

TEMPLATE-DRIVEN INTERACTION FOR HYPERSPECTRAL VIDEO OBJECT TRACKING

Yuzeng Chen¹, Qiangqiang Yuan¹, Member IEEE, Yi Xiao¹

¹School of Geodesy and Geomatics, Wuhan University, Wuhan, China

ABSTRACT

Hyperspectral video data, rich in spectral, spatial, and temporal information, offers significant advantages over traditional tracking paradigms. However, existing hyperspectral trackers face challenges with feature interaction and the underuse of temporal cues. To address these issues, HyperTDI, a template-driven interaction approach, is introduced. It integrates a template-driven interaction module (TDIM) to enhance cross-modal interaction between hyperspectral and false-color features and an adaptive feature fusion module (AFFM) for efficient contextual dependency capture. A temporal decoding network (TDN) is also devised to ensure temporal coherence by incorporating historical states. Experimental results on nine hyperspectral datasets show that HyperTDI outperforms 21 state-of-the-art trackers, proving its effectiveness for hyperspectral video object tracking.

Index Terms—Template-driven interaction, hyperspectral video data, object tracking

1. INTRODUCTION

Hyperspectral video imaging enables the simultaneous capture of spatial, temporal, and spectral information, enhancing material-based object tracking [1]. However, exploiting the rich multidimensional information remains challenging due to band gaps between hyperspectral and conventional modalities, complicating data integration [2]. Some methods convert hyperspectral data to a false-color modality, causing spectral distortion, while others fragment spectral reflectance by generating multiple false-color modalities. Parallel processing of multiple modalities also affects efficiency and accuracy during inference [3]. Recent approaches [4, 5] seek to integrate hyperspectral and false-color modalities, but limited inter-modality feature interaction hinders complementary information extraction. Additionally, the underutilization of temporal information restricts dynamic object modeling. While some works [1, 6] address temporal information, splitting hyperspectral data increases computational overhead and complicates feature interaction. These challenges highlight the need for efficient strategies to interact with multidimensional information and extract temporal cues for enhanced performance.

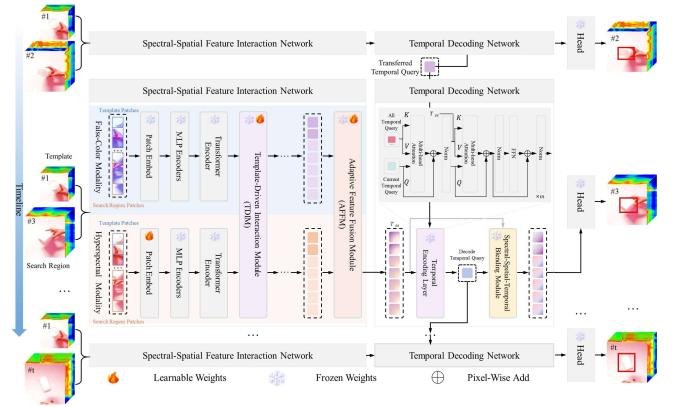


Figure 1. The architecture of the HyperTDI tracker.

Based on prior analysis, we introduce HyperTDI, a template-driven interaction method for hyperspectral video tracking (Figure 1) that enhances multidimensional information synergy. The tracker integrates spectral, spatial, and temporal features to improve tracking performance through inter-modality interactions and temporal correlations. HyperTDI uses a template-driven interaction module (TDIM) to extract complementary features between hyperspectral and false-color data for efficient cross-modal interaction. An adaptive feature fusion module (AFFM), utilizing a HiViT [7] block with trainable adapters, captures contextual dependencies, while a temporal decoding network (TDN) ensures temporal coherence by incorporating historical states. Evaluations on hyperspectral benchmarks, using consistent configurations and weights, show that HyperTDI outperforms 21 advanced trackers.

2. METHODOLOGY

2.1. Overview of the HyperTDI

Figure 1 depicts the HyperTDI method, where the hyperspectral modality is converted into a false-color modality and jointly processed with the original hyperspectral data. The search regions and template images are split into tokens and passed through the shared HiViT backbone with MLP and Transformer encoders. The TDIM enables cross-modal interaction by transferring multimodal

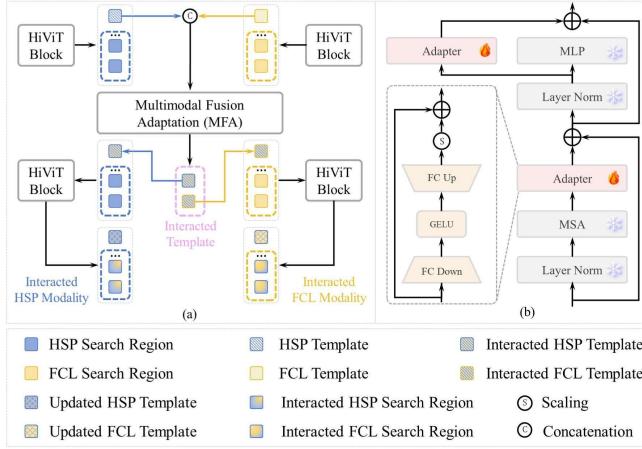


Figure 2. Architecture of the (a) Template-driven interaction module and (b) Multimodal fusion adaptation. HSP: Hyperspectral. FCL: False-color.

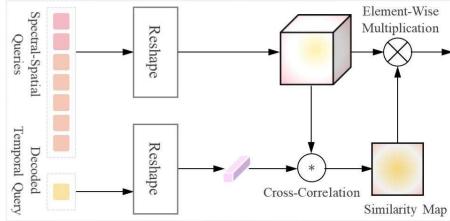


Figure 3. The architecture of the spectral-spatial-temporal blending module.

context from template tokens to search region tokens via self-attention. After the backbone, the AFFM fuses hyperspectral and false-color tokens. The TDN integrates spectral-spatial features from the AFFM with learnable autoregressive queries from the previous frame, combining static and dynamic appearance changes before the head network.

2.2. Template-driven interaction module (TDIM)

The TDIM, shown in Figure 2, enables cross-modal interaction between hyperspectral and false-color modalities by enhancing object features through mutual complementation. It consists of template enhancement and search interaction. The TDIM, shown in Figure 2, enables cross-modal interaction between hyperspectral and false-color modalities by enhancing object features through mutual complementation. It consists of template enhancement and search interaction.

For template enhancement, a multimodal template sequence is formed by concatenating template tokens, which combine object cues from both modalities, facilitating context mining. The multimodal fusion adaptation (Figure 2) utilizes pre-trained HiViT blocks with adapters to capture contextual dependencies. In the search interaction phase, the search region and refined template tokens are entered into the next block. Through iterative TDIM across consecutive blocks, object areas in both modalities are progressively interacted with and enhanced by complementary context.

2.3. Adaptive feature fusion module (AFFM)

Building on ViT's ability to capture contextual dependencies, we propose the AFFM, which reuses the final HiViT blocks

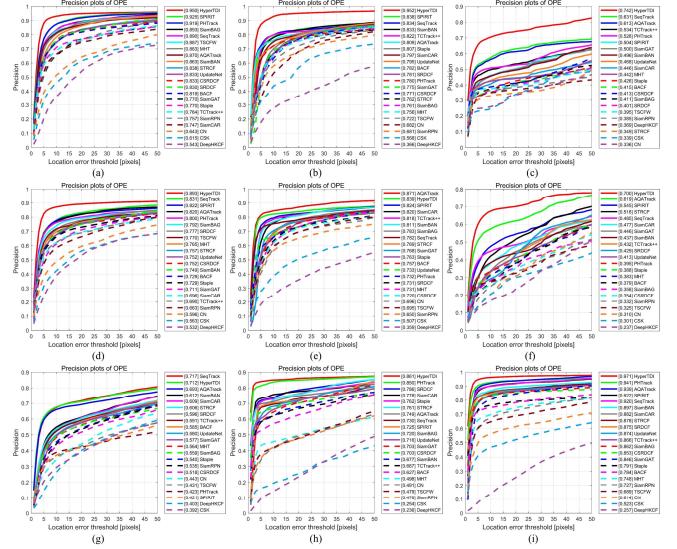


Figure 4. Subfigures (a) through (i) represent precision plots for HOTC20, NIR23, RedNIR23, VIS23, NIR24, RedNIR24, VIS24, MSSOT, and MSVT hyperspectral video tracking benchmarks, respectively.

and introduces trainable adapters for efficient performance with fewer tunable parameters. The fused token, integrating spatial and spectral features, is then fed into the TDN to incorporate temporal cues for enhanced tracking.

2.4. Temporal decoding network (TDN)

The TDN, shown in Figure 1, integrates temporal context with spectral-spatial features through a temporal encoding layer and a spectral-spatial-temporal blending module (Figure 3), which filters relevant information as in [6]. The processed output is then passed to the head network.

3. EXPERIMENT AND DISCUSSION

3.1. Dataset and criteria

The proposed method is evaluated on nine hyperspectral tracking benchmarks: HOTC20 [8], NIR23, RedNIR23, VIS23, NIR24, RedNIR24, VIS24, MSVT [9], and MSSOT [10]. HOTC20 is the 2020 Hyperspectral Object Tracking Challenge dataset, while NIR23, RedNIR23, and VIS23 refer to the 2023 dataset (<https://www.hsitracking.com/>), and NIR24, RedNIR24, and VIS24 is the 2024 dataset. Performance is assessed based on two metrics: AUC from the success plot and DP@20 from the precision plot [11].

3.2. Implementation detail

The model is pre-trained on RGB datasets for 150 epochs, followed by 15 epochs of prompt-tuning on the hyperspectral datasets. The sizes of the search region and template are configured to 256 pixels and 128 pixels, respectively, with a batch size of 32. HyperTDI is implemented using PyTorch and trained on an NVIDIA 3090 GPU.

