22 2018). Table 2 shows the results and properties of these trackers,
23 while Fig. 10 depicts the precision plot and success plot. Notably, in
24 the precision plot analysis, SENSE, SiamHYPER, SEE-Net, MMF-
25 Net, SPIRiT, and SST-Net stand out as the top performers, achieving
26 Pre scores of 0.952, 0.947, 0.934, 0.933, 0.925, and 0.917,
27 respectively. In contrast, SSTtrack showcases even superior results,
28 achieving a score of 0.956, surpassing them by margins of 0.4%,
29 0.9%, 2.2%, 2.3%, 3.1%, and 3.9%, respectively. Regarding the
30 success plot analysis, MMF-Net, SENSE, SPIRiT, SiamHYPER,
31 SEE-Net, and BAHT show promising performance, with Suc scores
32 of 0.691, 0.690, 0.679, 0.678, 0.666, and 0.665, respectively. Our
33 tracker obtains notable improvements over them by up to 2.2%, 2.3%,
34 3.4%, 3.5%, 4.7%, and 4.8%, respectively. In summary, SSTtrack
35 demonstrates superior performance in both Pre and Suc scores,
36 highlighting the effectiveness of modeling spectral-spatial-temporal
37 conditions in the HS tracking domain.

1 **Table 2**

2 Results and properties of comparative HS trackers tested on HS videos. Trackers are listed in chronological order.

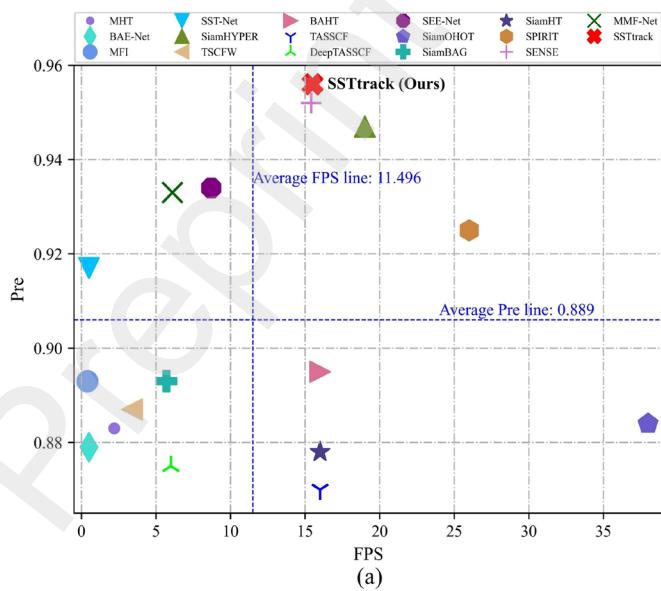
No.	Tracker	Venue	Framework	Feature	UFB	UTC	MOP	FPS	Pre	Suc
1	CNHT (Qian et al., 2018)	ICSM 2018	KCF	Deep feature	✓	-	CPU	2.6	0.336	0.171
2	DeepHKCF (Uzkent et al., 2019)	TGRS 2019	KCF	Deep feature	-	-	CPU	0.9	0.543	0.303
3	MHT (Xiong et al., 2020)	TIP 2019	KCF	Hand-crafted feature	✓	-	CPU	2.2	0.883	0.586
4	BAE-Net (Li et al., 2020b)	ICIP 2020	VITAL	Deep feature	✓	-	GPU	0.5	0.879	0.606
5	MFI (Zhang et al., 2021)	WISP 2021	KCF	Deep + Hand-crafted features	-	-	CPU	0.4	0.893	0.601
6	SST-Net (Li et al., 2021)	WISP 2021	VITAL	Deep feature	✓	✓	GPU	0.5	0.917	0.623
7	SiamHYPER (Liu et al., 2022)	TIP 2022	SiamFC	Deep feature	✓	-	GPU	19.0	0.947	0.678
8	TSCFW (Hou et al., 2022)	TGRS 2022	KCF	Hand-crafted features	✓	-	CPU	3.4	0.887	0.604
9	BAHT (Tang et al., 2022b)	GRSL 2022	SiamFC	Deep feature	-	-	GPU	16.0	0.895	0.665
10	TASSCF (Tang et al., 2022a)	CVIU 2022	KCF	Hand-crafted features	-	-	CPU	16.0	0.870	0.587
11	DeepTASSCF (Tang et al., 2022a)	CVIU 2022	KCF	Deep feature	-	-	CPU	6.0	0.875	0.602
12	SEE-Net (Li et al., 2023b)	TIP 2023	SiamFC	Deep feature	✓	-	GPU	8.7	0.934	0.666
13	SiamOHOT (Sun et al., 2023)	TGRS 2023	SiamFC	Deep feature	✓	-	GPU	38.0	0.884	0.634
14	SiamBAG (Li et al., 2023a)	TGRS 2023	SiamFC	Deep feature	✓	-	GPU	5.7	0.893	0.632
15	SiamHT (Tang et al., 2023)	NCA 2023	SiamFC	Deep feature	-	-	GPU	16.0	0.878	0.620
16	SPIRIT (Chen et al., 2024b)	TGRS 2024	SiamFC	Deep feature	✓	✓	GPU	26.0	0.925	0.679
17	SENSE (Chen et al., 2024a)	IF 2024	SiamFC	Deep feature	✓	✓	GPU	15.4	0.952	0.690
18	MMF-Net (Li et al., 2024)	TGRS 2024	SiamFC	Deep feature	✓	-	GPU	6.1	0.933	0.691
19	SSTtrack	Ours	SiamFC	Deep feature	✓	✓	GPU	15.5	0.956	0.713

3 For the Venue, WISP corresponds to the WHISPERS, and IF refers to the Information Fusion. Regarding the Framework, KCF (Henriques et al., 2015) signifies the kernelized
4 correlation filter. VITAL (Song et al., 2018) involves visual tracking via adversarial learning. SiamFC (Bertinetto et al., 2016b) is a fully convolutional Siamese network.
5 UFB signifies an effort to utilize the full band, while UTC denotes an exploration into leveraging temporal conditions. MOP stands for the main operation platform.

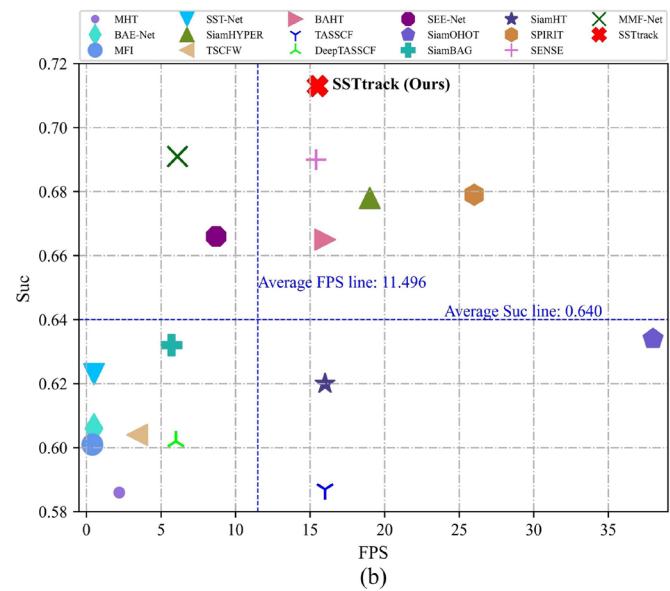
6 4.7. Accuracy vs. speed on HS videos

7 Fig. 11 depicts the trade-off between accuracy (Pre and Suc) and
8 speed (FPS) for various HS trackers. The competitors encompass
9 MHT (Xiong et al., 2020), BAE-Net (Li et al., 2020b), MFI (Zhang
10 et al., 2021), SST-Net (Li et al., 2021), SiamHYPER (Liu et al., 2022),
11 TSCFW (Hou et al., 2022), BAHT (Tang et al., 2022b), TASSCF
12 (Tang et al., 2022a), DeepTASSCF (Tang et al., 2022a), SEE-Net (Li
13 et al., 2023b), SiamOHOT (Sun et al., 2023), SiamBAG (Li et al.,
14 2023a), SiamHT (Tang et al., 2023), SPIRIT (Chen et al., 2024b),
15 SENSE (Chen et al., 2024a), MMF-Net (Li et al., 2024), and our
16 SSTtrack. In Fig. 11(a), it is evident that most HS trackers either

17 exhibit speeds lower than the average FPS value of 11.496 (e.g., SEE-
18 Net, MMF-Net, and SST-Net) or accuracy lower than the average Pre
19 value of 0.889 (e.g., BAHT, SiamOHOT, and SiamHT). Only
20 SSTtrack, SENSE, SiamHYPER, and SPIRIT surpass both
21 the average Pre and FPS levels. Furthermore, our SSTtrack yields the
22 highest Pre score of 0.956. In Fig. 11(b), SSTtrack, SENSE, SPIRIT,
23 SiamHYPER, and BAHT exceed the average Suc of 0.640 and FPS.
24 SSTtrack notably achieves the highest Suc of 0.713, surpassing
25 SENSE and SPIRIT by 2.3% and 3.4%, respectively. In summary,
26 SSTtrack strikes a favorable trade-off between accuracy and speed,
27 making it a solid candidate for HS video tracking.

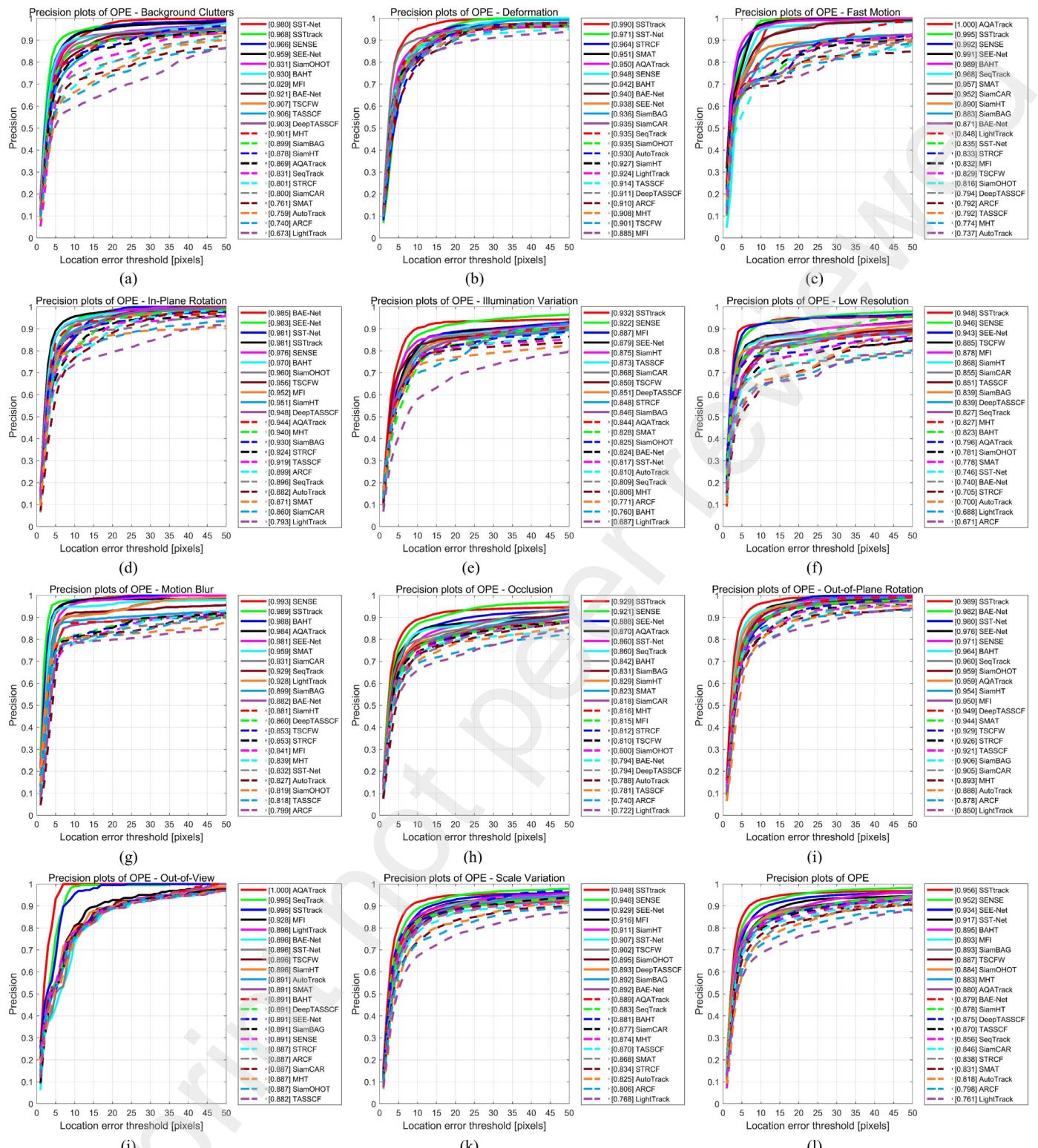


(a)



(b)

28 **Fig. 11.** Comparison of accuracy and speed with HS trackers on HS videos. (a) Pre score vs. FPS value. (b) Suc score vs. FPS value.

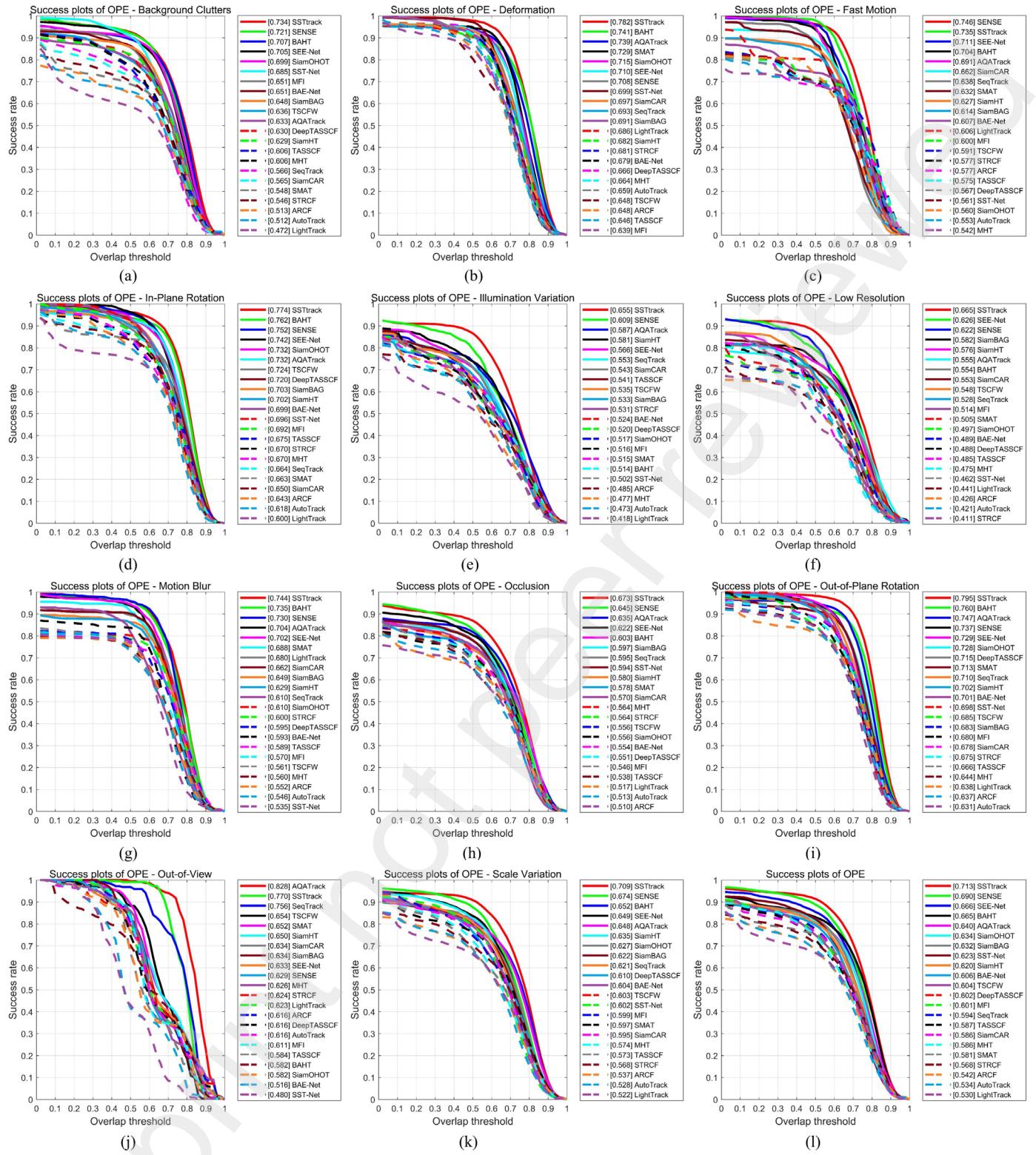


3 **Fig. 12.** Precision plots for each attribute and overall performance. (a) BC. (b) DEF. (c) FM. (d) IPR. (e) IV. (f) LR. (g) MB. (h) OCC. (i) OPR. (j) OV. (k) SV. (l) 4 OVE.

5 4.8. Attribute-based evaluation

6 We further conduct attribute-based comparisons involving eight
7 top-performing RGB trackers: STRCF (Li et al., 2018b), ARCF
8 (Huang et al., 2019), AutoTrack (Li et al., 2020a), LightTrack (Yan
9 et al., 2021b), SiamCAR (Cui et al., 2022), SeqTrack (Chen et al.,
10 2023b), SMAT (Yelluru Gopal et al., 2024), and AQATrack (Xie et
11 al., 2024), alongside 13 leading HS trackers: MHT (Xiong et al.,
12 2020), BAE-Net (Li et al., 2020b), MFI (Zhang et al., 2021), SST-

13 Net (Li et al., 2021), TSCFW (Hou et al., 2022), BAHT (Tang et al.,
14 2022b), TASSCF (Tang et al., 2022a), DeepTASSCF (Tang et al.,
15 2022a), SEE-Net (Li et al., 2023b), SiamOHOT (Sun et al., 2023),
16 SiamBAG (Li et al., 2023a), SiamHT (Tang et al., 2023), and SENSE
17 (Chen et al., 2024a). While the RGB tracker's evaluation is based on
18 false color videos generated from HS videos, others are directly
19 assessed on HS videos. Detailed breakdowns of Pre and Suc scores
20 for each attribute and overall (OVE) are provided in Table 3 and



3 **Fig. 13.** Success plots for each attribute and overall performance. (a) BC. (b) DEF. (c) FM. (d) IPR. (e) IV. (f) LR. (g) MB. (h) OCC. (i) OPR. (j) OV. (k) SV. (l) 4 OVE.

5 Table 4. Additionally, Fig. 12 illustrates the precision plots, while
6 Fig. 13 showcases the success plots. Notably, in terms of the Pre
7 score, as depicted in the detailed results, SSTtrack consistently ranks
8 in the top three across all 11 attributes and secures the first position
9 overall. In particular, SSTtrack holds the first place in six challenging
10 attributes: DEF, IV, LR, OCC, OPR, and SV. Regarding the Suc
11 score, SSTtrack ranks in the first position in nearly all attributes,
12 except for FM and OV. Leveraging efficient and effective

13 components, our tracker exhibits impressive performance in handling
14 challenging attributes, thereby offering a robust benchmark for the
15 HS tracking community.

16 4.9. Visual comparison

17 Fig. 14 displays six visual samples: *basketball*, *coke*, *excavator*,
18 *rider1*, *student*, and *trucker*, tracked by various methods including
19 SST-Net (Li et al., 2021), SENSE (Chen et al., 2024a), SMAT

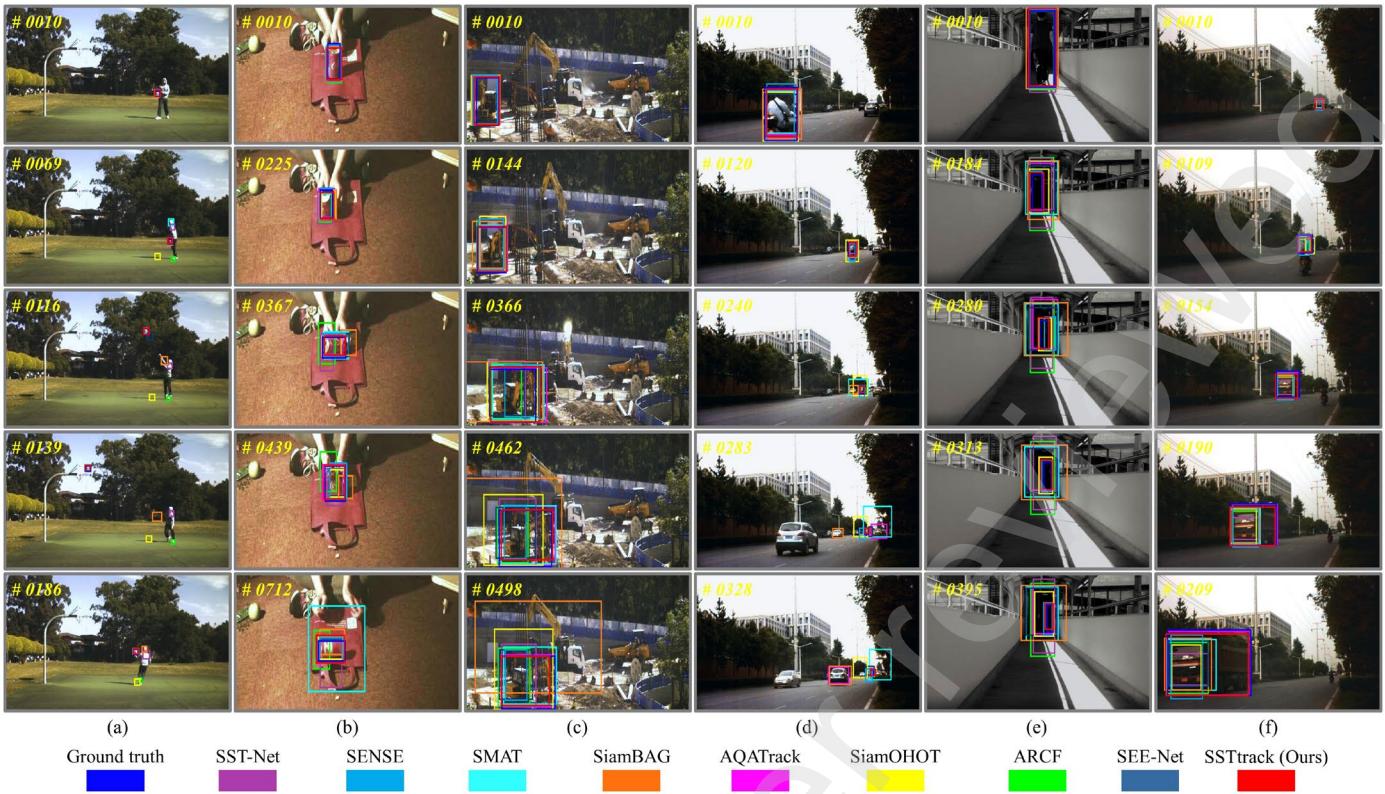


Fig. 14. Qualitative results. (a) *basketball*, attribute: MB, FM, LR, and OCC. (b) *coke*, attribute: IPR, BC, OPR, FM, and SV. (c) *excavator*, attribute: IPR, SV, OPR, OCC, and DEF. (d) *rider1*, attribute: LR, OCC, IV, and SV. (e) *student*, attribute: SV and IV. (f) *truck*, attribute: OV, OCC, IV, and SV. Bounding boxes are visualized in the false color image, with the current frame displayed in the top-left corner.

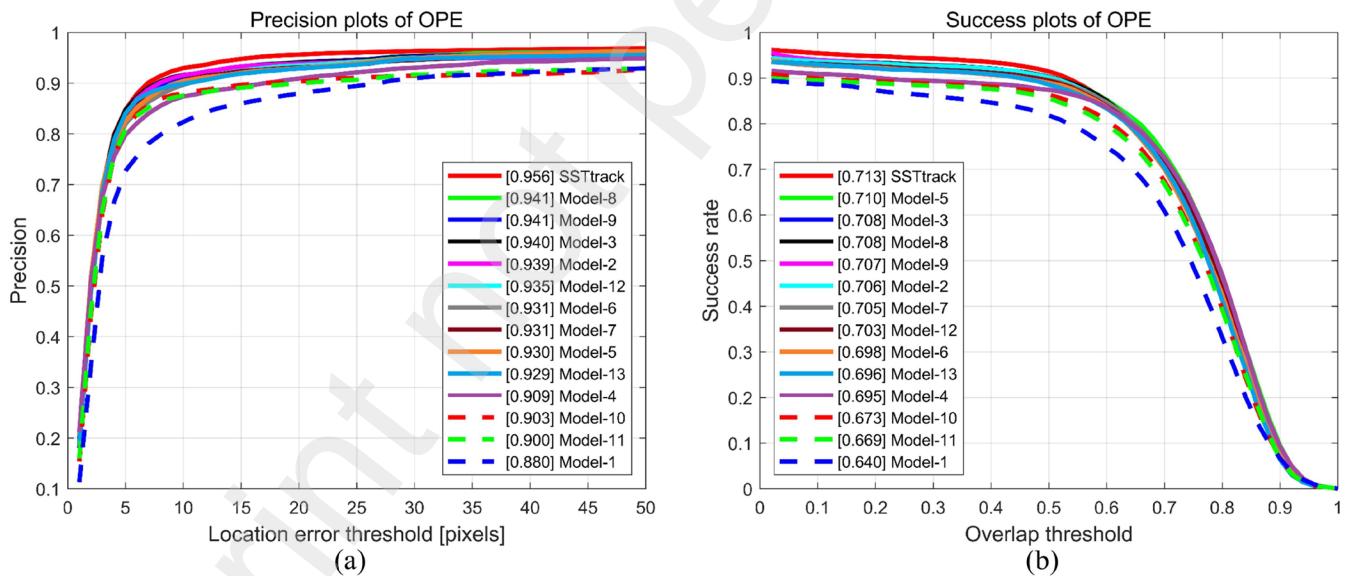


Fig. 15. The overall precision and success plots of ablation studies. (a) Precision plot. (b) Success plot.

(Yelluru Gopal et al., 2024), SiamBAG (Li et al., 2023a), AQATrack (Xie et al., 2024), SiamOHOT (Sun et al., 2023), ARCF (Huang et al., 2019), SEE-Net (Li et al., 2023b), and our SSTtrack. Benefiting from joint modeling of spectral, spatial, and temporal conditions, SSTtrack excels in mitigating interference and accurately determining object position and scale. The qualitative results affirm that SSTtrack maintains promising robustness and accuracy across various challenging scenarios. For instance, in scenarios like *excavator* with attributes IPR, SV, OPR, OCC, and DEF, and *truck* with attributes OV, OCC, IV, and SV, our SSTtrack shows ideal effectiveness in handling these complex situations.

4.10. Ablation studies

As introduced earlier, SSTtrack mainly comprises MGA, SSA, and TCA. To assess their respective contributions, we evaluate SSTtrack alongside 13 ablation models, denoted as Model-1 to Model-13. Notably, the baseline Model-1 is only capable of performing by converting HS video to false-color video, while the other models are evaluated directly on HS video.

The overall precision and success plots of ablation studies are shown in Fig. 15. The ablation results for each component are discussed as follows.

1 **Table 5**

2 Ablation studies on the effect of MGA.

No.	Model	S-MGA	R-MGA	MGA	Params	Pre	Suc	ParamsI	PreI	SucI
1	Model-1	-	-	-	71.980	0.880	0.640	n/a	n/a	n/a
2	Model-2	✓	-	-	<u>72.288</u>	0.939	0.706	<u>0.308</u>	5.9%	6.6%
3	Model-3	-	✓	-	72.188	<u>0.940</u>	<u>0.708</u>	0.208	<u>6.0%</u>	<u>6.8%</u>
4	SSTtrack	-	-	✓	72.555	0.956	0.713	0.575	7.6%	7.3%

3 S-MGA denotes the use of only single-scale feature prompts in the IPEM within MGA. R-MGA denotes the replacement of the IPEM within MGA with the spectral-spatial
4 self-expression module in SENSE. Params denotes the parameter volume (M). ParamsI indicates the improvement in the parameter volume relative to the foundation model
5 (Model-1). The first and second scores are highlighted in bold and underlined, respectively.

6 **Table 6**

7 Ablation studies on the effect of IMPM within SSA

No.	Model	L-IMPM	Params	Pre	Suc	ParamsI	PreI	SucI
1	Model-1	n/a	71.980	0.880	0.640	n/a	n/a	n/a
2	Model-4	null	72.379	0.909	0.695	0.399	2.9%	5.5%
3	Model-5	0	72.396	0.930	<u>0.710</u>	0.416	5.0%	<u>7.0%</u>
4	Model-6	7	72.396	0.931	0.698	0.416	5.1%	5.8%
5	Model-7	0-3	72.449	0.931	0.705	0.469	5.1%	6.5%
6	Model-8	4-7	72.449	<u>0.941</u>	0.708	0.469	<u>6.1%</u>	6.8%
7	Model-9	0-7	72.519	<u>0.941</u>	0.707	0.539	<u>6.1%</u>	6.7%
8	Model-10	0-14	<u>72.642</u>	0.903	0.673	<u>0.662</u>	2.3%	3.3%
9	Model-11	0-19	72.713	0.900	0.669	0.733	2.0%	2.9%
10	SSTtrack	0-9	72.555	0.956	0.713	0.575	7.6%	7.3%

8 L-IMPM indicates the layer index of the used IMPM in the M -layer transformer encoders. null denotes no use of IMPM in any layer, n denotes the use of IMPM at the n-th
9 layer, while d-m indicates that IMPM is used at layer d through layer m.

10 **Table 7**

11 Ablation studies on the effect of TCA.

No.	Model	TCA	V-TCA	Params	Pre	Suc	ParamsI	PreI	SucI
1	Model-1	-	-	71.980	0.880	0.640	n/a	n/a	n/a
2	Model-12	-	-	70.968	<u>0.935</u>	<u>0.703</u>	-1.012	<u>5.5%</u>	<u>6.3%</u>
3	Model-13	-	✓	74.919	0.929	0.696	2.939	4.9%	5.6%
4	SSTtrack	✓	-	<u>72.555</u>	0.956	0.713	<u>0.575</u>	7.6%	7.3%

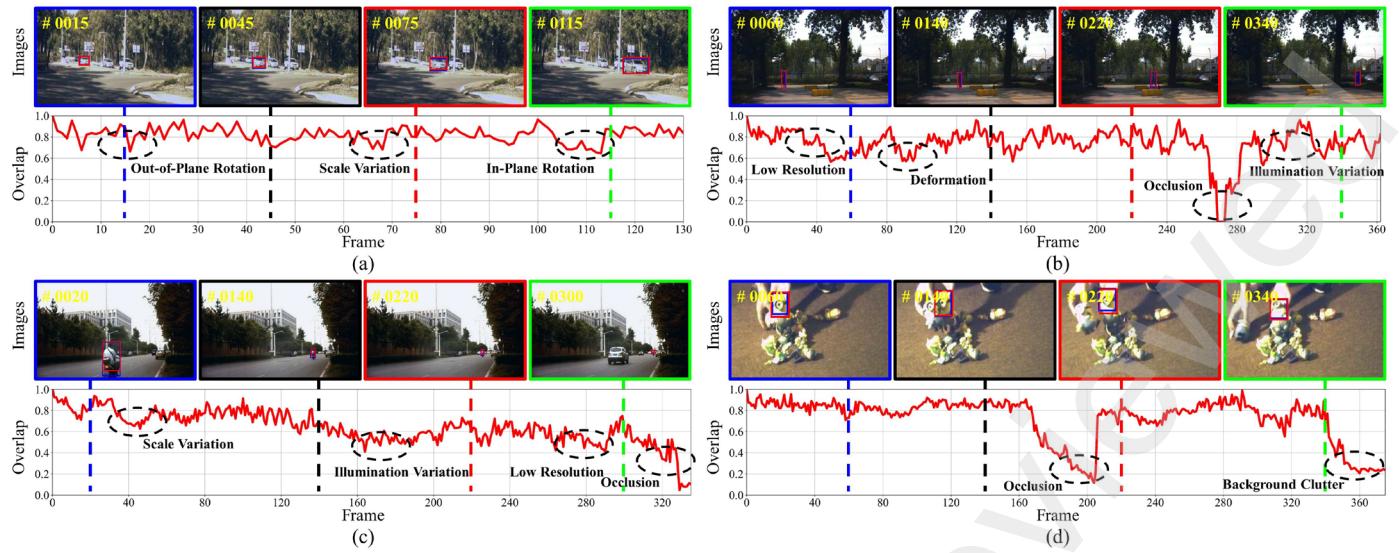
12 V-TCA means using a self-attention layer instead of SSTFM within TCA.

13 4.10.1. Effect of MGA

14 In SSTtrack, MGA is devised to learn the multi-modal generation
15 adaptively, comprising of the IPEM for solving the coefficient matrix
16 and the spectral contribution prompting module for generating band
17 contributions. To validate the effect of MGA, we conduct three
18 ablation experiments. Model-1 is the baseline. Model-2 entails
19 replacing the MGA in SSTtrack with S-MGA (i.e., employing only
20 single-scale feature prompts in the IPEM within MGA). Meanwhile,
21 Model-3 involves replacing the IPEM within MGA with R-MGA
22 (i.e., replacing the IPEM within MGA with the spectral-spatial self-
23 expression module in SENSE). The experimental results are
24 presented in Table 5. Compared to Model-1, SSTtrack brings
25 significant improvements (7.6% in Pre and 7.3% in Suc) with the
26 introduction of only 0.575M parameters. When compared with
27 Model-2 and Model-3, it is evident that our MGA also achieves
28 competitive gains, i.e., 1.7% and 1.6% in the Pre score and 0.7% and
29 0.5% in the Suc score, respectively. These results highlight the
30 effectiveness of MGA, enabling SSTtrack to fully leverage the rich
31 spectral condition in the HS data.

32 4.10.2. Effect of IMPM within SSA

33 To explore a more effective and efficient SSA architecture, we
34 embed our IMPM in a portion of the transformer encoder layer.
35 Specifically, we establish nine embedding types. Table 6 displays
36 experimental results, in which L-IMPM denotes the layer index of the
37 used IMPM in the N -layer (20 layers in total) transformer encoders
38 within HiViT (Zhang et al., 2022a). We observe that as the number
39 of IMPM increases, the model performance steadily improves, as
40 evident in Model-6, Model-7, Model-8, Model-9, and SSTtrack.
41 However, further increases in the number of IMPM lead to a
42 degradation in the magnitude of performance improvement, as
43 observed in Model-10 and Model-11. This phenomenon may be
44 attributed to the increased number of parameters causing the model
45 to lose critical information about the input modality, resulting in
46 inferior prediction results. SSTtrack achieves optimal performance
47 with appropriate learnable parameters (0.575M). Experimental
48 results also indicate the potential for further simplification and
49 enhancement of the SSA architecture.



2 **Fig. 16.** Overlap curves and tracking samples obtained using SSTtrack. (a) *car2*, attribute: SV, IPR, and OPR. (b) *pedestrian2*, attribute: OCC, LR, DEF, and IV. 3 (c) *rider1*, attribute: LR, OCC, IV, and SV. (d) *toy1*, attribute: BC and OCC. Ground truth and prediction are marked in blue and red bounding boxes, respectively.

4 4.10.3. Effect of TCA

5 TCA is devised to inject the temporal condition to guide spectral-
6 spatial feature representations. It comprises the Transformer decoder
7 for temporal condition transfer and the SSTFM for filtering spectral,
8 spatial, and temporal conditions. To demonstrate the role of TCA, the
9 ablation experiments are conducted, as depicted in Table 7. Model-
10 12 represents the elimination of TCA from SSTtrack. Compared with
11 Model-12, SSTtrack obtains consistent improvements by 2.1% in Pre
12 and 1.0% in Suc. Fig. 16 visualizes the overlap curves and tracking
13 results of SSTtrack, showcasing its ability to maintain high overlap
14 scores even with challenging attributes.

15 To further validate the effect of TCA, we use a self-attention layer
16 (Devlin et al., 2018) to replace the SSTFM within TCA for spectral-
17 spatial-temporal condition filtering, yielding the Model-13.
18 Experimental results reveal that SSTtrack outperforms Model-13 by
19 2.7% in Pre and 1.7% in Suc. In comparison to Model-12, it is
20 observed that employing an expensive fusion strategy (Model-13)
21 may not obtain significant improvement, or may even result in a
22 slight degradation. It may be because the complex network introduces
23 a higher-dimensional feature space, making it more challenging for
24 the model to learn and generalize accurately. The SSTFM within
25 TCA implements the filtering and fusion of spectral, spatial, and
26 temporal conditions in an almost parameter-less manner, with results
27 validating its performance.

28 5. Conclusions

29 This study proposes a unified HS video tracking framework
30 (SSTtrack) via modeling spectral-spatial-temporal conditions
31 simultaneously in an end-to-end fashion. First, a multi-modal
32 generation adapter is designed to learn multi-modal generation and
33 bridge the band gap. Following this, a spectral-spatial adapter is
34 designed to dynamically transfer and interact with multiple
35 modalities adaptively. Finally, to capture static and instantaneous
36 object properties, we design a temporal condition adapter to inject the
37 temporal condition for guiding spectral and spatial feature
38 representations. The proposed framework strikes a favorable trade-
39 off between accuracy and speed with the introduction of only 0.575M
40 learnable parameters, making it a reliable candidate for the HS video

41 tracking domain. Nevertheless, the current model struggles to capture
42 longer-term spectral-spatial-temporal conditions that could enhance
43 performance. Future work will focus on importing longer-term
44 conditions in our framework to address the memory restrictions.

45 Acknowledgments

46 Sincerely thanks to the editors and reviewers for their valuable
47 comments and constructive suggestions. Thanks to the related authors
48 for providing codes and results. This work was supported by the
49 National Natural Science Foundation of China (42230108, 423B2104,
50 and 42271411).

51 References

- 52 Bertinetto, L., Valmadre, J., Golodetz, S., Miksik, O., Torr, P.H.S., 2016a. Staple:
53 Complementary Learners for Real-Time Tracking, in Proc. IEEE Conf. Comput. Vis.
54 Pattern Recognit. (CVPR), pp. 1401-1409.
- 55 Bertinetto, L., Valmadre, J., Henriques, J.F., Vedaldi, A., Torr, P.H.S., 2016b. Fully-
56 Convolutional Siamese Networks for Object Tracking, Proc. Eur. Conf. Comput.
57 Vis. (ECCV) /IEEE Trans. Signal Process., pp. 850-865.
- 58 Bhat, G., Danelljan, M., Van Gool, L., Timofte, R., 2020. Know Your Surroundings:
59 Exploiting Scene Information for Object Tracking. Proc. Eur. Conf. Comput. Vis.
60 (ECCV).
- 61 Bhat, G., Danelljan, M., Van Gool, L., Timofte, R., Ieee, 2019. Learning Discriminative
62 Model Prediction for Tracking, Ieee I Conf Comp Vis, pp. 6181-6190.
- 63 Cai, Y., Zhang, Z., Liu, X., Cai, Z., 2020. Efficient Graph Convolutional Self-
64 Representation for Band Selection of Hyperspectral Image. IEEE J. Sel. Top. Appl.
65 Earth Obs. Remote Sens. 13, 4869-4880.
- 66 Cao, B., Guo, J., Zhu, P., Hu, Q., 2024. Bi-directional Adapter for Multi-modal Tracking.
67 in Proc. AAAI Conf. Artif. Intell. (AAAI).
- 68 Cao, Z., Huang, Z., Pan, L., Zhang, S., Liu, Z., Fu, C., 2023. Towards Real-World Visual
69 Tracking with Temporal Contexts. IEEE Trans. Pattern Anal. Mach. Intell.
- 70 Chen, B., Li, P., Bai, L., Qiao, L., Shen, Q., Li, B., Gan, W., Wu, W., Ouyang, W., 2022.
71 Backbone is All Your Need: A Simplified Architecture for Visual Object Tracking.
72 in Proc. Eur. Conf. Comput. Vis. (ECCV), 375-392.
- 73 Chen, L.L., Zhao, Y.Q., Kong, S.G., 2023a. SFA-guided mosaic transformer for tracking
74 small objects in snapshot spectral imaging. ISPRS J. Photogramm. Remote Sens.
75 204, 223-236.
- 76 Chen, X., Peng, H., Wang, D., Lu, H., Hu, H., 2023b. SeqTrack: Sequence to Sequence
77 Learning for Visual Object Tracking, 2023 IEEE/CVF Conference on Computer
78 Vision and Pattern Recognition (CVPR), pp. 14572-14581.
- 79 Chen, Y., Yuan, Q., Tang, Y., Xiao, Y., He, J., Liu, Z., 2024a. SENSE: Hyperspectral
80 video object tracker via fusing material and motion cues. Inf. Fusion 109, 102395.
- 81 Chen, Y.Z., Yuan, Q.Q., Tang, Y.Q., Xiao, Y., He, J., Zhang, L.P., 2024b. SPIRIT:
82 Spectral Awareness Interaction Network With Dynamic Template for Hyperspectral
83 Object Tracking. IEEE Trans. Geosci. Remote Sens. 62.
- 84 Chen, Z.D., Zhong, B.N., Li, G.R., Zhang, S.P., Ji, R.R., Ieee, 2020. Siamese Box
85 Adaptive Network for Visual Tracking, Proc. IEEE Conf. Comput. Vis. Pattern
86 Recognit. (CVPR), pp. 6667-6676.
- 87 Cui, Y., Guo, D.Y., Shao, Y.Y., Wang, Z.H., Shen, C.H., Zhang, L.Y., Chen, S.Y., 2022.
88 Joint Classification and Regression for Visual Tracking with Fully Convolutional
89 Siamese Networks. Int. J. Comput. Vis. 130, 550-566.

SSTtrack: A Unified Hyperspectral Video Tracking Framework via Modeling Spectral-Spatial-Temporal Conditions

- 1 Danelljan, M., Bhat, G., Khan, F.S., Felsberg, M., Soc, I.C., 2019. ATOM: Accurate
 2 Tracking by Overlap Maximization, Proc. IEEE Conf. Comput. Vis. Pattern
 3 Recognit. (CVPR), Long Beach, CA, pp. 4655-4664.
- 4 Danelljan, M., Hager, G., Khan, F.S., Felsberg, M., 2015. Learning Spatially Regularized
 5 Correlation Filters for Visual Tracking, in Proc. IEEE Int. Conf. Comput. Vis.
 6 (ICCV), pp. 4310-4318.
- 7 Danelljan, M., Hager, G., Khan, F.S., Felsberg, M., 2017. Discriminative Scale Space
 8 Tracking, IEEE Trans. Pattern Anal. Mach. Intell. 39, 1561-1575.
- 9 Danelljan, M., Khan, F.S., Felsberg, M., van de Weijer, J., 2014. Adaptive Color
 10 Attributes for Real-Time Visual Tracking, Proc. IEEE Conf. Comput. Vis. Pattern
 11 Recognit. (CVPR), pp. 1090-1097.
- 12 Danelljan, M., Van Gool, L., Timofte, R., Ieee, 2020. Probabilistic Regression for Visual
 13 Tracking, Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Electr
 14 Network, pp. 7181-7190.
- 15 Devlin, J., Chang, M.-W., Lee, K., Toutanova, K., 2018. BERT: Pre-training of Deep
 16 Bidirectional Transformers for Language Understanding. arXiv e-prints,
 17 arXiv:1810.04805.
- 18 Fan, H., Lin, L.T., Yang, F., Chu, P., Deng, G., Yu, S.J., Bai, H.X., Xu, Y., Liao, C.Y.,
 19 Ling, H.B., Soc, I.C., 2019. LaSOT: A High-quality Benchmark for Large-scale
 20 Single Object Tracking, Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR),
 21 pp. 5369-5378.
- 22 Galoogahi, H.K., Fagg, A., Lucey, S., 2017. Learning Background-Aware Correlation
 23 Filters for Visual Tracking, Ieee I Conf Comp Vis, pp. 1144-1152.
- 24 Gao, L., Chen, L., Liu, P., Jiang, Y., Xie, W., Li, Y., 2023a. A Transformer-Based
 25 Network for Hyperspectral Object Tracking, IEEE Trans. Geosci. Remote Sens. 61.
- 26 Gao, L., Liu, P., Jiang, Y., Xie, W., Lei, J., Li, Y., Du, Q., 2023b. CBFF-Net: A New
 27 Framework for Efficient and Accurate Hyperspectral Object Tracking, IEEE Trans.
 28 Geosci. Remote Sens. 61, 1-14.
- 29 Gao, S., Zhou, C., Zhang, J., 2023c. Generalized Relation Modeling for Transformer
 30 Tracking, in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 18686-
 31 18695.
- 32 Guo, D.Y., Shao, Y.Y., Cui, Y., Wang, Z.H., Zhang, L.Y., Shen, C.H., Ieee Comp, S.O.C.,
 33 2021. Graph Attention Tracking, Proc. IEEE Conf. Comput. Vis. Pattern Recognit.
 34 (CVPR), Electr Network, pp. 9538-9547.
- 35 Han, Y., Huang, K., 2024. ACTrack: Adding Spatio-Temporal Condition for Visual
 36 Object Tracking, arXiv e-prints, arXiv:2403.07914.
- 37 Henriques, J.F., Caseiro, R., Martins, P., Batista, J., 2012. Exploiting the Circulant
 38 Structure of Tracking-by-Detection with Kernels, Proc. Eur. Conf. Comput. Vis.
 39 (ECCV), pp. 702-715.
- 40 Henriques, J.F., Caseiro, R., Martins, P., Batista, J., 2015. High-Speed Tracking with
 41 Kernelized Correlation Filters, IEEE Trans. Pattern Anal. Mach. Intell. 37, 583-596.
- 42 Hou, Z., Li, W., Zhou, J., Tao, R., 2022. Spatial-Spectral Weighted and Regularized
 43 Tensor Sparse Correlation Filter for Object Tracking in Hyperspectral Videos, IEEE
 44 Trans. Geosci. Remote Sens. 60.
- 45 Huang, L.H., Zhao, X., Huang, K.Q., 2021. GOT-10k: A Large High-Diversity
 46 Benchmark for Generic Object Tracking in the Wild, IEEE Trans. Pattern Anal.
 47 Mach. Intell. 43, 1562-1577.
- 48 Huang, Z.Y., Fu, C.H., Li, Y.M., Lin, F.L., Lu, P., 2019. Learning Aberrance Repressed
 49 Correlation Filters for Real-Time UAV Tracking, in Proc. IEEE Int. Conf. Comput.
 50 Vis. (ICCV), 2891-2900.
- 51 Javed, S., Danelljan, M., Khan, F.S., Khan, M.H., Felsberg, M., Matas, J., 2022. Visual
 52 Object Tracking with Discriminative Filters and Siamese Networks: A Survey and
 53 Outlook, IEEE Trans. Pattern Anal. Mach. Intell. PP.
- 54 Jiang, X., Wang, X., Sun, C., Zhu, Z., Zhong, Y., 2024. A Channel Adaptive Dual
 55 Siamese Network for Hyperspectral Object Tracking, IEEE Trans. Geosci. Remote
 56 Sens. 62, 1-12.
- 57 Lei, X., Cheng, W., Xu, C., Yang, W., 2024. Joint target and background temporal
 58 propagation for aerial tracking, ISPRS J. Photogramm. Remote Sens. 211, 121-134.
- 59 Li, B., Wu, W., Wang, Q., Zhang, F., Xing, J., Yan, J., 2019. SiamRPN++: Evolution of
 60 Siamese Visual Tracking With Very Deep Networks, in Proc. IEEE Conf. Comput.
 61 Vis. Pattern Recognit. (CVPR), pp. 4277-4286.
- 62 Li, B., Yan, J., Wu, W., Zhu, Z., Hu, X., 2018a. High Performance Visual Tracking with
 63 Siamese Region Proposal Network, in Proc. IEEE Conf. Comput. Vis. Pattern
 64 Recognit. (CVPR), pp. 8971-8980.
- 65 Li, C.L., Xue, W.L., Jia, Y.Q., Qu, Z.C., Luo, B., Tang, J., Sun, D.D., 2022. LasHeR: A
 66 Large-Scale High-Diversity Benchmark for RGBT Tracking, IEEE Trans. Image
 67 Process. 31, 392-404.
- 68 Li, F., Tian, C., Zuo, W., Zhang, L., Yang, M.H., 2018b. Learning Spatial-Temporal
 69 Regularized Correlation Filters for Visual Tracking, Proc. IEEE Conf. Comput. Vis.
 70 Pattern Recognit. (CVPR), pp. 4904-4913.
- 71 Li, W., Hou, Z.F., Zhou, J., Tao, R., 2023a. SiamBAG: Band Attention Grouping-Based
 72 Siamese Object Tracking Network for Hyperspectral Videos, IEEE Trans. Geosci.
 73 Remote Sens. 61.
- 74 Li, Y., Fu, C., Ding, F., Huang, Z., Lu, G., 2020a. AutoTrack: Towards High-
 75 Performance Visual Tracking for UAV with Automatic Spatio-Temporal
 76 Regularization, Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR).
- 77 Li, Y., Zhu, J., 2015. A Scale Adaptive Kernel Correlation Filter Tracker with Feature
 78 Integration, Proc. Eur. Conf. Comput. Vis. (ECCV), pp. 254-265.
- 79 Li, Z., Xiong, F., Zhou, J., Lu, J., Qian, Y., 2023b. Learning a Deep Ensemble Network
 80 With Band Importance for Hyperspectral Object Tracking, IEEE Trans. Image
 81 Process. 32, 2901-2914.
- 82 Li, Z., Xiong, F., Zhou, J., Lu, J., Zhao, Z., Qian, Y., 2024. Material-Guided Multiview
 83 Fusion Network for Hyperspectral Object Tracking, IEEE Trans. Geosci. Remote
 84 Sens. 62, 1-15.
- 85 Li, Z., Xiong, F., Zhou, J., Wang, J., Lu, J., Qian, Y., 2020b. BAE-Net: A Band Attention
 86 Aware Ensemble Network for Hyperspectral Object Tracking, 2020 IEEE
 87 International Conference on Image Processing (ICIP), pp. 2106-2110.
- 88 Li, Z., Ye, X., Xiong, F., Lu, J., Zhou, J., Qian, Y., 2021. Spectral-Spatial-Temporal
 89 Attention Network for Hyperspectral Tracking, 11th Workshop on Hyperspectral
 90 Imaging and Signal Processing: Evolution in Remote Sensing (WHISPERS), pp. 1-
 91 5.
- 92 Lin, L., Fan, H., Zhang, Z., Xu, Y., Ling, H., 2022. SwinTrack: A Simple and Strong
 93 Baseline for Transformer Tracking, in NeurIPS, 16, 743-716,754.
- 94 Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P., 2017. Focal Loss for Dense Object
 95 Detection, 2017 IEEE International Conference on Computer Vision (ICCV), pp.
 96 2999-3007.
- 97 Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollar, P., Zitnick,
 98 C.L., 2014. Microsoft COCO: Common Objects in Context, Proc. Eur. Conf.
 99 Comput. Vis. (ECCV), Zurich, SWITZERLAND, pp. 740-755.
- 100 Liu, P.F., Yuan, W.Z., Fu, J.L., Jiang, Z.B., Hayashi, H., Neubig, G., 2023. Pre-train,
 101 Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural
 102 Language Processing, ACM Computing Surveys 55.
- 103 Liu, Z., Wang, X., Zhong, Y., Shu, M., Sun, C., 2022. SiamHYPER: Learning a
 104 Hyperspectral Object Tracker From an RGB-Based Tracker, IEEE Trans. Image
 105 Process. 31, 7116-7129.
- 106 Loshchilov, I., Hutter, F., 2018. Decoupled Weight Decay Regularization, ICLR, 1-19.
- 107 Lukezic, A., Vojir, T., Zajc, L.C., Matas, J., Kristan, M., 2018. Discriminative
 108 Correlation Filter Tracker with Channel and Spatial Reliability, Int. J. Comput. Vis.
 109 126, 671-688.
- 110 Mayer, C., Danelljan, M., Pani Paudel, D., Van Gool, L., 2021. Learning Target
 111 Candidate Association to Keep Track of What Not to Track, in Proc. IEEE Int. Conf.
 112 Comput. Vis. (ICCV).
- 113 Muller, M., Bibi, A., Giancola, S., Alsbabih, S., Ghanem, B., 2018. TrackingNet: A
 114 Large-Scale Dataset and Benchmark for Object Tracking in the Wild, 15th European
 115 Conference on Computer Vision (ECCV), Munich, GERMANY, pp. 310-327.
- 116 Ouyang, E., Wu, J., Li, B., Zhao, L., Hu, W., 2022. Band Regrouping and Response-
 117 Level Fusion for End-to-End Hyperspectral Object Tracking, IEEE Geosci. Remote
 118 Sens. Lett. 19.
- 119 Possegger, H., Mauthner, T., Bischof, H., Ieee, 2015. In Defense of Color-based Model-
 120 free Tracking, Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 2113-
 121 2120.
- 122 Qian, K., Zhou, J., Xiong, F., Zhou, H., Du, J., 2018. Object Tracking in Hyperspectral
 123 Videos with Convolutional Features and Kernelized Correlation Filter, International
 124 Conference on Smart Multimedia., pp. 308-319.
- 125 Rezatofighi, H., Tsai, N., Gwak, J., Sadeghian, A., Reid, I., Savarese, S., 2019. Generalized Intersection Over Union: A Metric and a Loss for Bounding Box
 126 Regression, 2019 IEEE/CVF Conference on Computer Vision and Pattern
 127 Recognition (CVPR), pp. 658-666.
- 129 Shi, L., Zhong, B., Liang, Q., Li, N., Zhang, S., Li, X., 2024. Explicit Visual Prompts for
 130 Visual Object Tracking, arXiv e-prints, arXiv:2401.03142.
- 131 Song, Y.B., Ma, C., Wu, X.H., Gong, L.J., Bao, L.C., Zuo, W.M., Shen, C.H., Lau,
 132 R.W.H., Yang, M.H., 2018. VITAL: Visual Tracking via Adversarial Learning, Proc.
 133 IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 8990-8999.
- 134 Su, Y., Mei, S., Zhang, G., Wang, Y., He, M., Du, Q., Ieee, 2022. Gaussian Information
 135 Entropy Based Band Reduction for Unsupervised Hyperspectral Video Tracking,
 136 IEEE International Geoscience and Remote Sensing Symposium (IGARSS), pp.
 137 791-794.
- 138 Sun, C., Wang, X., Liu, Z., Wan, Y., Zhang, L., Zhong, Y., 2023. SiamOHOT: A
 139 Lightweight Dual Siamese Network for Onboard Hyperspectral Object Tracking via
 140 Joint Spatial-Spectral Knowledge Distillation, IEEE Trans. Geosci. Remote Sens.,
 141 1-1.
- 142 Tang, Y., Huang, H., Liu, Y., Li, Y., 2023. A Siamese network-based tracking framework
 143 for hyperspectral video, Neural Comput. Appl. 35, 2381-2397.
- 144 Tang, Y., Liu, Y., Huang, H., 2022a. Target-aware and spatial-spectral discriminant
 145 feature joint correlation filters for hyperspectral video object tracking, Comput. Vis.
 146 Image Underst. 223.
- 147 Tang, Y., Liu, Y., Ji, L., Huang, H., 2022b. Robust Hyperspectral Object Tracking by
 148 Exploiting Background-Aware Spectral Information With Band Selection Network,
 149 IEEE Geosci. Remote Sens. Lett. 19.
- 150 Uzkent, B., Rangnekar, A., Hoffman, M.J., 2019. Tracking in Aerial Hyperspectral
 151 Videos Using Deep Kernelized Correlation Filters, IEEE Trans. Geosci. Remote
 152 Sens. 57, 449-461.
- 153 Wang, H., Liu, X., Li, Y., Sun, M., Yuan, D., Liu, J., 2024. Temporal Adaptive RGBT
 154 Tracking with Modality Prompt, arXiv e-prints, arXiv:2401.01244.
- 155 Wang, Q., Zhang, L., Bertinetto, L., Hu, W., Torr, P.H.S., 2019. Fast Online Object
 156 Tracking and Segmentation: A Unifying Approach, in Proc. IEEE Conf. Comput.
 157 Vis. Pattern Recognit. (CVPR), pp. 1328-1338.
- 158 Wang, S., Qian, K., Chen, P., Ieee, 2022. BS-SiamRPN: Hyperspectral Video Tracking
 159 based on Band Selection and the Siamese Region Proposal Network, 12th Workshop
 160 on Hyperspectral Imaging and Signal Processing: Evolution in Remote Sensing
 161 (WHISPERS).
- 162 Wei, X., Bai, Y., Zheng, Y., Shi, D., Gong, Y., 2023. Autoregressive Visual Tracking,
 163 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR),
 164 pp. 9697-9706.
- 165 Wu, Y., Lim, J., Yang, M.H., 2015. Object Tracking Benchmark, IEEE Trans. Pattern
 166 Anal. Mach. Intell. 37, 1834-1848.
- 167 Xie, F., Wang, C., Wang, G., Cao, Y., Yang, W., Zeng, W., 2022. Correlation-Aware
 168 Deep Tracking, in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp.
 169 8741-8750.
- 170 Xie, J., Zhong, B., Mo, Z., Zhang, S., Shi, L., Song, S., Ji, R., 2024. Autoregressive
 171 Queries for Adaptive Tracking with Spatio-Temporal Transformers, in Proc. IEEE
 172 Conf. Comput. Vis. Pattern Recognit. (CVPR).
- 173 Xiong, F., Zhou, J., Qian, Y., 2020. Material Based Object Tracking in Hyperspectral
 174 Videos, IEEE Trans. Image Process. 29, 3719-3733.
- 175 Xu, Y.D., Wang, Z.Y., Li, Z.X., Yuan, Y., Yu, G., Assoc Advancement Artificial, I.,
 176 2020. SiamFC plus plus : Towards Robust and Accurate Visual Tracking with
 177 Target Estimation Guidelines, in Proc. AAAI Conf. Artif. Intell. (AAAI), pp. 12549-
 178 12556.

- 1 Xuefeng Zhu, Tianyang Xu, Zongtao Liu, Zhangyong Tang, Xiaojun Wu, Kittler, J., 2024.
2 UniMod1K: Towards a More Universal Large-Scale Dataset and Benchmark for
3 Multi-modal Learning. Int. J. Comput. Vis., 1-16.
4 Yan, B., Peng, H., Fu, J., Wang, D., Lu, H., 2021a. Learning Spatio-Temporal
5 Transformer for Visual Tracking. in Proc. IEEE Conf. Comput. Vis. Pattern
6 Recognit. (CVPR).
7 Yan, B., Peng, H., Wu, K., Wang, D., Fu, J., Lu, H., 2021b. LightTrack: Finding
8 Lightweight Neural Networks for Object Tracking via One-Shot Architecture
9 Search. Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 15175-15184.
10 Yang, J., Li, Z., Zheng, F., Leonardis, A., Song, J., 2022. Prompting for Multi-Modal
11 Tracking.
12 Ye, B., Chang, H., Ma, B., Shan, S., 2022. Joint Feature Learning and Relation Modeling
13 for Tracking: A One-Stream Framework. in Proc. Eur. Conf. Comput. Vis. (ECCV).
14 Yelluru Gopal, G., Amer, M.A., 2024. Separable Self and Mixed Attention Transformers
15 for Efficient Object Tracking. in Proceedings of the IEEE/CVF Winter Conference
16 on Applications of Computer Vision (WACV), 6708-6717.
17 Zhang, L.C., Gonzalez-Garcia, A., Van De Weijer, J., Danelljan, M., Khan, F.S., Ieee,
18 2019a. Learning the Model Update for Siamese Trackers, Ieee I Conf Comp Vis,
19 Seoul, SOUTH KOREA, pp. 4009-4018.
20 Zhang, X., Tian, Y., Huang, W., Ye, Q., Dai, Q., Xie, L., Tian, Q., 2022a. HiViT:
21 Hierarchical Vision Transformer Meets Masked Image Modeling. arXiv e-prints,
22 arXiv:2205.14949.
23 Zhang, Y., Li, X., Wang, F., Wei, B., Li, L., Ieee, 2022b. A Fast Hyperspectral Object
24 Tracking Method Based On Channel Selection Strategy. 12th Workshop on
25 Hyperspectral Imaging and Signal Processing: Evolution in Remote Sensing
26 (WHISPERS).
27 Zhang, Z., Qian, K., Du, J., Zhou, H., 2021. Multi-Features Integration Based
28 Hyperspectral Videos Tracker, 2021 11th Workshop on Hyperspectral Imaging and
29 Signal Processing: Evolution in Remote Sensing (WHISPERS), pp. 1-5.
30 Zhang, Z.P., Peng, H.W., Soc, I.C., 2019b. Deeper and Wider Siamese Networks for
31 Real-Time Visual Tracking. Proc. IEEE Conf. Comput. Vis. Pattern Recognit.
32 (CVPR), Long Beach, CA, pp. 4586-4595.
33 Zhao, C., Liu, H., Su, N., Yan, Y., 2022. TFTN: A Transformer-Based Fusion Tracking
34 Framework of Hyperspectral and RGB. IEEE Trans. Geosci. Remote Sens. 60.
35 Zhu, J., Lai, S., Chen, X., Wang, D., Lu, H., 2023a. Visual Prompt Multi-Modal Tracking.
36 in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR).
37 Zhu, X.G., Zhang, H.R., Hu, B., Huang, K.P., Arun, P.V., Jia, X.P., Zhao, D., Wang, Q.,
38 Zhou, H.X., Yang, S.W., 2023b. DSP-Net: A Dynamic Spectral-Spatial Joint
39 Perception Network for Hyperspectral Target Tracking. IEEE Geosci. Remote Sens.
40 Lett. 20.
41 Zhu, Z., Wang, Q., Li, B., Wu, W., Yan, J., Hu, W., 2018. Distractor-Aware Siamese
42 Networks for Visual Object Tracking, Proc. Eur. Conf. Comput. Vis. (ECCV), pp.
43 103-119.