

1 Appendix

2 A Positional Encoding

3 Transformer requires a positional encoding to identify the position of the current processing token [8].
 4 Through a series of comparison experiments, we choose *untied positional encoding*, which is proposed
 5 in TUPE [6], as the positional encoding solution of our tracker. In addition, we generalize the *untied*
 6 *positional encoding* to arbitrary dimensions to fit with other components in our tracker.

7 The original transformer [8] proposes a absolute positional encoding method to represent the position:
 8 a fixed or learnable vector p_i is assigned to each position i . Starting from the basic attention module,
 9 we have:

$$\text{Atten}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}V\right), \quad (1)$$

10 where Q, K, V are the *query* vector, *key* vector and *value* vector, which are the parameters of the
 11 attention function, d_k is the dimension of *key*. Introducing the linear projection matrix and multi-head
 12 attention to the attention module (1), we get the multi-head variant defined in [8]:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W_O, \quad (2)$$

13 where $\text{head}_i = \text{Atten}(QW_i^Q, KW_i^K, VW_i^V)$, $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^V \in$
 14 $\mathbb{R}^{d_{\text{model}} \times d_v}$, $W_i^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$ and h is the number of heads. For simplicity, as in [6], we assume
 15 that $d_k = d_v = d_{\text{model}}$, and use the single-head version of self-attention module. Denoting the input
 16 sequence as $x = x_1, x_2, \dots, x_n$, where n is the length of sequence, x_i is the i -th token in the input
 17 data. Denoting the output sequence as $z = (z_1, z_2, \dots, z_n)$. Self-attention module can be rewritten
 18 as

$$z_i = \sum_{j=1}^n \frac{\exp(\alpha_{ij})}{\sum_{j'=1}^n \exp(\alpha_{ij'})} (x_j W^V), \quad (3)$$

$$\text{where } \alpha_{ij} = \frac{1}{\sqrt{d}} (x_i W^Q)(x_j W^K)^T. \quad (4)$$

19 Obviously, the self-attention module is permutation-invariance. Thus it can not "understand" the
 20 order of input tokens.

21 **Untied absolute positional encoding.** By adding a learnable positional encoding [8] to the single-
 22 head self-attention module, we can obtain the following equation:

$$\begin{aligned} \alpha_{ij}^{\text{Abs}} &= \frac{((w_i + p_i)W^Q)((w_j + p_j)W^K)^T}{\sqrt{d}} \\ &= \frac{(w_i W^Q)(w_j W^K)^T}{\sqrt{d}} + \frac{(w_i W^Q)(p_j W^K)^T}{\sqrt{d}} \\ &\quad + \frac{(p_i W^Q)(w_j W^K)^T}{\sqrt{d}} + \frac{(p_i W^Q)(p_j W^K)^T}{\sqrt{d}}. \end{aligned} \quad (5)$$

23 The equation (5) is expanded into four terms: token-to-token, token-to-position, position-to-token,
 24 position-to-position. [6] discuss the problems exists in the equation and proposes the *untied absolute*
 25 *positional encoding*, which unties the correlation between tokens and positions by removing the
 26 token-position correlation terms in equation (5), and using an isolated pair of projection matrices U^Q
 27 and U^K to perform linear transformation upon positional embedding vector. The following is the
 28 new formula for obtaining α_{ij} using the *untied absolute positional encoding* in the l -th layer:

$$\begin{aligned} \alpha_{ij} &= \frac{1}{\sqrt{2d}} (x_i^l W^{Q,l})(x_j^l W^{K,l})^T \\ &\quad + \frac{1}{\sqrt{2d}} (p_i U^Q)(p_j U^K)^T. \end{aligned} \quad (6)$$

29 where p_i and p_j is the positional embedding at position i and j respectively, $U^Q \in \mathbb{R}^{d \times d}$ and
 30 $U^K \in \mathbb{R}^{d \times d}$ are learnable projection matrices for the positional embedding vector. When extending

31 to the multi-head version, the positional embedding p_i is shared across different heads, while U^Q and
 32 U^K are different for each head.

33 **Relative positional bias.** According to [7], relative positional encoding is a necessary supplement to
 34 absolute positional encoding. In [6], a relative positional encoding is applied by adding a relative
 35 positional bias to equation (6):

$$\begin{aligned} \alpha_{ij} = & \frac{1}{\sqrt{2d}}(x_i^l W^{Q,l})(x_j^l W^{K,l})^T \\ & + \frac{1}{\sqrt{2d}}(p_i U^Q)(p_j U^K)^T + b_{j-i}, \end{aligned} \quad (7)$$

36 where for each $j - i$, b_{j-i} is a learnable scalar. The *relative positional bias* is also shared across
 37 layers. When extending to the multi-head version, b_{j-i} is different for each head.

38 **Generalize to multiple dimensions.** Before working with our tracker’s encoder and decoder
 39 network, we need to extend the *untied positional encoding* to a multidimensional version. One
 40 straightforward method is allocating a positional embedding matrix for every dimension and summing
 41 up all embedding vectors from different dimensions at the corresponding index to represent the final
 42 embedding vector. Together with *relative positional bias*, for an n -dimensional case, we have:

$$\begin{aligned} \underbrace{\alpha_{ij} \dots}_{n} \underbrace{mn \dots}_{n} = & \frac{1}{\sqrt{2d}}(\underbrace{x_{ij} \dots}_{n} W^Q)(\underbrace{x_{mn} \dots}_{n} W^K)^T \\ & + \frac{1}{\sqrt{2d}}[\underbrace{(p_i^1 + p_j^2 + \dots)}_n U^Q][\underbrace{(p_m^1 + p_n^2 + \dots)}_n U^K]^T \\ & + \underbrace{b_{m-i, n-j, \dots}}_n. \end{aligned} \quad (8)$$

43 **Generalize to concatenation-based fusion.** In order to work with *concatenation-based fusion*, the
 44 *untied absolute positional encoding* is also concatenated to match the real position, the indexing tuple
 45 of *relative positional bias* now appends with a pair of indices to reflect the origination of *query* and
 46 *key* involved currently.

47 Taking l -th layer in the encoder as the example:

$$\begin{aligned} \alpha_{ij, mn, g, h} = & \frac{1}{\sqrt{2d}}(x_{ij, g}^l W^{Q, l})(x_{mn, h}^l W^{K, l})^T \\ & + \frac{1}{\sqrt{2d}}[(p_{i, g}^1 + p_{j, g}^2) U_g^Q][(p_{m, h}^1 + p_{n, h}^2) U_h^K]^T \\ & + b_{m-i, n-j, g, h}, \end{aligned} \quad (9)$$

48 where g and h are the index of the origination of *query* and *key* respectively, for instance, 1 for the
 49 tokens from the template image, 2 for the tokens from the search image. The form in the decoder is
 50 similar, except that g is fixed. In our implementation, the parameters of *untied positional encoding*
 51 are shared inside the encoder and the decoder, respectively.

52 B Figures on LaSOT Test set

53 Fig. 1 and Fig. 2 show the success plot and the precision plot respectively. The comparison includes
 54 our SwinTrack-T-224, our SwinTrack-B-384, DiMP [1], STMTrack[5], TiDiMP[9], TransT[2] and
 55 STARK[10].

56 C Response Visualization

57 We provide the heatmap visualization of the response map generated by the IoU-aware classification
 58 branch head in our SwinTrack-B-384 in Fig. 3. The visualized sequences are from LaSOT_{ext} [3],
 59 with challenges include fast motion, full occlusion, hard distractor, *etc.* The results demonstrate the
 60 great discriminative power of our tracker. Many trackers will show a multi-peak on the response
 61 map when the object is occluded or multiple similar objects exist. With the vision-motion integrated
 62 Transformer architecture, our tracker eases such phenomenon.

63 References

- 64 [1] Bhat, G., Danelljan, M., Gool, L.V., Timofte, R., 2019. Learning discriminative model prediction for
65 tracking, in: ICCV.
- 66 [2] Chen, X., Yan, B., Zhu, J., Wang, D., Yang, X., Lu, H., 2021. Transformer tracking, in: CVPR.
- 67 [3] Fan, H., Bai, H., Lin, L., Yang, F., Chu, P., Deng, G., Yu, S., Huang, M., Liu, J., Xu, Y., et al., 2021. Lasot:
68 A high-quality large-scale single object tracking benchmark. IJCV 129, 439–461.
- 69 [4] Fan, H., Lin, L., Yang, F., Chu, P., Deng, G., Yu, S., Bai, H., Xu, Y., Liao, C., Ling, H., 2019. Lasot: A
70 high-quality benchmark for large-scale single object tracking, in: CVPR.
- 71 [5] Fu, Z., Liu, Q., Fu, Z., Wang, Y., 2021. Stmtrack: Template-free visual tracking with space-time memory
72 networks, in: CVPR.
- 73 [6] Ke, G., He, D., Liu, T.Y., 2021. Rethinking positional encoding in language pre-training, in: International
74 Conference on Learning Representations. URL: <https://openreview.net/forum?id=09-528y2Fgf>.
- 75 [7] Shaw, P., Uszkoreit, J., Vaswani, A., 2018. Self-attention with relative position representations. arXiv .
- 76 [8] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.,
77 2017. Attention is all you need, in: NeurIPS.
- 78 [9] Wang, N., Zhou, W., Wang, J., Li, H., 2021. Transformer meets tracker: Exploiting temporal context for
79 robust visual tracking, in: CVPR.
- 80 [10] Yan, B., Peng, H., Fu, J., Wang, D., Lu, H., 2021. Learning spatio-temporal transformer for visual tracking,
81 in: ICCV.

Figure 1: Comparison with state-of-the-art trackers on LaSOT [4] Test set using success (SUC) AUC score.

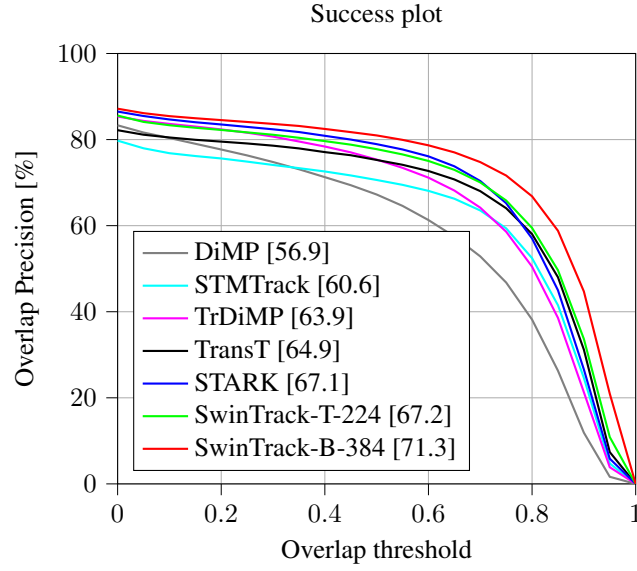


Figure 2: Comparison with state-of-the-art trackers on LaSOT [4] Test set using precision (PRE) AUC score.

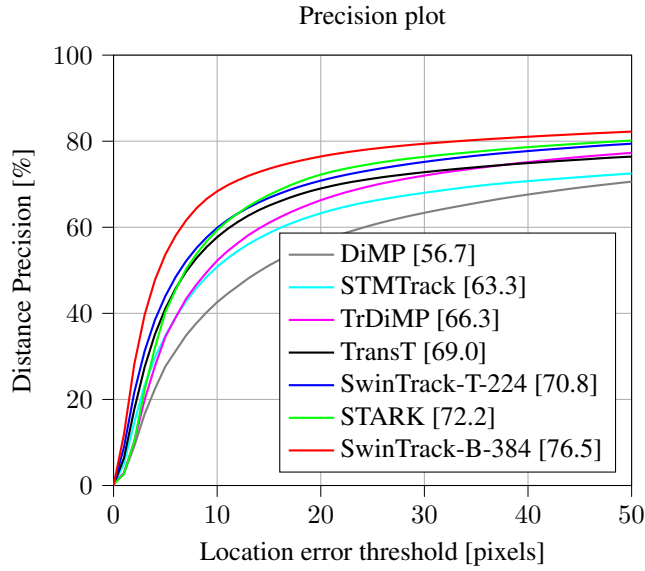


Figure 3: Heatmap visualization of the tracking response map of our SwinTrack-B-384 on LaSOT_{ext} [3]. The odd rows visualize the search region patches with ground-truth bounding box (in red rectangles). The even rows visualize the search region patches blended with the heatmap visualization of the response map.

