

ATCTrack: Aligning Target-Context Cues with Dynamic Target States for Robust Vision-Language Tracking

Supplementary Material

A. Target Words Annotation Pipeline

Given the inherently flexible and diverse nature of textual descriptions, it is challenging for trackers to accurately identify target words and context words. In our work, we approach the identification of target words as a multi-label binary classification task, enhancing the model’s ability to recognize target words through supervised learning. However, existing benchmarks [20, 33, 61, 76] provide only textual descriptions without labeled information on the types of target words (*i.e.*, target words or context words). For such a natural language processing task, we leverage the powerful text understanding capabilities of the large language models [37, 70] to construct an automated target words annotation pipeline. Specifically, we employ the widely-used multimodal large language model, GPT-4o [37], and have devised a specific core prompt to guide GPT-4o in recognizing target words (as shown in Fig. A1).

Leveraging our automated annotation pipeline, we complete the labeling of target words in textual data from the MGIT [33], TNL2K [76], LaSOT [19], RefCOCOg [57], OTB99-Lang [49] and Vasttrack [61] datasets. We conduct a random sampling of the labeled results, inspect 50 sentences, and find that the annotations are entirely accurate. This ensures the reliability of our supervised models in classifying target words. In the future, we will open source both the target words label information and our code.

B. Evaluation of Target Words Identification

In this section, we discuss the specific implementation methods for the target words classification accuracy results shown in Fig. 2 (a). Recent studies, such as QueryNLT [65], TTCTrack [58] and OSDT [89], have utilized vision-text similarity metrics to identify target words. Although this is one of their main contributions, they have not provided quantitative evaluation results. For this, we conduct a quantitative analysis based on the target words label information obtained from Sec. A.

B.1. Similarity-Based Target Words Identification

Considering that QueryNLT [65], TTCTrack [58] and OSDT [89] have not open-sourced their code, we employ JointNLT [98], a representative vision-language tracker, as a proxy model for evaluation. The core insight of JointNLT is the use of a one-stream network to jointly model the feature extraction and interaction of text, template images, and search images. The extensive feature interaction among

these elements can, to some extent, represent the feature interaction operations conducted in the aforementioned works for measuring vision-text similarity.

Specifically, at time step t ($t \geq 0$), after the feature encoding by the JointNLT’s backbone network, we obtain the visual features $f_V^t \in \mathbb{R}^{400 \times 512}$ and the textual features $f_L^t \in \mathbb{R}^{L \times 512}$. Here, the length of the visual tokens is fixed at 400, while the length of the textual tokens, L , is determined by the number of words in the sentence. The similarity between them is obtained through the following operations:

$$att_{vl}^t = (f_L^t)^T \cdot f_V^t, \quad (A1)$$

where $att_{vl}^t \in \mathbb{R}^{L \times 400}$ represents the similarity between each visual and textual token. By averaging along the dimension of the search tokens, we can determine the attention each textual token receives at the current time step t , denoted as $att_l^t \in \mathbb{R}^L$.

By concatenating att_l^t at each time step in a video sequence along the time dimension, we can obtain a heatmap of textual feature information for this sequence, denoted as $Att_l \in \mathbb{R}^{L \times T}$, where T represents the number of frames in the video sequence.

For a more intuitive understanding, we conduct a visualization analysis using two video sequences as examples. The related results are depicted in Fig. A2, which serves as a supplement to Fig. 2 (b) and (c) in the main text. As shown in Fig. A2 (a), the target being tracked in this sequence is “plane”. In the corresponding Att_l heatmap, the target word “plane” receives significant attention, indicating that the tracker correctly understands the intent embedded in the text prompt, and this text cue aids in the tracking process. For the example in Fig. A2 (b), the intended tracking target is “yellow people”, but the tracker primarily focuses on the word “the light”. This indicates that the tracker did not correctly focus on the target words, which could mislead the tracking process.

B.2. Evaluation of Target Words Identification Accuracy

In addition to qualitatively demonstrating the tracker’s ability to distinguish each word in the text as described above, we also need to conduct a quantitative evaluation. First, to analyze the tracker’s attention to each word throughout the entire video sequence, we average Att_l along the time dimension, resulting in $Res_l \in \mathbb{R}^L$. Each element in Res_l reflects the amount of attention the tracker gives to the word at the corresponding position.

You are an expert in linguistic analysis for dynamic visual tracking. Your task is to analyze a text description of a target object in a video and identify which phrases describe the **target's intrinsic attributes** (stable properties that remain consistent with the object's physical essence) vs. **contextual attributes** (dynamic properties that may change with scene evolution). Finally, output the phrases of the target's intrinsic attributes in a structured format.

Rules:

1. Target's intrinsic attributes must satisfy:

- Directly describe the target's inherent physical properties (e.g., category, color, material, shape, brand)
- Remain valid even if the target changes pose, location, or interacts with other objects
- Examples: "red", "car", "striped", "round glasses"

2. Contextual attributes must be:

- Related to the target's temporary state or environment (e.g., position, motion, relative relationships)
- Likely to become invalid due to scene dynamics
- Examples: "on the left", "jumping", "next to a chair"

Examples:

1. Input: "a white van parked beside a traffic light"

Output: [{"phrase": "white", "reason": "color is a stable property"}, {"phrase": "van", "reason": "object category"}]

2. Input: "the running black cat with a collar"

Output: [{"phrase": "black", "reason": "color attribute"}, {"phrase": "cat", "reason": "object category"}, {"phrase": "collar", "reason": "physical accessory"}]

3. Input: "the second man from left to right direction"

Output: [{"phrase": "man", "reason": "object category"}]

Question: The input is {xx}, what should the corresponding output be ?



Figure A1. **Prompt used to guide GPT-4o in identifying target words information.** This prompt primarily consists of two parts: task requirement descriptions and example guidance. Replace {xx} with the sentence to be identified to achieve output results similar to the example format.

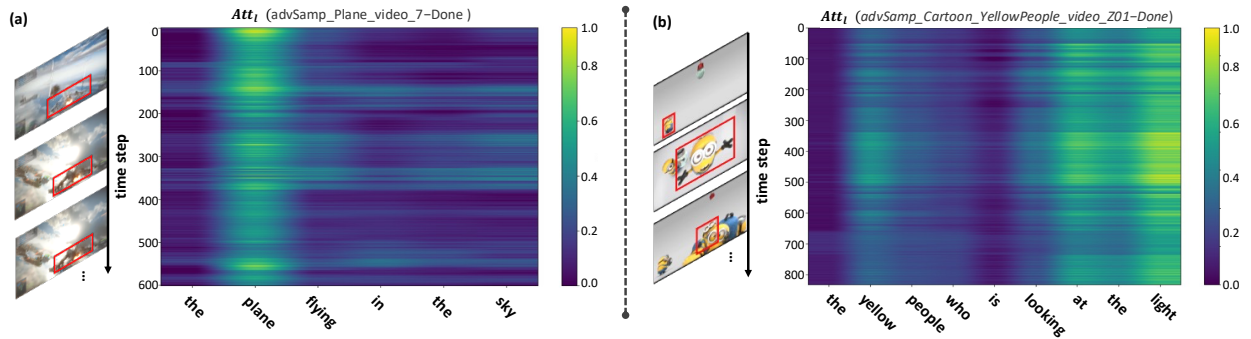


Figure A2. **Visualization results of Att_l across two video sequences.** (a) In the sequence 'advSamp_Plane_video_7-Done', the target "plane" receives significant attention during the tracking process, which aligns with our intended effect. (b) In the sequence 'advSamp_Cartoon_YellowPeople_video_Z01-Done', the target "yellow people" is intended to be tracked, but the tracker primarily focuses on the text "the light". This indicates that the tracker did not correctly focus on the target words, which could mislead the tracking process.

Based on this information, we can map to obtain the tracker's final prediction results for target words classification, $p \in \{0, 1\}^L$. Specifically, based on the target words

label information obtained from Sec. A, we can determine the number of target words k in the sentence. Then, we calculate the top k elements and their indices in Res_l . Sub-

sequently, we set the elements at these indices in p to 1, while all other elements are set to 0.

Additionally, utilizing the target words label information provided in Sec. A, we can obtain a ground truth label $g \in \{0, 1\}^L$. In this label, 0 indicates that the word token at that position is a context word, and 1 indicates it is a target word. We then establish two accuracy assessment metrics, namely, Acc_{all} and $\text{Acc}_{\text{target}}$, by performing different calculations on p and g to evaluate the tracker’s accuracy in classifying target words. Here, Acc_{all} represents the overall classification accuracy of the model for both target and context words; while $\text{Acc}_{\text{target}}$ focuses on the classification accuracy specifically for target words.

$$\text{Acc}_{\text{all}} = \frac{\sum_{i=1}^L \mathbf{1}(p_i = g_i)}{L}, \quad (\text{A2})$$

$$\text{Acc}_{\text{target}} = \frac{\sum_{i=1}^N \mathbf{1}(p_i = 1 \wedge g_i = 1)}{\sum_{i=1}^N \mathbf{1}(g_i = 1)}. \quad (\text{A3})$$

Here, $\mathbf{1}(\cdot)$ is an indicator function that returns 1 if the condition within the parentheses is satisfied.

Similarly, for our proposed ATCTrack and its predictions about target words p_{it} (see Eq. (1)), we can use the same method to map it to p , and then use the above formula for accuracy measurement. The corresponding accuracy results are displayed in Fig. 2 (a). It is evident that our method significantly outperforms methods based on vision-text similarity in both metrics.

B.3. Analysis of Evaluation Results

Fig. 2 (a) shows the target words identification accuracy of our method compared to the existing vision-text similarity-based method [58, 65, 89]. As can be seen, our method achieves an impressive 96.7% in the $\text{Acc}_{\text{target}}$ metric, significantly surpassing the latter’s 29.9%. This precise target word awareness lays a solid foundation for subsequent text cue adjustment and utilization. This demonstrates that our lightweight multilayer perceptron (Eq. (1)) effectively transfers the LLMs’ target word distinguishing capability into the tracker. Although existing LLMs have good target word sensing capabilities, integrating LLMs directly into the tracker incurs substantial computational costs, which is detrimental to practical applications. Additionally, there are lightweight text component analysis tools in the field of natural language processing, such as the widely used Scene Graph Parser [64]. We evaluated the Scene Graph Parser’s accuracy in identifying target words in sentences and found it to be only 21.0%. This indicates that these tools are not yet capable of meeting our target word identification needs in a plug-and-play manner.

C. More Details on the ATCTrack

Due to space constraints, we focus primarily on the main contributions of our paper in the Sec. 3, specifically the textual target-context guidance module (see Sec. 3.2) and the visual target-context guidance module (see Sec. 3.3). For other components of our tracker, such as the prediction head and memory storage module, we provide a brief introduction using current mainstream methods, supplemented by relevant references. In this section, we offer an additional explanation of these components.

C.1. Prediction Head

The prediction head is used to predict the final bbox b^t . We employ a CNN-based tracking head [80, 86], which is widely adopted in tracker design. Firstly, for the search feature $f_R^t \in \mathbb{R}^{N_x \times D}$ that integrates both textual and visual cues, we transform it into a 2D spatial feature map. Subsequently, after passing through L_h stacked Conv-BN-ReLU layers, we obtain a classification score map $P \in [0, 1]^{1 \times H_s \times W_s}$, the size of the bbox $B \in [0, 1]^{2 \times H_s \times W_s}$, and the offset size $O \in [0, 1]^{2 \times H_s \times W_s}$. Then, the position with the highest classification score is considered to be the target position, *i.e.*, $(x_d, y_d) = \arg \max_{(x, y)} P_{xy}$. The final target bbox is obtained as:

$$x = x_d + O(0, x_d, y_d), \quad (\text{A4})$$

$$y = y_d + O(1, x_d, y_d), \quad (\text{A5})$$

$$w = S(0, x_d, y_d), \quad (\text{A6})$$

$$h = S(1, x_d, y_d). \quad (\text{A7})$$

C.2. Memory Storage Module

As introduced in Sec. 3.4, we employ the sliding windows method [7, 80] to update memory units, a method widely used in recent vision trackers focused on temporal modeling. The visual memory feature M in MSM consists of a list of L_m memory units m , denoted as $M = \{m_i\}_{i=1}^{L_m}$. Below, we will illustrate how the sliding windows memory storage method is implemented.

For a video sequence with T frames ($0 \leq t \leq T-1$), the memory units in M need to be initialized when processing the first frame (*i.e.*, $t = 0$). Specifically, after encoding the visual input information via a vision encoder, we obtain the feature $f_{[C]}^0$ encoded from the [CLS] token. Considering that the [CLS] token can represent global visual features [17], we use $f_{[C]}^0$ to initialize the L_m memory units. During the time interval $t \in [1, T-1]$, after tracking each search frame, we obtain the updated memory unit m^t . We pop the memory unit with index 0 from M and append m^t to the end of M .

| Model | Params | Speed | AUC | P |
|---------------|--------|-------|------|------|
| JointNLT [98] | 153M | 31FPS | 56.9 | 58.1 |
| MMTrack [94] | 177M | 37FPS | 58.6 | 59.4 |
| MemVLT [25] | 175M | 32FPS | 63.3 | 67.4 |
| ATCTrack-B | 160M | 35FPS | 67.5 | 73.6 |
| ATCTrack-L | 340M | 30FPS | 68.6 | 75.0 |

Table A1. Results of efficiency analysis.

D. More Details on Model Implementation

Due to space constraints, only core model implementation details are provided in Sec. 4.1. Here, we supplement some additional details. First, regarding the model structure, when performing context words calibration, we use two stacked modules consisting of Eq. (3) and Eq. (4). When executing visual memory representation, we use two stacked modules consisting of Equations Eq. (6) and Eq. (7). It is important to note that we only use the FFN in the visual memory representation part. Considering that the computational cost of FFN in Transformer modules is higher than that of Attention [71], our module design helps reduce the overall parameters and computation of the model.

Additionally, for model training, we use the AdamW optimizer [53] to optimize our model. The text encoder remains frozen, the learning rate is set to 10^{-5} for the vision encoder, 10^{-4} for the remaining unfrozen modules, and the weight decay is set to 10^{-4} . We train for a total of 150 epochs and reduce the learning rate by a factor of 10 after 120 epochs. Finally, during the model inference stage, dynamic template updating follows the implementation of STARK [82]. We set the update interval to 25 and the update confidence threshold to 0.8.

E. Experimental Details of Ablation Studies

In Sec. 4.3, we conduct detailed ablation analyses to investigate the properties of the various modules in ATCTrack. Due to space limitations, we do not fully elaborate on the specific implementation of the ablation experiments. In this section, we provide additional details.

E.1. Ablation Study on important model components

Tab. 2 presents the ablation study results of two core components in our approach: the textual and the visual target-context guidance modules. The specific implementations are as follows:

Tab. 2 (#1) demonstrates the baseline results without our textual and visual object-context guidance modules. In this setup, textual features are processed as a whole entity, an approach widely adopted by recent trackers such as

SNLT [74] and MMTrack [94]. Specifically, we employ a transformer-based decoder to facilitate interaction between textual features f_L and search features f_X^t :

$$f_R^t = Trans_{Dec}(f_X^t, f_L), \quad (A8)$$

where $Trans_{Dec}$ represents the standard transformer decoder layer [71], primarily consisting of attention operations and feed-forward networks. f_R^t denotes the search features embedded with textual cues, which are subsequently fed into the prediction head to obtain final tracking results. To ensure fair comparison, we configure the transformer decoder with four layers, matching the parameter count with the visual and textual object-context guidance module.

Tab. 2 (#2) shows the results using only the textual object-context guidance module. In this implementation, we omit the visual memory guidance process and directly feed the output features f_{XL}^t from the textual target-context guidance module into the prediction head to obtain final results.

Tab. 2 (#3) presents the results using only our visual object-context guidance module. In this implementation, we employ a transformer-based decoder to guide the search features with textual information, which is formulated as:

$$f_{XL}^t = Trans_{Dec}(f_X^t, f_L), \quad (A9)$$

For fair comparison, we implement a two-layer decoder architecture.

Tab. 2 (#4) demonstrates the results of our complete ATCTrack model.

E.2. Ablation Study on Textual Target-Context Modeling

Tab. 3 shows different ways of utilizing textual cues, with the specific implementations for each setting as follows:

Naive method. This setting is consistent with that of Tab. 2 (#1).

+ Target words awareness. This refers to the incorporation of target words awareness method based on the “naive method” setting. Specifically, we concatenate the f_{LT} with f_L to obtain context features f_{LC} for subsequent textual guidance.

+ Context words calibration. This refers to the incorporation of context words calibration operations based on the “+ target words awareness” setting. This is the approach adopted by our ATCTrack.

- Dual-type textual guidance. This approach utilizes only the calibrated single-type text features $f_{L'}$ for textual guidance, where $f_{LC} = f_{L'}$.

| Method | MGIT (Action) | | | TNL2K | | | LaSOT | | | LaSOT _{ext} | | |
|--------------------------------------|---------------|-------------------|-------------|-------------|-------------------|-------------|-------------|-------------------|-------------|----------------------|-------------------|-------------|
| | AUC | P _{Norm} | P | AUC | P _{Norm} | P | AUC | P _{Norm} | P | AUC | P _{Norm} | P |
| <i>Basic Variants</i> | | | | | | | | | | | | |
| Wang [75] | - | - | - | - | - | - | 27.7 | - | 30.4 | - | - | - |
| Feng [21] | - | - | - | 25.0 | 34.0 | 27.0 | 50.0 | - | 56.0 | - | - | - |
| Feng [22] | - | - | - | 25.0 | 33.0 | 27.0 | 35.0 | - | 35.0 | - | - | - |
| GTI [85] | - | - | - | - | - | - | 47.8 | - | 47.6 | - | - | - |
| TNL2K-II [76] | - | - | - | 42.0 | 50.0 | 42.0 | 51.3 | - | 55.4 | - | - | - |
| SNLT [23] | 3.6 | 22.6 | 0.4 | - | - | - | 54.0 | 63.6 | 57.4 | - | - | - |
| VLT _{TT} [30] | 46.8 | 60.2 | 31.8 | 54.7 | 71.8 | 55.3 | 67.3 | 80.2 | 71.5 | 48.4 | 59.9 | 54.3 |
| TransVLT [91] | - | - | - | 56.0 | 61.7 | - | 66.4 | - | 70.8 | - | - | - |
| JointNLT [98] | 61.0 | 78.6 | 44.5 | 56.9 | 73.6 | 58.1 | 60.4 | 69.4 | 63.6 | - | - | - |
| TransNLT [74] | - | - | - | 57.0 | 75.0 | 57.0 | 60.0 | - | 63.0 | - | - | - |
| DecoupleTNL [54] | - | - | - | 56.7 | - | 56.0 | 71.2 | - | 75.3 | - | - | - |
| All-in-One [87] | - | - | - | 55.3 | - | 57.2 | 71.7 | 82.4 | 78.5 | 54.5 | 63.5 | - |
| MMTrack [94] | - | - | - | 58.6 | 75.2 | 59.4 | 70.0 | 82.3 | 75.7 | 49.4 | 59.9 | 55.3 |
| QueryNLT [65] | - | - | - | 56.9 | 73.6 | 58.1 | 59.9 | 69.6 | 63.5 | - | - | - |
| TTCTrack [58] | - | - | - | 58.1 | - | - | 67.6 | - | - | 48.8 | - | - |
| OSDT [89] | - | - | - | 59.3 | 76.2 | 61.5 | 64.3 | 73.4 | 68.6 | - | - | - |
| OneTracker [32] | - | - | - | 58.0 | - | 59.1 | 70.5 | 79.9 | 76.5 | - | - | - |
| UVLTrack-B [56] | - | - | - | 62.7 | - | 65.4 | 69.4 | - | 74.9 | 49.2 | - | 55.8 |
| CTVLT [24] | 69.2 | - | 62.9 | 62.2 | - | 79.5 | 72.3 | - | 79.7 | - | - | - |
| ChatTracker-B [67] | - | - | - | 59.6 | 76.3 | 62.1 | 71.7 | 80.9 | 77.5 | - | - | - |
| MemVLT [25] | 69.4 | 81.3 | 63.7 | 63.3 | 80.9 | 67.4 | 72.9 | 85.7 | 80.5 | 52.1 | 63.3 | 59.8 |
| SUTrack-B224 [14] | - | - | - | 65.0 | - | 67.9 | 73.2 | 83.4 | 80.5 | 53.1 | 64.2 | 60.5 |
| SUTrack-B384 [14] | - | - | - | 65.6 | - | 69.3 | 74.4 | 83.9 | 81.9 | 52.9 | 63.6 | 60.1 |
| ATCTrack-B | 73.7 | 84.5 | 70.1 | 67.5 | 85.3 | 73.6 | 74.6 | 87.0 | 82.1 | 54.6 | 65.7 | 62.8 |
| <i>Performance-oriented Variants</i> | | | | | | | | | | | | |
| ChatTracker-L [67] | - | - | - | 65.4 | 76.5 | 70.2 | 74.1 | 83.8 | 81.2 | - | - | - |
| UVLTrack-L [56] | - | - | - | 64.8 | - | 68.8 | 71.3 | - | 78.3 | 51.2 | - | 59.0 |
| SUTrack-L224 [14] | - | - | - | 66.7 | - | 70.3 | 73.5 | 83.3 | 80.9 | 54.0 | 65.3 | 61.7 |
| SUTrack-L384 [14] | - | - | - | 67.9 | - | 72.1 | 75.2 | 84.9 | 83.2 | 53.6 | 64.2 | 60.5 |
| ATCTrack-L | 74.0 | 86.5 | 76.1 | 68.6 | 85.8 | 75.0 | 74.7 | 87.1 | 82.3 | 55.4 | 66.8 | 64.0 |

Table A2. Comparison with state-of-the-art vision-language trackers on four popular benchmarks: MGIT [33], TNL2K [76], LaSOT [19], and LaSOT_{ext} [20]. The best two results are highlighted in red and blue, respectively.

E.3. Ablation Study on Visual Target-Context Modeling

Tab. 4 shows different ways of utilizing visual cues, with the specific implementations for each setting as follows:

Naive method. This setting is consistent with that of Tab. 2 (#1).

+ ROI. This represents the augmentation of the “naive method” by incorporating explicit visual memory features for tracking assistance. Specifically, we employ the Region of Interest (RoI) approach [62], which is widely adopted in recent Visual-Language Trackers (VLTs) such as JointNLT

[98] and TrDiMP [73]. We apply RoI processing to the search features f_X^t using the predicted bounding box scaled by 1.5 to obtain localized search features $f_{X'}^t \in \mathbb{R}^{36 \times D}$. Subsequently, the visual memory representation process is implemented through the following computations:

$$f_{[C]M'} = \text{Norm}(f_{[C]M} + \Phi_{CA}(f_{[C]M}, f_{X'}^t)), \quad (\text{A10})$$

$$f_{[C]M''} = \text{Norm}(f_{[C]M'} + \text{FFN}(f_{[C]M'})). \quad (\text{A11})$$

+ Search + crop mask. This setting involves using a local mask to construct the object-context indication map. Specifically, for the global object-context indication map

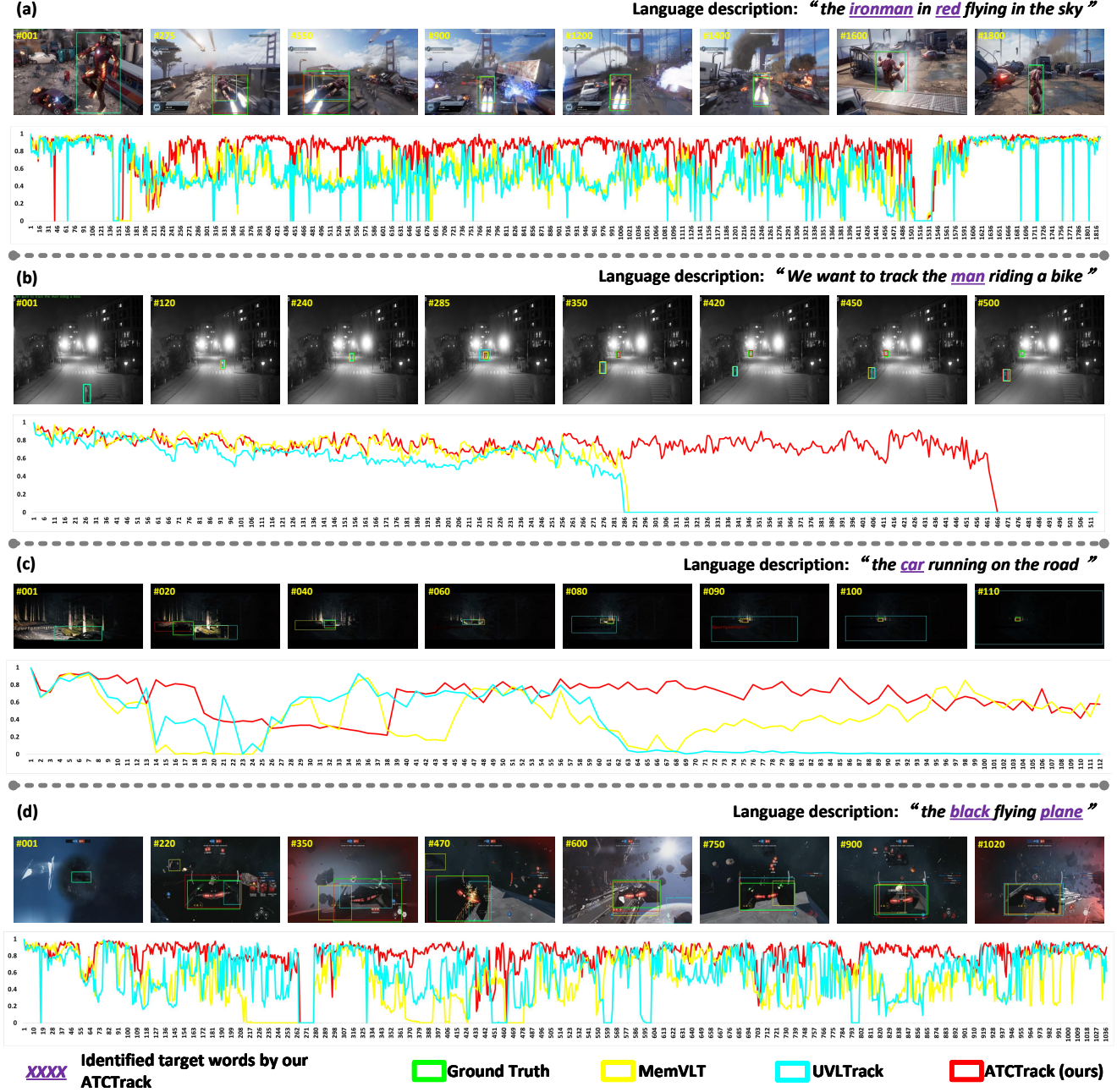


Figure A3. Qualitative comparison results of our tracker with other two state-of-the-art vision-language trackers (*i.e.*, MemVLT and UVLTrack) on four challenging cases. For each video case, we select representative frames to illustrate the predicted bounding boxes of each model and plot the curves of the IOU predictions across the entire video. Better viewed in color with zoom-in.

h^t , we retain only the values within the area corresponding to 1.5 times the predicted bbox, while setting the values in all other areas to zero, resulting in h_i^t . Then, the visual memory representation process is implemented through the following computations:

$$f_{[C]M'} = \text{Norm}(f_{[C]M} + \Phi_{CA}(f_{[C]M}, h_i^t \odot f_X^t)), \quad (\text{A12})$$

$$f_{[C]M''} = \text{Norm}(f_{[C]M'} + \text{FFN}(f_{[C]M'})). \quad (\text{A13})$$

+ **Search + global mask**. This setting involves using a global mask to construct the object-context indication map,

| Method | TNL2K | | | LaSOT | | | LaSOT _{ext} | | |
|--------------------------------------|-------------|-------------------|-------------|-------------|-------------------|-------------|----------------------|-------------------|-------------|
| | AUC | P _{Norm} | P | AUC | P _{Norm} | P | AUC | P _{Norm} | P |
| <i>Basic Variants</i> | | | | | | | | | |
| SiamFC [5] | - | - | - | 29.5 | 45.0 | 28.6 | 33.6 | 42.0 | 33.9 |
| SiamRPN++ [43] | - | - | - | 41.3 | 48.2 | 41.2 | 49.6 | 56.9 | 49.1 |
| SiamBAN [15] | - | - | - | 41.0 | 48.5 | 41.7 | 51.4 | 59.8 | 52.1 |
| TransT [12] | - | - | - | 64.9 | 73.8 | 69.0 | - | - | - |
| Stark [82] | - | - | - | 67.1 | 77.0 | - | - | - | - |
| KeepTrack [59] | - | - | - | 67.1 | 77.2 | 70.2 | - | - | - |
| Mixformer [16] | - | - | - | 69.2 | 78.7 | 74.7 | - | - | - |
| TransInMo [31] | 52.0 | 58.5 | 52.7 | 65.7 | 76.0 | 70.7 | - | - | - |
| OSTrack-256 [86] | 54.3 | - | - | 69.1 | 78.7 | 75.2 | 47.4 | 57.3 | 53.3 |
| OSTrack-384 [86] | 55.9 | - | - | 71.1 | 81.1 | 77.6 | 50.5 | 61.3 | 57.6 |
| AiATrack [28] | - | - | - | 69.0 | 79.4 | 73.8 | 47.7 | 55.6 | 55.4 |
| SimTrack [10] | - | - | - | 69.3 | 78.5 | - | - | - | - |
| GRM [29] | - | - | - | 69.9 | 79.3 | 75.8 | - | - | - |
| SeqTrack-B256 [13] | 54.9 | - | - | 69.9 | 79.7 | 76.3 | 49.5 | 60.8 | 56.3 |
| SeqTrack-B384 [13] | 56.4 | - | - | 71.5 | 81.1 | 77.8 | 50.5 | 61.6 | 57.5 |
| ARTrack-256 [77] | 57.5 | - | - | 70.4 | 79.5 | 76.6 | 46.4 | 56.5 | 52.3 |
| ARTrack-384 [77] | 59.8 | - | - | 72.6 | 81.7 | 79.1 | 51.9 | 62.0 | 58.5 |
| OSTrack-Zoom [39] | 56.5 | - | 57.3 | 70.2 | - | 76.2 | 50.5 | - | 57.4 |
| DropTrack [78] | 56.9 | - | 57.9 | 71.8 | 81.8 | 78.1 | 52.7 | 63.9 | 60.2 |
| ROMTrack-256 [8] | - | - | - | 69.3 | 78.8 | 75.6 | 48.9 | 59.3 | 55.0 |
| ROMTrack-384 [8] | - | - | - | 71.4 | 81.4 | 78.2 | 51.3 | 62.4 | 58.6 |
| F-BDMTrack-256 [84] | 56.4 | - | 56.5 | 69.9 | 79.4 | 75.8 | 47.9 | 57.9 | 54.0 |
| F-BDMTrack-384 [84] | 57.8 | - | 59.4 | 72.0 | 81.5 | 77.7 | 50.8 | 61.3 | 57.8 |
| EVPTTrack-224 [66] | 57.5 | - | 58.8 | 70.4 | 80.9 | 77.2 | 48.7 | 59.5 | 55.1 |
| EVPTTrack-384 [66] | 59.1 | - | 62.0 | 72.7 | 82.9 | 80.3 | 53.7 | 65.5 | 61.9 |
| ODTrack-B [95] | 60.9 | - | - | 73.2 | 83.2 | 80.6 | 52.4 | 63.9 | 60.1 |
| AQATrack-256 [80] | 57.8 | - | 59.4 | 71.4 | 81.9 | 78.6 | 51.2 | 62.2 | 58.9 |
| AQATrack-384 [80] | 59.3 | - | 62.3 | 72.7 | 82.9 | 80.2 | 52.7 | 64.2 | 60.8 |
| ARTrackV2-256 [3] | - | - | - | 71.6 | 80.2 | 77.2 | 50.8 | 61.9 | 57.7 |
| ARTrackV2-384 [3] | - | - | - | 73.0 | 82.0 | 79.6 | 52.9 | 63.4 | 59.1 |
| HIPTrack [6] | - | - | - | 72.7 | 82.9 | 79.5 | 53.0 | 64.3 | 60.6 |
| OneTracker [32] | 58.0 | - | 59.1 | 70.5 | 79.9 | 76.5 | - | - | - |
| LoRAT-B224 [50] | 58.8 | - | 61.3 | 71.7 | 80.9 | 77.3 | 50.3 | 61.6 | 57.1 |
| LoRAT-B378 [50] | 59.9 | - | 63.7 | 72.9 | 81.9 | 79.1 | 53.1 | 64.8 | 60.6 |
| SUTrack-B224 [14] | 65.0 | - | 67.9 | 73.2 | 83.4 | 80.5 | 53.1 | 64.2 | 60.5 |
| SUTrack-B384 [14] | 65.6 | - | 69.3 | 74.4 | 83.9 | 81.9 | 52.9 | 63.6 | 60.1 |
| ATCTrack-B | 67.5 | 85.3 | 73.6 | 74.6 | 87.0 | 82.1 | 54.6 | 65.7 | 62.8 |
| <i>Performance-oriented Variants</i> | | | | | | | | | |
| ODTrack-L [95] | 61.7 | - | - | 74.0 | 84.2 | 82.3 | 53.9 | 65.4 | 61.7 |
| LoRAT-L224 [50] | 61.1 | - | 65.1 | 74.2 | 83.6 | 80.9 | 52.8 | 64.7 | 60.0 |
| LoRAT-L378 [50] | 62.3 | - | 67.0 | 75.1 | 84.1 | 82.0 | 56.6 | 69.0 | 65.1 |
| SUTrack-L224 [14] | 66.7 | - | 70.3 | 73.5 | 83.3 | 80.9 | 54.0 | 65.3 | 61.7 |
| SUTrack-L384 [14] | 67.9 | - | 72.1 | 75.2 | 84.9 | 83.2 | 53.6 | 64.2 | 60.5 |
| ACTrack-L | 68.6 | 85.8 | 75.0 | 74.7 | 87.1 | 82.3 | 55.4 | 66.8 | 64.0 |

Table A3. Comparison with state-of-the-art vision-only trackers on three popular benchmarks: TNL2K [76], LaSOT [19], and LaSOT_{ext} [20]. The best two results are highlighted in red and blue, respectively.

which is used to obtain explicit visual memory features. This is the approach adopted by our ATCTrack.

E.4. Ablation Study on the Contribution of different modules

w/o HiViT backbone. This setting refers to replacing the HiViT backbone [69, 90] with the ViT backbone typically used in conventional trackers [16, 86].

w/o dynamic template. This setting refers to using only the original static template for visual input, without the sparse dynamic template [82].

w/o Textual_{TC} & Visual_{TC}. This setting is the same as setting in Tab. 2 (#1), meaning that the visual and textual target-context guidance mechanism we designed is not utilized.

w/o target words label. This setting, with the model structure unchanged, refers to not using target words supervision signals, thus excluding L_{bce} loss.

F. Additional Experimental Results

F.1. Efficiency Analysis

In Tab. A1, we compare ATCTrack with the latest VLTs (*i.e.*, JointNLT [98], MMTrack [94], and MemVLT [25]) in terms of efficiency (Params and Speed) and performance (AUC and P on TNL2K). For ATCTrack-B, the parameters and tracking speed are comparable to recent trackers, but it shows significant performance advantages, such as a 4.2% improvement in AUC compared to MemVLT. For ATCTrack-L, the parameter scale is considerably larger than ATCTrack-B, which leads to a further performance improvement.

F.2. Comparison with More Trackers

In Tab. 1 of Sec. 4.2, due to space constraints, we compare ATCTrack with several recent high-performance vision-language trackers. As a supplement, Tab. A2 presents the performance of a broader range of vision-language trackers. Additionally, in line with the prevailing paradigm of vision-language tracking models [25, 94, 98], Tab. A3 provides additional comparisons with vision-only trackers. The strong performance of our model among these trackers further demonstrates the effectiveness of our approach.

G. More Qualitative Results

Due to space limitations, Fig. 4 only presents four cases for the qualitative comparison between our model and the latest SOTA models. In this section, we provide additional qualitative comparison results, as illustrated in Fig. A3.

References

- [1] Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Visual prompting: Modifying pixel space to adapt pre-trained models. *arXiv preprint arXiv:2203.17274*, 3:11–12, 2022. 3
- [2] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 2, 4
- [3] Yifan Bai, Zeyang Zhao, Yihong Gong, and Xing Wei. Artrackv2: Prompting autoregressive tracker where to look and how to describe. *arXiv preprint arXiv:2312.17133*, 2023. 7
- [4] Moshe Bar. Visual objects in context. *Nature Reviews Neuroscience*, 5(8):617–629, 2004. 4
- [5] Luca Bertinetto, Jack Valmadre, Joao F Henriques, Andrea Vedaldi, and Philip HS Torr. Fully-convolutional siamese networks for object tracking. In *European conference on computer vision*, pages 850–865. Springer, 2016. 7
- [6] Wenrui Cai, Qingjie Liu, and Yunhong Wang. Learning historical status prompt for accurate and robust visual tracking. *arXiv preprint arXiv:2311.02072*, 2023. 7
- [7] Wenrui Cai, Qingjie Liu, and Yunhong Wang. Hip-track: Visual tracking with historical prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19258–19267, 2024. 6, 3
- [8] Yidong Cai, Jie Liu, Jie Tang, and Gangshan Wu. Robust object modeling for visual tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9589–9600, 2023. 7
- [9] Monica S Castelano and Chelsea Heaven. The relative contribution of scene context and target features to visual search in scenes. *Attention, Perception, & Psychophysics*, 72(5):1283–1297, 2010. 4
- [10] Boyu Chen, Peixia Li, Lei Bai, Lei Qiao, Qihong Shen, Bo Li, Weihao Gan, Wei Wu, and Wanli Ouyang. Backbone is all your need: A simplified architecture for visual object tracking. In *European Conference on Computer Vision*, pages 375–392. Springer, 2022. 7
- [11] Honghao Chen, Yurong Zhang, Xiaokun Feng, Xi-angxiang Chu, and Kaiqi Huang. Revealing the dark secrets of extremely large kernel convnets on robustness. *arXiv preprint arXiv:2407.08972*, 2024. 7
- [12] 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. 7

- [13] 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. 7
- [14] Xin Chen, Ben Kang, Wanting Geng, Jiawen Zhu, Yi Liu, Dong Wang, and Huchuan Lu. Sutrack: Towards simple and unified single object tracking. *arXiv preprint arXiv:2412.19138*, 2024. 6, 8, 5, 7
- [15] Zedu Chen, Bineng Zhong, Guorong Li, Shengping Zhang, and Rongrong Ji. Siamese box adaptive network for visual tracking. In *CVPR*, pages 6668–6677, 2020. 7
- [16] Yutao Cui, Cheng Jiang, Limin Wang, and Gangshan Wu. Mixformer: End-to-end tracking with iterative mixed attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13608–13618, 2022. 7, 8
- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3, 4, 5
- [18] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. 5
- [19] 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 5369–5378, 2019. 2, 4, 6, 7, 1, 5
- [20] 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. 2, 6, 1, 5, 7
- [21] Qi Feng, Vitaly Ablavsky, Qinxun Bai, and Stan Sclaroff. Robust visual object tracking with natural language region proposal network. *arXiv preprint arXiv:1912.02048*, 1(7):8, 2019. 5
- [22] 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. 5
- [23] 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 5847–5856, 2021. 3, 5
- [24] Xiaokun Feng, Dailing Zhang, Shiyu Hu, Xuchen Li, Meiqi Wu, Jing Zhang, Xiaotang Chen, and Kaiqi Huang. Enhancing vision-language tracking by effectively converting textual cues into visual cues. *arXiv preprint arXiv:2412.19648*, 2024. 6, 5
- [25] Xiaokun Feng, Xuchen Li, Shiyu Hu, Dailing Zhang, Jing Zhang, Xiaotang Chen, Kaiqi Huang, et al. Memvlt: Vision-language tracking with adaptive memory-based prompts. *Advances in Neural Information Processing Systems*, 37:14903–14933, 2025. 3, 4, 5, 6, 7, 8
- [26] Xiaokun Feng, Haiming Yu, Meiqi Wu, Shiyu Hu, Jintao Chen, Chen Zhu, Jiahong Wu, Xiangxiang Chu, and Kaiqi Huang. Narrlv: Towards a comprehensive narrative-centric evaluation for long video generation models. *arXiv preprint arXiv:2507.11245*, 2025. 7
- [27] Xiaokun Feng, Dailing Zhang, Shiyu Hu, Xuchen Li, Meiqi Wu, Jing Zhang, Xiaotang Chen, and Kaiqi Huang. Cstrack: Enhancing rgb-x tracking via compact spatiotemporal features. *arXiv preprint arXiv:2505.19434*, 2025. 7
- [28] Shenyuan Gao, Chunluan Zhou, Chao Ma, Xinggang Wang, and Junsong Yuan. Aiatrack: Attention in attention for transformer visual tracking. In *European Conference on Computer Vision*, pages 146–164. Springer, 2022. 7
- [29] Shenyuan Gao, Chunluan Zhou, and Jun Zhang. Generalized relation modeling for transformer tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18686–18695, 2023. 7
- [30] Mingzhe Guo, Zhipeng Zhang, Heng Fan, and Liping Jing. Divert more attention to vision-language tracking. In *Proceedings of the Advances in Neural Information Processing Systems*, pages 4446–4460, 2022. 2, 5
- [31] 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. *arXiv preprint arXiv:2201.02526*, 2022. 7
- [32] Lingyi Hong, Shilin Yan, Renrui Zhang, Wanyun Li, Xinyu Zhou, Pinxue Guo, Kaixun Jiang, Yiting Chen, Jinglun Li, Zhaoyu Chen, et al. Onetracker: Unifying visual object tracking with foundation models and efficient tuning. *arXiv preprint arXiv:2403.09634*, 2024. 3, 6, 5, 7
- [33] Shiyu Hu, Dailing Zhang, Meiqi Wu, Xiaokun Feng, Xuchen Li, Xin Zhao, and Kaiqi Huang. A multi-modal global instance tracking benchmark (mgit): Better locating target in complex spatio-temporal and

- causal relationship. In *the 37th Conference on Neural Information Processing Systems*, pages 25007–25030, 2023. [1](#), [2](#), [4](#), [6](#), [7](#), [5](#)
- [34] Shiyu Hu, Xin Zhao, Lianghua Huang, and Kaiqi Huang. Global instance tracking: Locating target more like humans. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1):576–592, 2023. [7](#)
- [35] Shiyu Hu, Xin Zhao, and Kaiqi Huang. Sotverse: A user-defined task space of single object tracking. *International Journal of Computer Vision*, 132:872–930, 2024. [2](#)
- [36] 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. [7](#)
- [37] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. [1](#)
- [38] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. [5](#)
- [39] Yutong Kou, Jin Gao, Bing Li, Gang Wang, Weiming Hu, Yizheng Wang, and Liang Li. Zoomtrack: Target-aware non-uniform resizing for efficient visual tracking. *Advances in Neural Information Processing Systems*, 36:50959–50977, 2023. [7](#)
- [40] Dana Lahat, Tülay Adalı, and Christian Jutten. Multimodal data fusion: an overview of methods, challenges, and prospects. *Proceedings of the IEEE*, 103(9):1449–1477, 2015. [1](#)
- [41] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European conference on computer vision (ECCV)*, pages 734–750, 2018. [6](#)
- [42] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021. [3](#)
- [43] Bo Li, Wei Wu, Qiang Wang, Fangyi Zhang, Junliang Xing, and Junjie Yan. SiamRPN++: Evolution of siamese visual tracking with very deep networks. In *CVPR*, pages 4282–4291, 2019. [7](#)
- [44] Xuchen Li, Xiaokun Feng, Shiyu Hu, Meiqi Wu, Dailing Zhang, Jing Zhang, and Kaiqi Huang. Dtlm-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. [2](#)
- [45] Xuchen Li, Shiyu Hu, Xiaokun Feng, Dailing Zhang, Meiqi Wu, Jing Zhang, and Kaiqi Huang. Dtlvlt: A multi-modal diverse text benchmark for visual language tracking based on llm. *arXiv preprint arXiv:2410.02492*, 2024. [7](#)
- [46] Xuchen Li, Shiyu Hu, Xiaokun Feng, Dailing Zhang, Meiqi Wu, Jing Zhang, and Kaiqi Huang. How texts help? a fine-grained evaluation to reveal the role of language in vision-language tracking. *arXiv preprint arXiv:2411.15600*, 2024.
- [47] Xuchen Li, Shiyu Hu, Xiaokun Feng, Dailing Zhang, Meiqi Wu, Jing Zhang, and Kaiqi Huang. Visual language tracking with multi-modal interaction: A robust benchmark. *arXiv preprint arXiv:2409.08887*, 2024. [7](#)
- [48] Yihao Li, Jun Yu, Zhongpeng Cai, and Yuwen Pan. Cross-modal target retrieval for tracking by natural language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4931–4940, 2022. [2](#)
- [49] Zhenyang Li, Ran Tao, Efstratios Gavves, Cees GM Snoek, and Arnold WM Smeulders. Tracking by natural language specification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6495–6503, 2017. [1](#), [2](#), [3](#), [4](#), [7](#)
- [50] Liting Lin, Heng Fan, Zhipeng Zhang, Yaowei Wang, Yong Xu, and Haibin Ling. Tracking meets lora: Faster training, larger model, stronger performance. In *ECCV*, 2024. [7](#)
- [51] Xinran Ling, Chen Zhu, Meiqi Wu, Hangyu Li, Xiaokun Feng, Cundian Yang, Aiming Hao, Jiashu Zhu, Jiahong Wu, and Xiangxiang Chu. Vmbench: A benchmark for perception-aligned video motion generation. *arXiv preprint arXiv:2503.10076*, 2025. [7](#)
- [52] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. [4](#), [6](#)
- [53] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. [7](#), [4](#)
- [54] Ding Ma and Xiangqian Wu. Tracking by natural language specification with long short-term context decoupling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14012–14021, 2023. [6](#), [5](#)
- [55] Guoqing Ma, Haoyang Huang, Kun Yan, Liangyu Chen, Nan Duan, Shengming Yin, Changyi Wan, Ranchen Ming, Xiaoniu Song, Xing Chen, et al. Step-video-t2v technical report: The practice, challenges,

- and future of video foundation model. *arXiv preprint arXiv:2502.10248*, 2025. 5
- [56] 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. 3, 6, 7, 5
- [57] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 11–20, 2016. 7, 1
- [58] Zhongjie Mao, Yucheng Wang, Xi Chen, and Jia Yan. Textual tokens classification for multi-modal alignment in vision-language tracking. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8025–8029. IEEE, 2024. 2, 4, 1, 3, 5
- [59] Christoph Mayer, Martin Danelljan, Danda Pani Paudel, and Luc Van Gool. Learning target candidate association to keep track of what not to track. In *ICCV*, pages 13444–13454, 2021. 7
- [60] 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. 7
- [61] Liang Peng, Junyuan Gao, Xinran Liu, Weihong Li, Shaohua Dong, Zhipeng Zhang, Heng Fan, and Libo Zhang. Vasttrack: Vast category visual object tracking. *Advances in Neural Information Processing Systems*, 37:130797–130818, 2025. 7, 1
- [62] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 3, 5
- [63] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 658–666, 2019. 6
- [64] Sebastian Schuster, Ranjay Krishna, Angel Chang, Li Fei-Fei, and Christopher D Manning. Generating semantically precise scene graphs from textual descriptions for improved image retrieval. In *Proceedings of the fourth workshop on vision and language*, pages 70–80, 2015. 8, 3
- [65] Yanyan Shao, Shuting He, Qi Ye, Yuchao Feng, Wenhao Luo, and Jiming Chen. Context-aware integration of language and visual references for natural language tracking. *arXiv preprint arXiv:2403.19975*, 2024. 2, 3, 4, 6, 1, 5
- [66] Liangtao Shi, Bineng Zhong, Qihua Liang, Ning Li, Shengping Zhang, and Xianxian Li. Explicit visual prompts for visual object tracking. *arXiv preprint arXiv:2401.03142*, 2024. 6, 7
- [67] Yiming Sun, Fan Yu, Shaoxiang Chen, Yu Zhang, Junwei Huang, Yang Li, Chenhui Li, and Changbo Wang. Chattracker: Enhancing visual tracking performance via chatting with multimodal large language model. *Advances in Neural Information Processing Systems*, 37:39303–39324, 2025. 6, 5
- [68] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning? *Advances in neural information processing systems*, 33:6827–6839, 2020. 3
- [69] 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. 6, 8
- [70] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 2, 4, 1
- [71] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Proceedings of the Advances in Neural Information Processing Systems*, 30, 2017. 5, 4
- [72] Hongyu Wang, Xiaotao Liu, Yifan Li, Meng Sun, Dian Yuan, and Jing Liu. Temporal adaptive rgbt tracking with modality prompt. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5436–5444, 2024. 3
- [73] Ning Wang, Wengang Zhou, Jie Wang, and Houqiang Li. Transformer meets tracker: Exploiting temporal context for robust visual tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1571–1580, 2021. 3, 8, 5
- [74] Rong Wang, Zongheng Tang, Qianli Zhou, Xiaoqian Liu, Tianrui Hui, Quange Tan, and Si Liu. Unified transformer with isomorphic branches for natural language tracking. *IEEE Transactions on Circuits and Systems for Video Technology*, 2023. 4, 5
- [75] Xiao Wang, Chenglong Li, Rui Yang, Tianzhu Zhang, Jin Tang, and Bin Luo. Describe and attend to track: Learning natural language guided structural represen-

- tation and visual attention for object tracking. *arXiv preprint arXiv:1811.10014*, 2018. 5
- [76] 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. 2, 4, 6, 7, 1, 5
- [77] Xing Wei, Yifan Bai, Yongchao Zheng, Dahu Shi, and Yihong Gong. Autoregressive visual tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9697–9706, 2023. 7
- [78] Qiangqiang Wu, Tianyu Yang, Ziquan Liu, Baoyuan Wu, Ying Shan, and Antoni B Chan. Dropmae: Masked autoencoders with spatial-attention dropout for tracking tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14561–14571, 2023. 7
- [79] Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. Object tracking benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(09):1834–1848, 2015. 7
- [80] Jinxia Xie, Bineng Zhong, Zhiyi Mo, Shengping Zhang, Liangtao Shi, Shuxiang Song, and Rongrong Ji. Autoregressive queries for adaptive tracking with spatio-temporal transformers. *arXiv preprint arXiv:2403.10574*, 2024. 6, 8, 3, 7
- [81] Chenlong Xu, Bineng Zhong, Qihua Liang, Yaozong Zheng, Guorong Li, and Shuxiang Song. Less is more: Token context-aware learning for object tracking. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8824–8832, 2025. 2
- [82] 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. 5, 6, 8, 4, 7
- [83] 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. 3
- [84] Dawei Yang, Jianfeng He, Yinchao Ma, Qianjin Yu, and Tianzhu Zhang. Foreground-background distribution modeling transformer for visual object tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10117–10127, 2023. 7
- [85] Zhengyuan Yang, Tushar Kumar, Tianlang Chen, Jingsong Su, and Jiebo Luo. Grounding-tracking-integration. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(9):3433–3443, 2021. 3, 7, 5
- [86] Botao Ye, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Joint feature learning and relation modeling for tracking: A one-stream framework. In *Proceedings of the European Conference on Computer Vision*, pages 341–357, 2022. 3, 5, 6, 8, 7
- [87] 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. 3, 6, 7, 5
- [88] Dailing Zhang, Shiyu Hu, Xiaokun Feng, Xuchen Li, Jing Zhang, Kaiqi Huang, et al. Beyond accuracy: Tracking more like human via visual search. *Advances in Neural Information Processing Systems*, 37:2629–2662, 2025. 3, 5
- [89] Guangtong Zhang, Bineng Zhong, Qihua Liang, Zhiyi Mo, Ning Li, and Shuxiang Song. One-stream step-wise decreasing for vision-language tracking. *IEEE Transactions on Circuits and Systems for Video Technology*, 2024. 3, 1, 5
- [90] 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*, 2022. 6, 8
- [91] Haojie Zhao, Xiao Wang, Dong Wang, Huchuan Lu, and Xiang Ruan. Transformer vision-language tracking via proxy token guided cross-modal fusion. *Pattern Recognition Letters*, 168:10–16, 2023. 5
- [92] Xin Zhao, Shiyu Hu, Yipei Wang, Jing Zhang, Yimin Hu, Rongshuai Liu, Haibin Ling, Yin Li, Renshu Li, Kun Liu, and Jiadong Li. Biodrone: A bionic drone-based single object tracking benchmark for robust vision. *International Journal of Computer Vision*, 132:1659–1684, 2024. 2
- [93] Yaozong Zheng, Bineng Zhong, Qihua Liang, Zhenjun Tang, Rongrong Ji, and Xianxian Li. Leveraging local and global cues for visual tracking via parallel interaction network. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(4):1671–1683, 2022. 2
- [94] 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. 3, 5, 6, 7, 4, 8
- [95] Yaozong Zheng, Bineng Zhong, Qihua Liang, Zhiyi Mo, Shengping Zhang, and Xianxian Li. Odtrack: Online dense temporal token learning for visual tracking. *arXiv preprint arXiv:2401.01686*, 2024. 3, 7

- [96] Yaozong Zheng, Bineng Zhong, Qihua Liang, Ning Li, and Shuxiang Song. Decoupled spatio-temporal consistency learning for self-supervised tracking. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 10635–10643, 2025. [3](#)
- [97] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16816–16825, 2022. [3](#)
- [98] 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](#), [3](#), [6](#), [7](#), [1](#), [4](#), [5](#), [8](#)