# Dynamic Updates for Language Adaptation in Visual-Language Tracking

Xiaohai Li[1], Bineng Zhong[1,*], Qihua Liang[1,†], Zhiyi Mo[2], Jian Nong[2], Shuxiang Song[1]

[1]Key Laboratory of Education Blockchain and Intelligent Technology,
Ministry of Education,Guangxi Normal University, Guilin 541004, China

[2]Guangxi Colleges and Universities Key Laboratory of Intelligent Software,
Wuzhou University, Wuzhou 543002, China

bruc_0619@stu.gxnu.edu.cn, bnzhong@gxnu.edu.cn, qhliang@gxnu.edu.cn

zhiyim@gxuwz.edu.cn, nongjian@gxuwz.edu.cn,songshuxiang@mailbox.gxnu.edu.cn

## Abstract

*The consistency between the semantic information provided by the multi-modal reference and the tracked object is crucial for visual-language (VL) tracking. However, existing VL tracking frameworks rely on static multi-modal references to locate dynamic objects, which can lead to semantic discrepancies and reduce the robustness of the tracker. To address this issue, we propose a novel vision-language tracking framework, named DUTrack, which captures the latest state of the target by dynamically updating multi-modal references to maintain consistency. Specifically, we introduce a Dynamic Language Update Module, which leverages a large language model to generate dynamic language descriptions for the object based on visual features and object category information. Then, we design a Dynamic Template Capture Module, which captures the regions in the image that highly match the dynamic language descriptions. Furthermore, to ensure the efficiency of description generation, we design an update strategy that assesses changes in target displacement, scale, and other factors to decide on updates. Finally, the dynamic template and language descriptions that record the latest state of the target are used to update the multi-modal references, providing more accurate reference information for subsequent inference and enhancing the robustness of the tracker. DUTrack achieves new state-of-the-art performance on four mainstream vision-language and two vision-only tracking benchmarks, including LaSOT, LaSOT_ext, TNL2K, OTB99-Lang, GOT-10K, and UAV123. Code and models are available at https://github.com/GXNU-ZhongLab/DUTrack.*
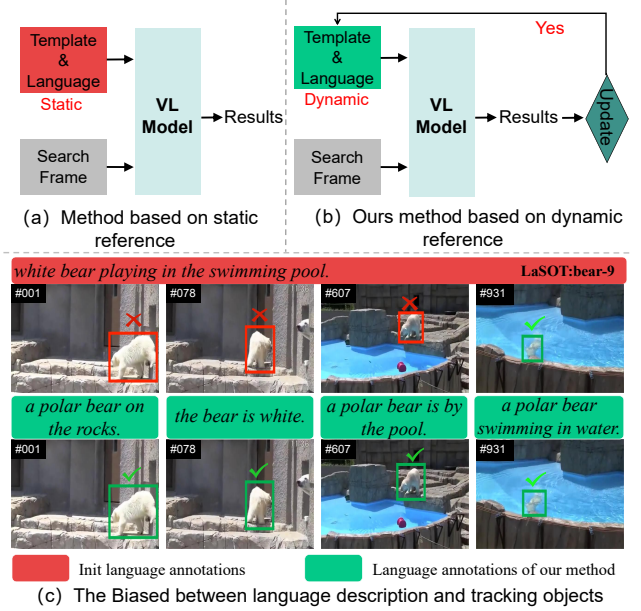
Figure 1. Comparison of different VL Tracking. (a) This vision-language tracking framework [16, 37] relies on static multi-modal references. (b) Our proposes VL framework with dynamically updating multi-modal references. (c) Compare the semantic bias between static annotations and those generated by our method.

## 1. Introduction

Visual-language (VL) tracking aims to effectively track objects in video sequences based on multi-modal reference information provided by natural language descriptions and template frames. In real-world scenarios, objects may undergo appearance changes. These changes can cause discrepancies between the multi-modal reference information and the target, making the tracking task more complex. Therefore, maintaining the consistency between multi-modal target reference information and the target's

*Corresponding Author
†Corresponding Author

1

state in the tracking scene during the tracking process has become a key challenge in this field.

Despite this challenge, a large number of current VL tracking methods [10, 16, 27, 44, 46] tend to overlook it, instead relying on static multi-modal references, as shown in Fig. 1(a), and focusing on establishing stronger multi-modal interaction mechanisms. These methods can be categorized into one-stream and two-stream VL tracking frameworks. The two-stream framework [16, 27, 44, 46] follows a typical three-stage paradigm: (*i*) first, two independent encoders are used to extract uni-modal features from the initial language description and image, respectively; (*ii*) then, multi-modal joint learning is employed to transfer the object information from the static multi-modal references to the search frame; (*iii*) finally, a bounding box prediction head is used to output the tracking result. In contrast, the one-stream framework [14, 37–39] differs in that it uses a unified encoder to simultaneously perform both multi-modal feature extraction and interaction, thereby simplifying the feature processing pipeline. These methods focus on designing effective multi-modal interaction mechanisms. Although they have achieved good performance, there is still a gap compared to the best vision-only trackers [33, 45]. We believe the key reason for this gap lies in *their excessive reliance on static multi-modal references composed of the initial template frame and language annotation*. As shown in Fig. 1(c), the initial language description can only provide the object's state at a specific moment and cannot continuously reflect the object's dynamic changes throughout the video sequence. Therefore, static multi-modal reference information easily diverges from the target's actual state, leading to lag or distortion in information during long-term tracking, which hampers accurate continuous localization and recognition of the object.

To address the above issues, we propose a novel VL tracking framework from the perspective of dynamically updating multi-modal reference information, as shown in Fig. 1(b). Our method can effectively reduce semantic discrepancies between natural language descriptions and the actual state of the target. Specifically, we design a Dynamic Template Capture Module (DTCM) and a Dynamic Language Update Module (DLUM) to update the visual and language references. The DTCM selects the top-k patches with the highest attention scores from the search image based on the attention map guided by the language annotations, and these patches represent the latest visual features of the object. The DLUM is based on a large language model. It generates updated language annotations using the search image and object category information. To improve update efficiency, we also design a strategy that adjusts the update frequency based on changes in the object's position, scale, and other factors. Extensive detailed experiments have demonstrated that our method, which combines

dynamic visual and language information to update multi-modal references, can effectively enhance the performance of VL trackers. The major contributions of our work are summarized as follows:

- We propose DUTrack, which enhances tracking capability by dynamically updating multi-modal references.
- We introduce two dynamic update modules, DTCM and DLUM, which can update visual and language reference.
- DUTrack sets a new state-of-the-art on four visual-language benchmarks and remains competitive on two vision-only benchmarks.

## 2. Related Work

**Vision Tracking.** The vision trackers initialize their tracking procedure based on the given bounding box in the first frame. Based on whether the reference information is updated during tracking, trackers can be divided into two categories: static-reference trackers [18–20, 33, 36] and dynamic-reference trackers [30, 34, 40, 41, 45]. Static-reference trackers rely heavily on the object's initial appearance features, focusing on modeling the relationship between the template frame and the search frame features. SiamFC [3] extracts features from both the search frame and the template frame through a weight-sharing network and then uses a cross-correlation operation to propagate the reference information to the search frame features. OS-Track [36] integrates the feature extraction and feature fusion stages, significantly improving the performance of the tracker. The advantage of these static-reference trackers is their simplicity and relatively low computational cost, but their drawback is the difficulty in adapting to changes in the object's appearance. In contrast, dynamic-reference trackers continuously update the reference frame or object model over time to adapt to changes in the object's appearance. [40] proposes an adaptive template update algorithm that updates the tracker's template based on changes in the object's appearance. STARK [34] introduces dynamic reference updates and an update controller to capture changes in the object's appearance. In this paper, we aim to propose a visual-language tracker that can automatically update multi-modal reference information, reducing the tracker's dependence on the initial language annotations.

**Vision-language Tracking.** In contrast to traditional vision tracking tasks, vision language trackers not only utilize RGB reference information but also incorporate natural language descriptions as additional reference input. The primary objective is to integrate both visual and linguistic modalities to achieve more accurate object tracking. Current vision language trackers can be categorized into one-stream [14, 37–39] and two-stream [13, 16, 44, 46] frameworks based on how they process multi-modal information. Two-stream visual-language tracking frameworks typically use two different models to extract visual and lan-
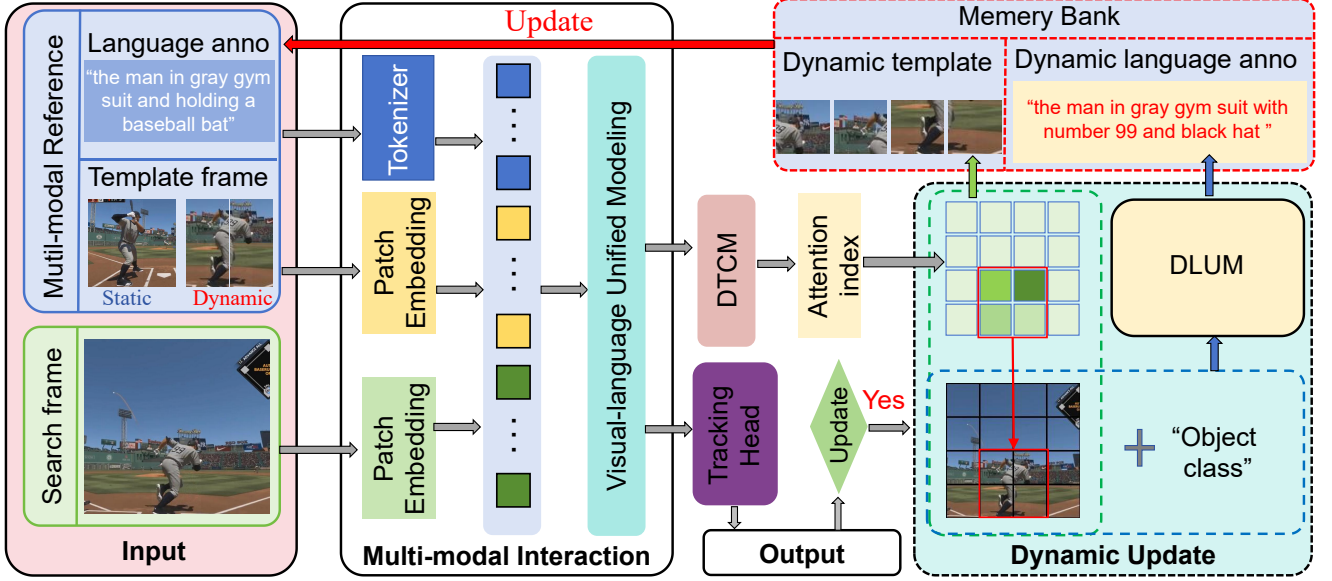
Figure 2. **Overall framework of the proposed DUTrack.** The input consists of two parts: search frame and multi-modal reference. The image and text information are transformed into tokens through Patch Embedding and Tokenizer processing, respectively. Then, these tokens enter the multi-modal interaction module for unified interaction. The resulting multi-modal features are processed through the tracking head to produce the final output. Based on this result, it is determined whether to update the multi-modal reference. The dynamic multi-modal reference is primarily responsible for generating a new reference according to the object's state in the current frame.

guage features (e.g., using ViT [1] to extract visual features and BERT [7] to extract language features), followed by multi-modal fusion to enable information interaction between the two modalities. UVLTrack [27] extracts features from both modalities separately at a shallow level and then performs feature fusion at a deeper level. The VLT [16] introduces a modality mixer during the independent extraction of different modality features, effectively promoting multi-modal interaction. In contrast, one-stream trackers use only one encoder to simultaneously extract features from both modalities and perform multi-modal feature fusion in a single stage. To enhance the object perception of visual-language trackers in complex scenes, Zhang [38] proposed an integrated framework that employs a unified transformer backbone to jointly learn feature extraction and interaction. ATTracker [14] uses an asymmetric multi-encoder to unify the learning of multi-modal features.

Although the aforementioned visual-language trackers have achieved significant success, they still have shortcomings: they rely on static multi-modal references. These multi-modal references consist of an initial template frame and language annotation. The initial language annotation typically provides an overview of the object's complete actions or describes its state at a specific moment. Similarly, the initial template frame describes the target's visual state at the beginning. Relying solely on static references can lead to discrepancies between the references and the target's state. In contrast, we propose a visual-language tracking

framework that can dynamically update multi-modal references to maintain consistency between the references and the object.

## 3. Method

In this section, we will provide a detailed introduction to the proposed DUTrack. First, we review the tracking process of DUTrack. Then, we introduce the multi-modal interaction module in detail. Finally, we delve into the dynamic template capture module and the dynamic language update module.

### 3.1. Overview

An overview of the DUTrack framework is shown in Fig. 2. This framework is very concise and consists of four main components: the multi-modal interaction module, the dynamic update module, the dynamic change capture module, and the tracking head. The input to this pipeline primarily consists of two parts: the search image $S \in \mathbb{R}^{3 \times H_S \times W_S}$ and the multi-modal reference, where the multi-modal reference is composed of the template frame $T \in \mathbb{R}^{3 \times H_T \times W_T}$ and language annotations $L$. $H$ and $W$ represent the height and width of the image. After the input enters the multi-modal interaction module, $S$ and $T$ are transformed into tokens through image patch embedding operations. $L$ is converted into a string of tokens through a tokenizer. Then, we concatenate these tokens and input them into the visual-

3

language unified modeling for multi-modal interaction. In the final stage of interaction, the module outputs the search feature map and the global attention map. The global attention map, processed through the dynamic template capture module, can index the tokens in $S$ that have a high match with $L$. The search feature map is input into the tracking head to obtain results. By analyzing changes in the object's position, scale, and other attributes, the search image and object category information are fed into the dynamic language update module to generate dynamic language annotations. Finally, the dynamic template and language annotations update the multi-modal reference.

## 3.2. Multi-modal Interaction

Existing visual tracking frameworks typically use two separate encoders to independently extract features from different modalities before fusing them. This approach results in a lack of deep correlations between the extracted features, leading to sub-optimal performance in complex scenarios. Inspired by recent advancements in joint feature learning and relation modeling, we adopt a one-stream framework for multi-modal feature learning. To capture rich spatial information, we use HiViT [42] for unified visual-language modeling. Unlike the vanilla ViT [1], which directly uses $16 \times 16$ embeddings, the search image and template image are transformed into tokens $S_t \in \mathbb{R}^{N_S \times D}$ and $T_t \in \mathbb{R}^{N_T \times D}$ through three stages of down-sampling ( a $4 \times 4$ embedding layer and two $2 \times 2$ merging layers). The language annotation $L$ is converted into $L_t \in \mathbb{R}^{N_L \times D}$ through BERT's [7] tokenizer, and $L_t$ begins with a $[CLS]$ token. Here $N$ represents the number of tokens, and $D$ represents the dimensionality, $N_S = H_S W_S / 16^2$, $N_T = H_T W_T / 16^2$, $N_L = 16$, $D = 512$. It is worth noting that $T_t$ includes not only the initial template patches but also the patches from the dynamic template. Formally, the operation of unified visual-language modeling can be expressed as:

$$
\begin{aligned}
Q = K = V &= [L_t; T_t; S_t], \\
feat, attn &= MHSA(Q, K, V), \\
f_{vl} &= [L_t; T_t; S_t] + LN(\lambda_1 \cdot feat), \\
[F_L; F_T; F_S] &= f_{vl} + LN(\lambda_2 \cdot MLP(f_{vl})),
\end{aligned}
\tag{1}
$$

where $LN$ represents the layer normalization, $MHSA$ stands for multi-head self attention and $[;]$ denotes the concatenation operation. $\lambda_1$ and $\lambda_2$ are two learnable parameters. Finally, the search feature $F_S$ is fed into the tracking head to obtain the result $[x, y, w, h]$.

## 3.3. Dynamic Template Capture Module

The object appearance information provided by a static template frame only captures the object's appearance at a specific moment, which is insufficient to offer the tracker
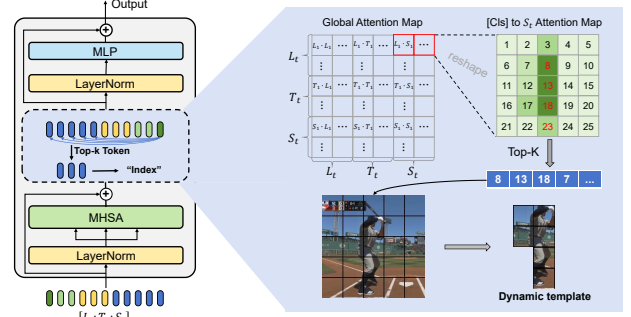


Figure 3. Illustration of the process of capturing dynamic templates. The left side shows the unified visual-language modeling generating a global attention map, while the right side captures dynamic templates based on the global attention map.

adequate spatial cues. To address this issue, we designed a simple yet efficient dynamic template capture module. This module extracts high-response patches from the search image and uses them as dynamic templates for the subsequent frame. Specifically, the multi-head self-attention operation can be seen as spatial aggregation of tokens with normalized importance, as shown in Fig. 3, This is measured by the dot product similarity between each pair of tokens. The calculation for each token is as follows:

$$
A = Softmax(\frac{QK^T}{\sqrt{d}}) \cdot V,
\tag{2}
$$

Where $A$ is the similarity matrix between tokens, based on Eq. 1, Eq. 2 can be extended as:

$$
A = Softmax(\frac{[Q_L; Q_T; Q_S][K_L; K_T; K_S]^T}{\sqrt{d}}) \cdot [V_L; V_T; V_S],
\tag{3}
$$

Where subscripts $L$, $T$, and $S$ denote matrix items belonging to language annotations, templates, and search regions. The module we designed aims to identify patches in the search area that highly match the language annotations. The $[CLS]$ token in the language annotations can comprehensively summarize the semantic information. Therefore, we select the attention map of $[CLS]$ towards the search area from $A$, which can be represented as follows:

$$
A_{l2s} = Softmax(\frac{[Q_{CLS}][K_S]^T}{\sqrt{d}}) \cdot [V_L; V_T; V_S].
\tag{4}
$$

Then, we select the top-k patches with the highest similarity from $A_{l2s}$ and record their indices. Based on these indices, we locate the corresponding regions in the image to serve as dynamic templates.

## 3.4. Dynamic Language Update Module

The official annotations provided by vision-language tracking benchmarks are limited to describing the short-

term state of the object, which poses challenges to frame-by-frame reasoning in tracking paradigms. To enhance the contribution of language annotations to tracking performance, we leverage large language models to dynamically generate descriptions of the target at specific time points. Typically, these annotations include semantic information about the object's position, scale, and color. Based on this, we designed a dynamic update strategy. Specifically, we introduce an object stamp $r_{stamp} : [x_1, y_1, w_1, h_1]$ that records target information from the last update, initialized with annotations from the target's first frame. The tracking result for the current frame is also represented as $r_i : [x_2, y_2, w_2, h_2]$. By comparing the displacement of the center point, changes in the target's scale, and alterations in the mean color within the bounding box between $r_{stamp}$ and $r_i$, we dynamically determine whether the language annotations should be updated. This process can be described as follows:

$$\Delta S = \frac{w_1 \times h_1}{w_2 \times h_2}, \tag{5}$$

$$\Delta D = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}, \tag{6}$$

$$\Delta C = \sqrt{(R_1 - R_2)^2 + (G_1 - G_2)^2 + (B_1 - B_2)^2}, \tag{7}$$

where $\Delta S$, $\Delta D$, and $\Delta C$ represent the changes in scale, position, and color mean of the current target relative to the object stamp. $R$, $G$, and $B$ represent the mean values of the RGB pixels, respectively. Then, we manually set three thresholds to control the update frequency. Finally, we leverage a large language model to simultaneously learn the relationship between the search image and the object category, effectively generating contextually relevant descriptions associated with the object category.

## 4. Experiments

### 4.1. Implementation Details

**Model.** We use BERT [7]'s tokenizer and the hierarchical patch embedding method from HiViT [42] to convert both language and images into tokens. Then, we initialize our ViT module as a unified vision-language encoder using Fast-ITPN [31]. The one-stream framework can effectively improve the efficiency of multi-modal interaction. As shown in Tab. 1, under the premise of using the same ViT-base as the backbone, our tracker demonstrates significant competitiveness in both speed and AUC performance.

Table 1. Comparison of performance, model parameters, and inference speed on TNL2K [32].

| Tracker | AUC(%)↑ | Params(M)↓ | Speed(fps)↑ | Device |
|---------|---------|------------|-------------|--------|
| JiontNLT [46] | 56.9 | 153.0 | 25.6 | RTX-2080Ti |
| MMTrack [44] | 58.6 | 176.9 | 36.2 | RTX-2080Ti |
| Ours | **64.9** | **69.9** | **43.5** | RTX-2080Ti |

**Training and Inference.** Based on different search frame input sizes, we trained two models, DUTrack-256 and DUTrack-384. The training process is as follows: we first use LaSOT [8], GOT-10K [21], COCO [26], TrackingNet [28], and TNL2k [32] for the first stage training. In this stage, we do not use language information, aiming to develop strong visual tracking capability. We employ AdamW to optimize the network parameters, with both the learning rate and weight decay set to $1 \times 10^{-4}$. During this stage, we train for 150 epochs, with a sample size of $60,000$ images. In the second stage, we use LaSOT [8], GOT-10K [21], and TNL2K [32] as the training benchmarks and introduce a dynamic update multi-modal reference mechanism. Since the training strategy is based on random sampling, to reduce training time, we directly use the language annotations provided by DTLLM-VLT [24] as our input. This stage involves 50 epochs of training, with the same training parameters as in the first stage. During inference, the top-k is set to 3, and the large language model used is BLIP [22].

### 4.2. State-of-the Art Comparisons

In this section, we evaluated DUTrack's performance on four vision-language tracking benchmarks, including TNL2K [32], LaSOT [8], LaSOT_ext [9], and OTB99-Lang [25]. The results are shown in Tab. 2. Furthermore, since DUTrack has the ability to generate language descriptions, it can also be evaluated on vision-only benchmarks that do not provide language annotations. We additionally evaluated it on the vision-only benchmarks GOT-10K [21] and the drone tracking dataset UAV123 [2]. Their results are shown in Tab. 3.

**LaSOT** is a large-scale vision-language tracking benchmark, where the language annotations primarily describe the object's behavior and state throughout the video. It consists of 1,400 video sequences, with the training and testing sets split in a 1,220/280 ratio. As shown in Tab. 2, our proposed DUTrack-256 outperforms the QueryNLT [29], published at CVPR 2024, with improvements of 14%, 14.2%, and 17.6% in AUC, $P_{Norm}$, and P, respectively. DUTrack-384 achieves new state-of-the-art results with an AUC of 74.1%, a $P_{Norm}$ of 84.9%, and a P of 82.9%. Compared to the top vision-only tracker ODTrack-B384, we achieve improvements of 0.9%, 1.7%, and 2.3% on these three metrics. LaSOT is a benchmark with relatively long average sequences, where the mismatch between the language annotations and the frames is more pronounced, which has prevented previous VL trackers from surpassing top vision-only trackers on this benchmark. However, DUTrack's performance demonstrates that our proposed dynamic update mechanism for multi-modal reference information significantly advances vision-language tasks.

**LaSOT_ext** is an extension of the LaSOT dataset, containing 150 sequences. As shown in Tab. 2, we observe that our

Table 2. Performance comparison on four benchmarks, including LaSOT [8], LaSOT$_{ext}$ [9], TNL2K [32], and OTB99-Lang [25]. We compare DUTrack with state-of-the-arts. These works can be mainly divided into visual-only trackers and VL trackers. The top three results are highlighted in red and blue, respectively.

| Type | Method | Source | LaSOT | | | LaSOT$_{ext}$ | | | TNL2K | | | OTB99-Lang | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | AUC | P$_{Norm}$ | P | AUC | P$_{Norm}$ | P | AUC | P$_{Norm}$ | P | AUC | P |
| Vision-only | SiamFC [3] | ECCV2016 | 33.6 | 42.0 | 33.9 | 23.0 | 31.1 | 26.9 | 29.5 | 45.0 | 28.6 | - | - |
| | SiamBAN [6] | CVPR2020 | 51.4 | 59.8 | 52.1 | - | - | - | 41.0 | 48.5 | 41.7 | - | - |
| | TransT [4] | CVPR2021 | 64.9 | 73.8 | 69.0 | - | - | - | 50.7 | 57.1 | 51.7 | - | - |
| | Stark [34] | ICCV2021 | 67.1 | 77.0 | - | - | - | - | - | - | - | - | - |
| | GTELT [47] | CVPR2022 | 67.7 | - | - | 45.0 | 54.2 | 52.2 | - | - | - | - | - |
| | TransInMo [17] | IJCAI2022 | 65.7 | 76.0 | 70.7 | - | - | - | 52.0 | 58.5 | 52.7 | - | - |
| | OSTrack-256 [36] | ECCV2022 | 69.1 | 78.7 | 75.2 | 47.4 | 57.3 | 53.3 | 54.3 | - | - | - | - |
| | OSTrack-384 [36] | ECCV2022 | 71.1 | 81.1 | 77.6 | 50.5 | 61.3 | 57.6 | 55.9 | - | - | - | - |
| | SeqTrack-B256 [5] | CVPR2023 | 69.9 | 79.7 | 76.3 | 49.5 | 60.8 | 56.3 | 54.9 | - | - | - | - |
| | SeqTrack-B384 [5] | CVPR2023 | 71.5 | 81.1 | 77.8 | 50.5 | 61.6 | 57.5 | 56.4 | - | - | - | - |
| | AQATrack-256 [33] | CVPR2024 | 71.4 | 81.9 | 78.6 | 51.2 | 62.2 | 58.9 | 57.8 | - | 59.4 | - | - |
| | AQATrack-384 [33] | CVPR2024 | 72.7 | 82.9 | 80.2 | 52.7 | 64.2 | 60.8 | 59.3 | - | 62.3 | - | - |
| | ODTrack-B384 [45] | AAAI2024 | 73.2 | 83.2 | 80.6 | 52.4 | 63.9 | 60.1 | 60.9 | - | - | - | - |
| Vision-Language | LSTMTrack [11] | WACV2020 | 35.0 | - | 35.0 | - | - | - | 25.0 | 33.0 | 27.0 | 61.0 | 79.0 |
| | SNLT [12] | CVPR2021 | 54.0 | 63.6 | 57.4 | - | - | - | - | - | - | 66.6 | 84.8 |
| | GTI [35] | TCSVT2021 | 47.8 | - | 47.6 | - | - | - | - | - | - | 58.1 | 73.2 |
| | TNL2K-II [32] | CVPR2021 | 51.3 | - | 55.4 | - | - | - | 42.0 | 50.0 | 42.0 | 68.0 | 88.0 |
| | TransVLT [43] | PRL2023 | 66.4 | - | 70.8 | - | - | - | 56.0 | 61.7 | - | 69.9 | 91.2 |
| | JointNLT [46] | CVPR2023 | 60.4 | 69.4 | 63.6 | - | - | - | 56.9 | 73.6 | 58.1 | 65.3 | 85.6 |
| | MMTrack-384 [44] | TCSVT2023 | 70.0 | 82.3 | 75.7 | 49.7 | 59.9 | 55.3 | 58.6 | 75.2 | 59.4 | 70.5 | 91.8 |
| | ATTrack [15] | MM2024 | 63.7 | - | 67.3 | - | - | - | 56.9 | 75.0 | 64.7 | 69.3 | 90.3 |
| | OSDT [38] | TCSVT2024 | 64.3 | 68.6 | 73.4 | - | - | - | 59.3 | - | 61.5 | 66.2 | 86.7 |
| | UVLTrack-B [27] | AAAI2024 | 69.4 | - | 74.9 | 49.2 | - | 55.8 | 62.7 | - | 65.4 | 60.1 | 79.1 |
| | UVLTrack-L [27] | AAAI2024 | 71.3 | - | 78.3 | 51.2 | - | 57.6 | 64.8 | - | 68.8 | 63.5 | 83.2 |
| | QueryNLT [29] | CVPR2024 | 59.9 | 69.6 | 63.5 | - | - | - | 57.8 | 75.6 | 58.7 | 66.7 | 88.2 |
| | DUTrack-256 | Ours | 73.0 | 83.8 | 81.1 | 50.5 | 61.5 | 58.1 | 64.9 | 82.9 | 70.6 | 70.9 | 93.9 |
| | DUTrack-384 | Ours | 74.1 | 84.9 | 82.9 | 52.5 | 63.6 | 60.5 | 65.6 | 83.2 | 71.9 | 71.3 | 95.7 |

DUtrack-256 achieves comparable results of 50.5%, 61.5%, and 58.1% in terms of success, P$_{Norm}$, and precision score. Compared to UVLTrack [27] with the same input size, we outperform by 1.3%, 1.6%, and 2.8% on the three metrics, respectively. DUTrack-384 outperforms MMTrack-384 [44] with an improvement of 2.8% in AUC, 3.7% in P$_{Norm}$, and 5.2% in P. Compared to the top vision-only trackers, DUTrack's performance is on par. We believe the lack of significant improvement is due to the relatively small test dataset in this benchmark, which leads to considerable fluctuations in the experimental results.

**TNL2K** is a tracking benchmark specifically designed for vision-language tracking. Similar to LaSOT, it only provides initial language annotations. As shown in Tab. 2, Our approach gets the best scores of 65.6%, 83.2%, and 71.9% in terms of AUC, P$_{Norm}$, and precision, respectively. Moreover, our DUTrack-256 uses a base encoder, yet it still slightly outperforms UVLTrack-L [27], which uses a large encoder with the same input size, across all three metrics. Compared to ODTrack [45], which currently holds the best scores on this benchmark, our DUTrack-384 achieves

a 4.6% improvement in AUC. Although the TNL2K dataset provides high-quality language annotations, these annotations are fixed and cannot offer more closely matched target information. This limitation hinders previous methods from achieving optimal performance. At the same time, it clearly demonstrates that our dynamic language annotation updates can significantly enhance the performance of current vision-language trackers.

**OTB99-Lang** is an extension of the OTB benchmark in terms of language description. Tracker performance is evaluated using the area under the curve (AUC) and precision (P) metrics, following the protocol established by the official OTB evaluation. As shown in Tab. 2, DUTrack-256 produced results that were competitive, with success and precision scores of 70.9% and 93.9%, respectively. DUTrack-384 gets the best scores of 71.3%, 95.7%, and 71.9% in terms of AUC and precision, respectively.

**GOT-10k** is a vision-only tracking benchmark, and its test set does not provide any language annotations. The dataset is evaluated using two metrics: Average Overlap (AO) and Success Rate (SR). As shown in Tab. 3, DUTrack-

Table 3. Comparison with state-of-the-art methods on GOT-10k and UAV123 benchmarks in AO and AUC score.We add a symbol * over GOT-10k to indicate that the corresponding models are only trained with the GOT-10k training set.

| | SiamFC | MDNet | Ocean | SiamPRN++ | TrDimp | TransT | SimaTrack | VideoTrack | SeqTrack | ODTrack | DUTrack-256 | DUTrack-384 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GOT-10k* | 34.8 | 29.9 | 61.1 | 51.7 | 67.1 | 67.1 | 68.6 | 72.9 | 74.7 | 77.0 | 76.7 | 77.8 |
| UAV123 | 46.8 | 52.8 | 57.4 | 61.0 | 69.1 | 69.1 | 69.8 | 69.7 | 69.2 | - | 69.3 | 70.1 |



Figure 4. Attribute-based evaluation on the LaSOT test set. AUC score is used to rank different trackers.

256 achieved AO 76.7%, which shows an improvement of 2.0% compared to SeqTrack [5]. Compared to ODTrack [45], DUTrack-384 exceeds it by 0.8% in AO.

**UAV123** is a benchmark for object tracking research, focusing on video sequences captured by drones, with its main evaluation metric being AUC. From Tab. 3, we can see that the two versions of DUTrack achieve AUC scores of 69.3% and 70.1%, respectively. The improvement over SimTracker is relatively small, only 0.2%. We believe that language descriptions provide limited assistance for small targets captured by low-altitude drones.

### 4.3. Ablation Study

In this section, we provide detailed experiments to verify the impact of our proposed dynamic update of multi-modal reference information on performance. We use DUTrack-256 without any multi-modal reference information update as the baseline for this experiment. Then, we incrementally add modules to this baseline and conduct more granular comparative experiments on each module.

**Study on the Dynamic Template Capture Module.** To analyze the impact of the DTCM on experimental results, we use the number of top-k patches where the seman-
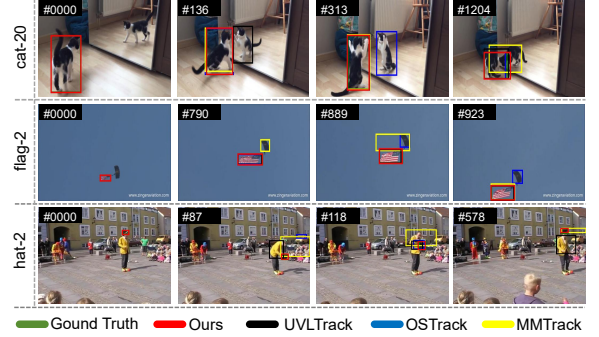


Figure 5. Qualitative comparison results of our tracker with two VL trackers(i.e UVLTrack and MMtrack) and one visual-only tracker OSTrak on three challenging sequences from the LaSOT benchmark. Better viewed in color with zoom-in.

Table 4. Study on different top-k for DTCM.

| Num | top-k | LaSOT | | | GOT-10K | |
|---|---|---|---|---|---|---|
| | | AUC | $P_{Norm}$ | P | AO | $SR_{0.5}$ |
| #1 | 0 | 71.0 | 79.7 | 75.9 | 72.2 | 82.8 |
| #2 | 1 | 71.1 | 79.2 | 77.6 | 73.4 | 83.9 |
| #3 | 2 | 71.5 | 81.0 | 77.8 | 74.2 | 85.0 |
| #4 | 3 | 71.7 | 81.9 | 78.1 | 74.6 | 84.7 |

tic features of the image and language annotations match most closely as a variable. We evaluate the performance of different top-k values on vision-language and vision-only benchmarks. As shown in Tab. 4, when top-k is set to 0, it represents the baseline performance. As the top-k number increases, the AUC on LaSOT increases from 71.0% to 71.7%, and the AO on GOT-10k improves from 72.2% to 74.6%. The notable AUC improvements on both vision-language and vision-only benchmarks clearly indicate that DTCM enhances the tracker's perception of target appearance by capturing new object spatial features.

Table 5. Study on different update parameters for DLUM

| Num | $\Delta S$ | $\Delta D$ | LaSOT | | |
|---|---|---|---|---|---|
| | | | AUC | $P_{Norm}$ | P |
| #1 | 0 | $16 \times$ stride | 72.4 | 83.2 | 80.3 |
| #2 | 0.5 | $2 \times$ stride | 72.5 | 83.2 | 80.4 |
| #3 | 0.8 | $1 \times$ stride | 72.7 | 83.4 | 80.6 |
| #4 | 1 | $0 \times$ stride | 73.0 | 83.8 | 81.6 |

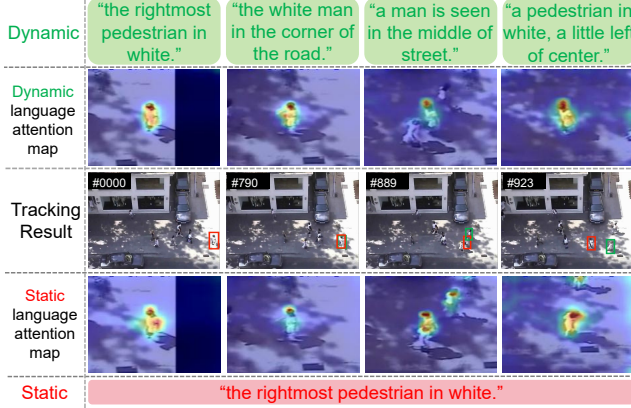**Study on the Dynamic Language Update Module.** To

Figure 6. Visualization of the attention map of the Cls token on the search region using dynamic and initial language annotations. Red represents dynamic annotations, green represents static annotations.

investigate the impact of language update frequency on DU-Track, we design the following experiment: building on the integration of the DTCM, we introduce language information and dynamically update it. In this set of experiments, we selected the variables $\Delta S$ and $\Delta D$ from Eq. 5 and Eq. 6 as the factors to study. The update strategy is as follows: if $\Delta S$ is smaller than the set threshold, or $\Delta D$ is larger than the set threshold, the language description will be updated. Here, the stride refers to the side length of the $16 \times 16$ patches into which the image is divided. In Tab. 5, #1 represents using only the initial language annotations. Comparing this with #4 in Tab. 4, we can see that the introduction of language information into the tracker results in improvements of 0.7%, 1.3%, and 2.3% in AUC, $P_{Norm}$, and P, respectively. By comparing #2, #3, and #4 in Tab. 5, we can observe that as the update frequency increases, the AUC improves from 72.5% to 73.0%. This trend indicates that updating the language descriptions is highly effective in reducing the discrepancy between the language and visual modalities, significantly enhancing the performance of the vision-language tracker.

Table 6. Study on the different LLM for DULM.

| Num | Type | LaSOT | | |
|-----|------|-------|---|---|
| | | AUC | $P_{Norm}$ | P |
| #1 | BLIP | 73.0 | 83.8 | 81.6 |
| #2 | BLIP-2 | 73.2 | 84.1 | 81.7 |
| #3 | DTLLM-Concise | 72.9 | 83.9 | 81.2 |
| #4 | DTLLM-Detailed | 72.5 | 83.3 | 80.6 |

**Study on the LLM.** To investigate the impact of LLM, we conduct detailed experiments based on the current best results. The variable in this set of experiments is the type of LLM, and we test four different language generation meth-

ods. As shown in Tab. 6, using BLIP [22] for language generation resulted in scores of AUC 73.0%, $P_{Norm}$ 83.8%, and P 81.6%. BLIP-2 [23], used in #2, improved AUC by 0.2% compared to #1. It can seamlessly map to the language model, significantly reducing computational costs while maintaining high performance. #3 and #4 use annotations generated by DTLLM [24] as language input, with the difference being the style of the generated language. The concise style annotations show results similar to #1, but the detailed style annotations led to a drop in performance, with AUC and $P_{Norm}$ decreasing by 0.5% and P by 1.0%. Our analysis suggests that the detailed annotations introduce too much learning pressure on the model, adding unnecessary noise that resulted in the performance decline.

## 4.4. Visualization.

To provide a more detailed demonstration of the effectiveness of our proposed method, Fig. 4 presents the AUC scores of DUTrack compared to current mainstream vision-language trackers and vision-only trackers across various challenges. Additionally, Fig. 5 visualizes the tracking results of three challenging sequences. Furthermore, in our experiments investigating DLUM, to verify the effectiveness of our proposed dynamic language annotation, we visualized the attention map of the $[CLS]$ token in the language description over the entire search area. As shown in Fig. 6, in the static language scenario, there is a noticeable misalignment between the tracker's attention and the tracking target. Although our dynamic language annotations may not provide highly precise semantic information, they eliminate this misalignment, thereby better assisting the tracker in locating the object.

## 5. Conclusion

In this paper, we introduce a new vision-language tracking framework, named DUTrack. Compared to previous vision-language trackers that rely solely on the initial template frame and language annotation, DUTrack addresses the issue of static multi-modal references deviating from the target's state by dynamically updating the multi-modal references. Extensive experiments demonstrate that, compared to previous SOTA vision-language trackers, our method effectively enhances vision-language tracking performance.

*Limitation.* Our proposed multi-modal reference update method allows control over the update frequency by setting a threshold. Although manually setting the threshold is simple and reliable, it often requires extensive experimentation to find the optimal value, making the process time-consuming and challenging to achieve consistent results. Enhancing the flexibility and robustness of the update strategy is also our next research direction.

## 6. Acknowledgments

## References

[1] Dosovitskiy Alexey. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv: 2010.11929*, 2020. 3, 4

[2] UT Benchmark. A benchmark and simulator for uav tracking. In *European Conference on Computer Vision*, 2016. 5

[3] Luca Bertinetto, Jack Valmadre, Joao F Henriques, Andrea Vedaldi, and Philip HS Torr. Fully-convolutional siamese networks for object tracking. In *Computer Vision–ECCV 2016 Workshops: Amsterdam, The Netherlands, October 8-10 and 15-16, 2016, Proceedings, Part II 14*, pages 850–865. Springer, 2016. 2, 6

[4] Xin Chen, Bin Yan, Jiawen Zhu, Dong Wang, Xiaoyun Yang, and Huchuan Lu. Transformer tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8126–8135, 2021. 6

[5] Xin Chen, Houwen Peng, Dong Wang, Huchuan Lu, and Han Hu. Seqtrack: Sequence to sequence learning for visual object tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14572–14581, 2023. 6, 7

[6] Zedu Chen, Bineng Zhong, Guorong Li, Shengping Zhang, and Rongrong Ji. Siamese box adaptive network for visual tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6668–6677, 2020. 6

[7] Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3, 4, 5

[8] Heng Fan, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Hexin Bai, Yong Xu, Chunyuan Liao, and Haibin Ling. Lasot: A high-quality benchmark for large-scale single object tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5374–5383, 2019. 5, 6

[9] Heng Fan, Hexin Bai, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Harshit, Mingzhen Huang, Juehuan Liu, et al. Lasot: A high-quality large-scale single object tracking benchmark. *International Journal of Computer Vision*, 129: 439–461, 2021. 5, 6

[10] Qi Feng, Vitaly Ablavsky, Qinxun Bai, Guorong Li, and Stan Sclaroff. Real-time visual object tracking with natural language description. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 700–709, 2020. 2

[11] Qi Feng, Vitaly Ablavsky, Qinxun Bai, Guorong Li, and Stan Sclaroff. Real-time visual object tracking with natural language description. In *WACV*, pages 689–698, 2020. 6

[12] Qi Feng, Vitaly Ablavsky, Qinxun Bai, and Stan Sclaroff. Siamese natural language tracker: Tracking by natural language descriptions with siamese trackers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5851–5860, 2021. 6

[13] Qi Feng, Vitaly Ablavsky, Qinxun Bai, and Stan Sclaroff. Siamese natural language tracker: Tracking by natural language descriptions with siamese trackers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5851–5860, 2021. 2

[14] Jiawei Ge, Jiuxin Cao, Xuelin Zhu, Xinyu Zhang, Chang Liu, Kun Wang, and Bo Liu. Consistencies are all you need for semi-supervised vision-language tracking. In *ACM Multimedia 2024*. 2, 3

[15] Jiawei Ge, Jiuxin Cao, Xuelin Zhu, Xinyu Zhang, Chang Liu, Kun Wang, and Bo Liu. Consistencies are all you need for semi-supervised vision-language tracking. In *ACM Multimedia 2024*, 2024. 6

[16] Mingzhe Guo, Zhipeng Zhang, Heng Fan, and Liping Jing. Divert more attention to vision-language tracking. *Advances in Neural Information Processing Systems*, 35:4446–4460, 2022. 1, 2, 3

[17] Mingzhe Guo, Zhipeng Zhang, Heng Fan, Liping Jing, Yilin Lyu, Bing Li, and Weiming Hu. Learning target-aware representation for visual tracking via informative interactions. In *IJCAI*, pages 927–934, 2022. 6

[18] Xiantao Hu, Bineng Zhong, Qihua Liang, Shengping Zhang, Ning Li, Xianxian Li, and Rongrong Ji. Transformer tracking via frequency fusion. *IEEE Transactions on Circuits and Systems for Video Technology*, 34(2):1020–1031, 2023. 2

[19] Xiantao Hu, Ying Tai, Xu Zhao, Chen Zhao, Zhenyu Zhang, Jun Li, Bineng Zhong, and Jian Yang. Exploiting multimodal spatial-temporal patterns for video object tracking. *arXiv preprint arXiv:2412.15691*, 2024.

[20] Xiantao Hu, Bineng Zhong, Qihua Liang, Shengping Zhang, Ning Li, and Xianxian Li. Toward modalities correlation for rgb-t tracking. *IEEE Transactions on Circuits and Systems for Video Technology*, 34(10):9102–9111, 2024. 2

[21] Lianghua Huang, Xin Zhao, and Kaiqi Huang. Got-10k: A large high-diversity benchmark for generic object tracking in the wild. *IEEE transactions on pattern analysis and machine intelligence*, 43(5):1562–1577, 2019. 5

[22] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pages 12888–12900. PMLR, 2022. 5, 8

[23] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 8

[24] Xuchen Li, Xiaokun Feng, Shiyu Hu, Meiqi Wu, Dailing Zhang, Jing Zhang, and Kaiqi Huang. Dtllm-vlt: Diverse text generation for visual language tracking based on llm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7283–7292, 2024. 5, 8

[25] Zhenyang Li, Ran Tao, Efstratios Gavves, Cees GM Snoek, and Arnold WM Smeulders. Tracking by natural language specification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6495–6503, 2017. 5, 6

[26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 5

[27] Yinchao Ma, Yuyang Tang, Wenfei Yang, Tianzhu Zhang, Jinpeng Zhang, and Mengxue Kang. Unifying visual and vision-language tracking via contrastive learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4107–4116, 2024. 2, 3, 6

[28] Matthias Muller, Adel Bibi, Silvio Giancola, Salman Alsubaihi, and Bernard Ghanem. Trackingnet: A large-scale dataset and benchmark for object tracking in the wild. In *Proceedings of the European conference on computer vision (ECCV)*, pages 300–317, 2018. 5

[29] Yanyan Shao, Shuting He, Qi Ye, Yuchao Feng, Wenhan Luo, and Jiming Chen. Context-aware integration of language and visual references for natural language tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19208–19217, 2024. 5, 6

[30] Liangtao Shi, Bineng Zhong, Qihua Liang, Ning Li, Shengping Zhang, and Xianxian Li. Explicit visual prompts for visual object tracking. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4838–4846, 2024. 2

[31] Yunjie Tian, Lingxi Xie, Jihao Qiu, Jianbin Jiao, Yaowei Wang, Qi Tian, and Qixiang Ye. Fast-itpn: Integrally pre-trained transformer pyramid network with token migration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 5

[32] Xiao Wang, Xiujun Shu, Zhipeng Zhang, Bo Jiang, Yaowei Wang, Yonghong Tian, and Feng Wu. Towards more flexible and accurate object tracking with natural language: Algorithms and benchmark. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13763–13773, 2021. 5, 6

[33] Jinxia Xie, Bineng Zhong, Zhiyi Mo, Shengping Zhang, Liangtao Shi, Shuxiang Song, and Rongrong Ji. Autoregressive queries for adaptive tracking with spatio-temporal transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19300–19309, 2024. 2, 6

[34] Bin Yan, Houwen Peng, Jianlong Fu, Dong Wang, and Huchuan Lu. Learning spatio-temporal transformer for visual tracking. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10448–10457, 2021. 2, 6

[35] Zhengyuan Yang, Tushar Kumar, Tianlang Chen, Jingsong Su, and Jiebo Luo. Grounding-tracking-integration. *IEEE Trans. Circuits Syst. Video Technol.*, pages 3433–3443, 2021. 6

[36] Botao Ye, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Joint feature learning and relation modeling for tracking: A one-stream framework. In *European Conference on Computer Vision*, pages 341–357. Springer, 2022. 2, 6

[37] Chunhui Zhang, Xin Sun, Yiqian Yang, Li Liu, Qiong Liu, Xi Zhou, and Yanfeng Wang. All in one: Exploring unified vision-language tracking with multi-modal alignment. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 5552–5561, 2023. 1, 2

[38] Guangtong Zhang, Bineng Zhong, Qihua Liang, Zhiyi Mo, Ning Li, and Shuxiang Song. One-stream stepwise decreasing for vision-language tracking. *IEEE Transactions on Circuits and Systems for Video Technology*, pages 1–1, 2024. 3, 6

[39] Huanlong Zhang, Jingchao Wang, Jianwei Zhang, Tianzhu Zhang, and Bineng Zhong. One-stream vision-language memory network for object tracking. *IEEE Transactions on Multimedia*, 26:1720–1730, 2023. 2

[40] Lichao Zhang, Abel Gonzalez-Garcia, Joost Van De Weijer, Martin Danelljan, and Fahad Shahbaz Khan. Learning the model update for siamese trackers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4010–4019, 2019. 2

[41] Qian Zhang, Zihao Wang, and Hong Liang. Siamrdt: An object tracking algorithm based on a reliable dynamic template. *Symmetry*, 14(4):762, 2022. 2

[42] Xiaosong Zhang, Yunjie Tian, Lingxi Xie, Wei Huang, Qi Dai, Qixiang Ye, and Qi Tian. Hivit: A simpler and more efficient design of hierarchical vision transformer. In *The Eleventh International Conference on Learning Representations*, 2023. 4, 5

[43] Haojie Zhao, Xiao Wang, Dong Wang, Huchuan Lu, and Xiang Ruan. Transformer vision-language tracking via proxy token guided cross-modal fusion. *Pattern Recognit. Lett.*, 168:10–16, 2023. 6

[44] Yaozong Zheng, Bineng Zhong, Qihua Liang, Guorong Li, Rongrong Ji, and Xianxian Li. Towards unified token learning for vision-language tracking. *IEEE Transactions on Circuits and Systems for Video Technology*, 2023. 2, 5, 6

[45] Yaozong Zheng, Bineng Zhong, Qihua Liang, Zhiyi Mo, Shengping Zhang, and Xianxian Li. Odtrack: Online dense temporal token learning for visual tracking. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7588–7596, 2024. 2, 6, 7

[46] Li Zhou, Zikun Zhou, Kaige Mao, and Zhenyu He. Joint visual grounding and tracking with natural language specification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23151–23160, 2023. 2, 5, 6

[47] Zikun Zhou, Jianqiu Chen, Wenjie Pei, Kaige Mao, Hongpeng Wang, and Zhenyu He. Global tracking via ensemble of local trackers. In *CVPR*, pages 8751–8760, 2022. 6