

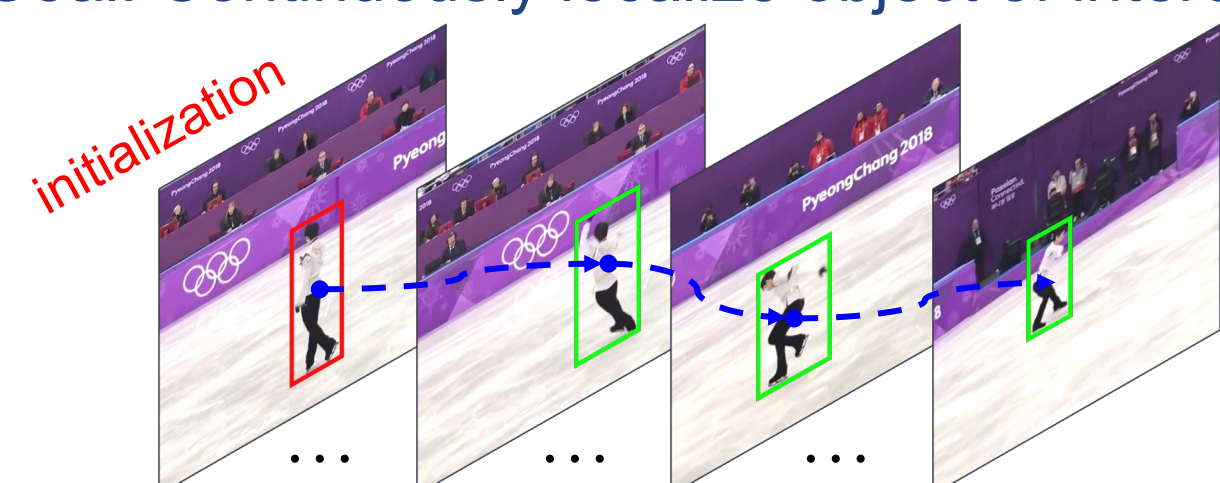


Code/Results

## Introduction

### Visual Object Tracking

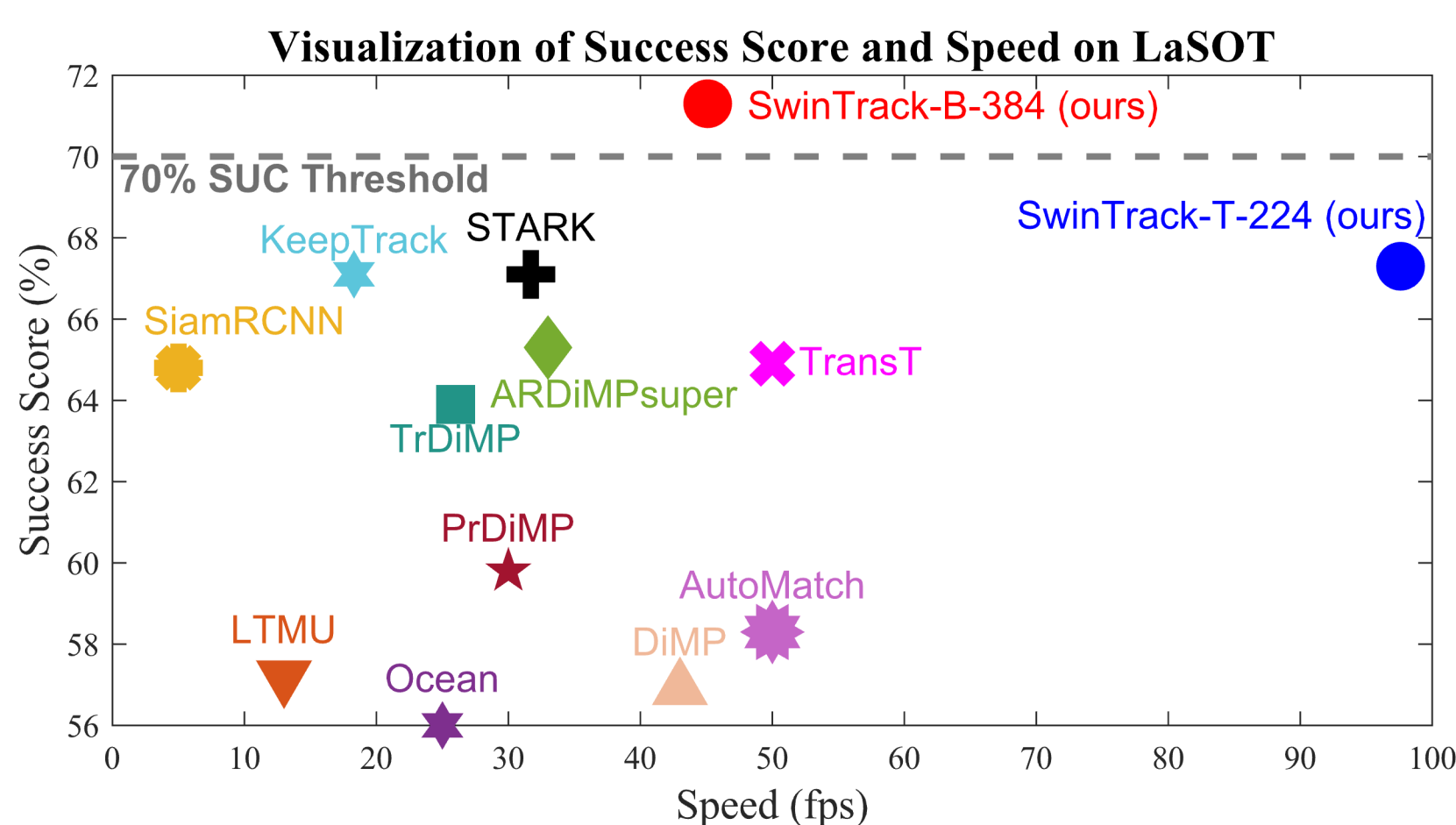
- Goal: Continuously localize object of interest in a video



### Motivation

- Transformer has greatly improved tracking performance
- Existing approaches usually adopt a hybrid CNN-Transformer architecture, i.e., CNN for feature extraction and Transformer for feature fusion
- A pure Transformer-based tracking architecture<sup>3</sup>, including Transformer-based feature extraction and fusion, is desired
- Motion information is crucial for temporal visual tracking.
- Swin Transformer shows SoTA results on various tasks.

## Contributions



- A simple but strong baseline, **SwinTrack**, is proposed with pure Transformer architecture
- We present a simple yet effective motion token in SwinTrack to enhance the robustness
- We conduct empirical studies on different components of SwinTrack, offering guidance for future tracker design
- SwinTrack shows SoTA results on multiple benchmarks, especially setting a new record with **0.713 SUC score** on the challenging LaSOT

## SwinTrack Framework

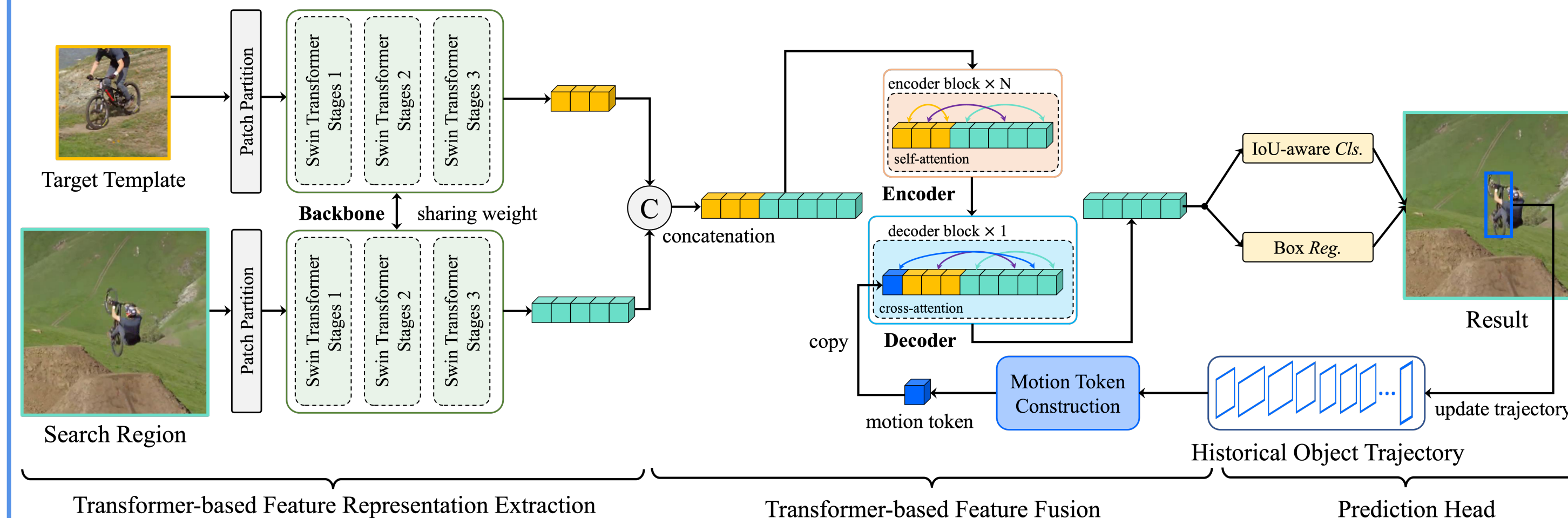


Figure: Architecture of SwinTrack.

### SwinTrack: Siamese Tracking via Vision-Motion Transformer

- Transformer-based Feature Representation Extraction
  - Swin Transformer as backbone*
  - Shared weight to extract template and search tokens*
- Transformer-based Feature Fusion
  - Encoder for fusing template and search tokens*
    - Joint vision representation learning
 
$$f_z^l, f_x^l = \text{DeConcat} \left( \text{SelfAttn} \left( \text{Concat} \left( f_z^{l-1}, f_x^{l-1} \right) \right) \right)$$
  - Decoder for fusing vision and motion information*
    - Motion token - Embedding of historical object trajectory
 
$$E_{\text{motion}} = \text{Concat} \left( E_{s(1)} + E_{s(2)} + \dots + E_{s(n)} \right)$$
    - Concatenation of object past bounding box embedding
 
$$s(i) = \max(t - i \times \Delta, 1)$$
    - Vision-motion representation learning
 
$$f_{vm} = \text{CrossAttn} \left( \text{Concat} \left( E_{\text{motion}}, f_z^l, f_x^l \right) \right)$$
    - Untied positional encoding with multi-dim multi-stream extension*
- Prediction Head & Loss Function
  - Response map generation by classification branch*
    - Three-layer perceptron with IoU-aware classification score
  - Box regression map generation by regression branch*
    - Three-layer perceptron with GloU loss

## Experiments

Table: Comparison with state-of-the-arts.

Tracker	LaSOT		LaSOT <sub>ext</sub>		TrackingNet		GOT-10k		TNL2k		
	SUC	P	SUC	P	SUC	P	AO	SR <sub>0.5</sub>	SR <sub>0.75</sub>	SUC	P
C-RPN	45.5	44.3	27.5	32.0	66.9	61.9	-	-	-	-	-
SiamPRN++	49.6	49.1	34.0	39.6	73.3	69.4	51.7	61.1	72.1	47.3	41.2
Ocean	56.0	56.6	-	-	-	-	-	-	-	-	-
DiMP	56.9	56.7	39.2	45.1	74.0	68.7	61.1	71.7	49.2	44.7	43.4
LTMU	57.2	57.2	41.4	47.3	-	-	-	-	-	48.5	47.3
SiamR-CNN	64.8	-	-	-	81.2	80.0	64.9	72.8	59.7	52.3	52.8
STMTrack	60.6	63.3	-	-	80.3	76.7	64.2	73.7	57.5	-	-
AutoMatch	58.3	59.9	37.6	43.0	76.0	72.6	65.2	76.6	54.3	-	-
TrDiMP	63.9	61.4	-	-	78.4	73.1	67.1	77.7	58.3	-	-
TransT	64.9	69.0	-	-	81.4	80.3	67.1	76.8	60.9	51.0	-
STARK	67.1	-	-	-	82.0	-	68.8	78.1	64.1	-	-
KeepTrack	67.1	70.2	48.2	-	-	-	-	-	-	-	-
SwinTrack-T-224	67.2	70.8	47.6	53.9	81.1	78.4	71.3	81.9	64.5	53.0	53.2
SwinTrack-B-384	71.3	76.5	49.1	55.6	84.0	82.8	72.4	80.5	67.8	55.9	57.1

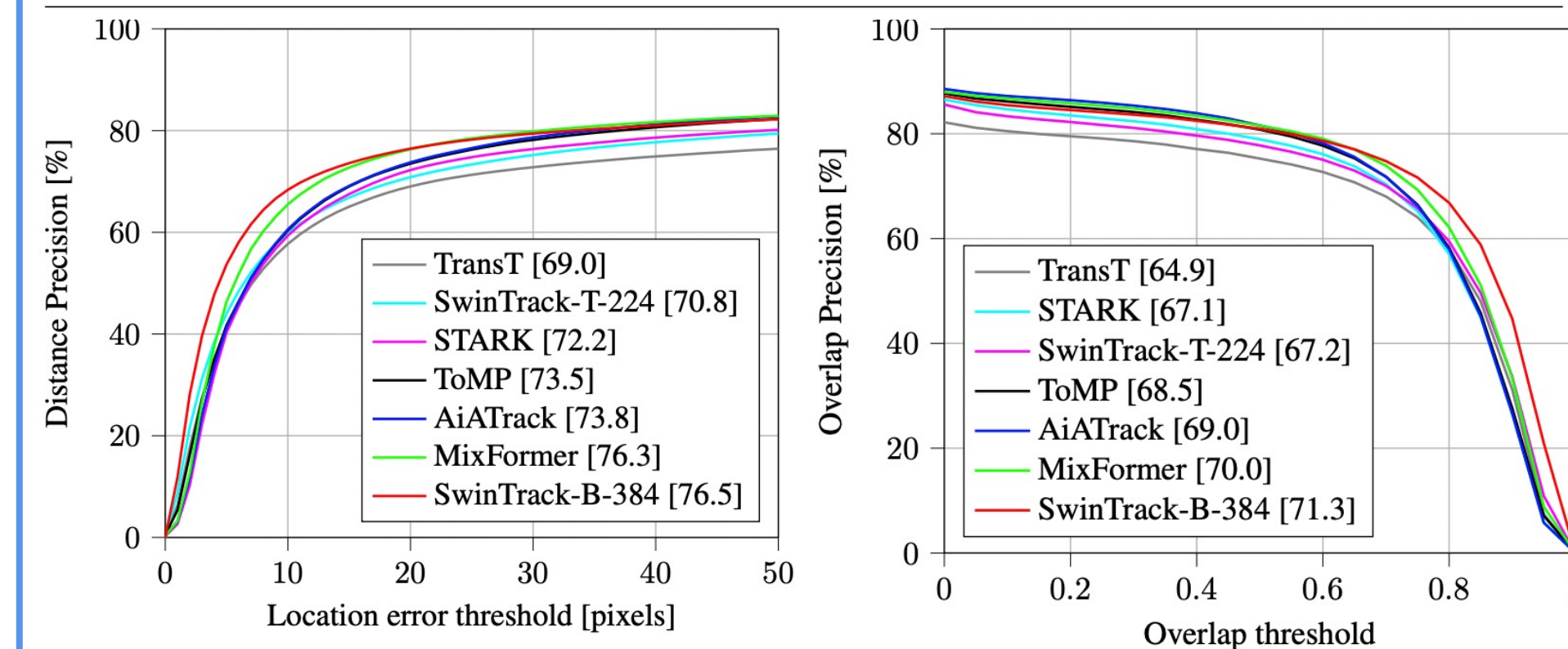


Figure: Comparison with latest Transformer-based trackers on LaSOT.

Table: Ablations with on SwinTrack-T-224 w/o motion token. ①: baseline; ②: replacing Transformer backbone w. ResNet-50; ③: replacing feature fusion w. cross attention-based fusion; ④: replacing decoder w. a target query-based; ⑤: replacing united positional encoding w. absolute sine position encoding; ⑥: replacing IoU-aware classification loss w. plain binary cross entropy loss; ⑦: removing the Hanning penalty window in inference.

	LaSOT SUC (%)	LaSOT <sub>ext</sub> SUC (%)	TrackingNet SUC (%)	GOT-10k mAO (%)	Speed fps	Params M
①	66.7	46.9	80.8	70.9	98	22.7
②	64.2	41.8	79.5	68.2	121	20.0
③	66.6	45.4	80.2	69.3	72	34.6
④	66.6	43.2	79.6	69.0	91	25.3
⑤	65.7	45.0	80.0	70.0	103	21.6
⑥	66.2	46.7	79.4	68.2	98	22.7
⑦	65.7	46.0	80.0	69.6	98	22.7

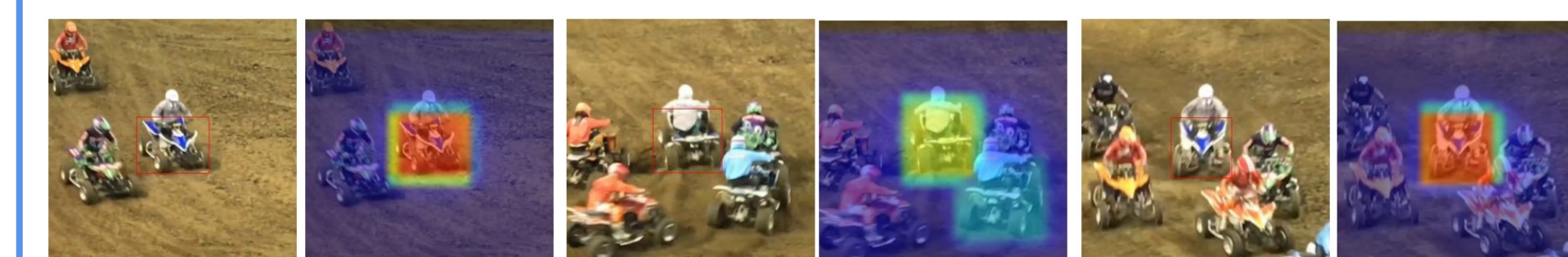


Figure: Visualization of tracking response maps of SwinTrack.

### Key References

- Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, B. Guo, Swin transformer: Hierarchical vision transformer using shifted windows, *ICCV*, 2021.
- X. Chen, B. Yan, J. Zhu, D. Wang, X. Yang, H. Lu, Transformer tracking, *CVPR*, 2021.
- A. Dosovitskiy, et al., An Image is worth 16x16 words: Transformers for image recognition at scale, *ICLR*, 2021.
- A. Vaswani. et al., Attention Is all you need, *NIPS*, 2017.