

# UTPTrack: Towards Simple and Unified Token Pruning for Visual Tracking

## Supplementary Material

To complement the manuscript, this appendix presents additional content and is organized as follows:

- Section A: Introduction of benchmarks.
- Section B: Practical efficiency.
- Section C: Overall performance (Sec. 4.2).
- Section D: Controlled-budget performance (Sec. 4.2).
- Section E: Additional ablation study and analysis.
- Section F: Visualization.

### A. Introduction of Benchmarks

#### A.1. RGB-based Tracking Benchmarks

**LaSOT.** LaSOT [9] is a densely annotated, class-balanced large-scale long-term tracking benchmark. The benchmark consists of high-quality 1400 training set videos and 280 test set videos with an average video length of approximately 2500 frames, containing challenges of various complex scenarios. The evaluation metrics include Success (AUC) and Precision (P and  $P_{\text{Norm}}$ ), among which AUC serves as the primary metric.

**LaSOT<sub>ext</sub>.** LaSOT<sub>ext</sub> [10] extends the LaSOT benchmark to include 150 videos of 15 additional new object classes, which were carefully selected outside of ImageNet. The evaluation metrics remain consistent with LaSOT.

**TrackingNet.** TrackingNet [30] is a large-scale short-term tracking benchmark containing over 30K video sequences and over 14M bounding box annotations. The test set of TrackingNet contains 511 video sequences without publicly available ground truth. Researchers evaluate trackers by submitting their results to the official online evaluation server, which provides AUC, P and  $P_{\text{Norm}}$ .

**GOT-10k.** GOT-10K [16] is a challenging, high-diversity, large-scale benchmark with more than 10K video sequences, covering most of more than 560 categories of moving objects and more than 80 categories of motion patterns. The test set has 180 video sequences, including 84 object categories and 32 motion categories. It is worth noting that GOT-10k advocates a one-shot protocol, which means that the object classes between the train set and the test set are zero-overlapped. Similar to TrackingNet, the tracking results need to be submitted to the official evaluation server to obtain the tracker’s performance. The evaluation metrics include Average Overlap (AO) and Success Rates at thresholds 0.5 and 0.75 ( $\text{SR}_{0.5}$  and  $\text{SR}_{0.75}$ ), among which AO serves as the primary metric.

#### A.2. Multimodal Tracking Benchmarks

**VOT-RGBD22.** VOT-RGBD22 [18] is a dataset designed for evaluating RGB-Depth (RGB-D) tracking algorithms, consisting of 127 video sequences. The evaluation follows the VOT protocol, including Accuracy, Robustness, and Expected Average Overlap (EAO), with EAO as the primary performance metric.

**LasHeR.** LasHeR [22] is a large-scale RGB-D dataset focused on long-term object tracking, containing 245 video sequences in its test set. The evaluation metrics include AUC and P, with the AUC as the primary metric.

**RGBT234.** RGBT234 [21] is an RGB-Thermal (RGB-T) tracking dataset, comprising 234 test video sequences that cover challenging scenarios such as illumination changes, occlusions, and low visibility. The evaluation metrics include Maximum Success Rate (MSR) and Maximum Precision Rate (MPR).

**VisEvent.** VisEvent [36] is a large-scale RGB-Event (RGB-E) object tracking dataset comprising 320 test video sequences. It emphasizes tracking performance under high-speed motion and low-light conditions. Researchers evaluate performance using the AUC and P metrics.

**TNL2K.** TNL2K [35] is a language-guided visual tracking dataset with 650 video sequences in the test set and covers a wide range of natural language descriptions. Researchers evaluate performance using AUC and P, with AUC as the primary metric.

**OTB99.** OTB99 [23] is a classic benchmark for short-term visual object tracking, with the test set consisting of 48 video sequences. Researchers assess tracking performance based on AUC and P scores.

### B. Practical Efficiency Comparisons

As shown in Table 8, UTPTrack achieves higher FPS and lower per-layer backbone latency on both GPU and CPU across resolutions for OSTrack and SUTrack, while reducing training hours via lower backbone computation.

Table 8. Efficiency comparison. Lat.: per-layer backbone latency. GPU: NVIDIA 1080Ti; CPU: Intel Xeon Gold 6226R@2.90GHz.

Model	Train Hours	FPS (GPU)	FPS (CPU)	Lat. (GPU, ms)	Lat. (CPU, ms)
OSTrack <sub>256</sub>	38.3	94.0	9.0	15.0	66.5
UTPTrack-O <sub>256</sub>	32.6	95.0	16.7	12.7	49.5
OSTrack <sub>384</sub>	101.4	39.8	3.2	36.1	340.8
UTPTrack-O <sub>384</sub>	80.6	47.3	6.0	30.2	212.6
SUTrack <sub>224</sub>	56.8	40.5	7.8	17.4	66.4
UTPTrack-S <sub>224</sub>	53.3	42.7	9.3	17.8	56.4
SUTrack <sub>384</sub>	209.2	27.4	3.4	63.3	340.2
UTPTrack-S <sub>384</sub>	186.6	30.7	4.6	53.0	261.2

Table 9. SOTA comparison on RGB-based tracker across RGB-Based tracking.

Method	LaSOT			LaSOT <sub>ext</sub>			TrackingNet			GOT-10k		
	AUC	P <sub>Norm</sub>	P	AUC	P <sub>Norm</sub>	P	AUC	P <sub>Norm</sub>	P	AO	SR <sub>0.5</sub>	SR <sub>0.75</sub>
<i>Trackers with Resolution 256</i>												
OTrack <sub>256</sub> (Baseline)	67.9	76.9	73.3	47.6	57.5	53.6	84.5	89.0	83.0	74.9	84.2	73.0
OTrack-CE <sub>256</sub>	67.4	76.5	<u>72.8</u>	<u>46.6</u>	<u>56.5</u>	52.3	<u>83.7</u>	88.3	82.4	<u>75.5</u>	<u>85.5</u>	73.8
OTrack-ToMe <sub>256</sub>	<u>67.6</u>	<u>76.6</u>	72.6	46.5	56.0	<u>52.5</u>	<b>84.0</b>	88.6	82.6	73.4	83.4	72.4
OTrack-EViT <sub>256</sub>	<b>68.4</b>	<b>77.7</b>	<b>74.1</b>	46.1	55.7	51.1	<b>84.0</b>	<u>88.7</u>	<b>82.9</b>	<b>75.8</b>	<b>85.9</b>	<b>74.6</b>
<b>UTPTrack-O<sub>256</sub></b> (ours)	67.4	76.3	72.5	<b>47.3</b>	<b>57.1</b>	<b>53.3</b>	<b>84.0</b>	<b>88.8</b>	<u>82.8</u>	75.2	85.3	<u>73.9</u>
<i>Trackers with Higher Resolution 384</i>												
OTrack <sub>384</sub> (Baseline)	70.7	80.6	77.2	50.0	61.1	56.8	84.0	88.7	83.0	76.6	86.9	75.3
OTrack-CE <sub>384</sub>	<u>70.6</u>	<u>80.1</u>	<b>77.0</b>	<u>50.6</u>	<u>61.2</u>	<u>57.0</u>	<u>84.1</u>	<u>88.8</u>	83.0	76.3	86.5	<u>74.9</u>
OTrack-ToMe <sub>384</sub>	<u>70.2</u>	<u>80.1</u>	76.6	50.4	61.1	56.6	<u>84.1</u>	<u>88.8</u>	<u>83.1</u>	<b>77.1</b>	<b>87.2</b>	74.7
OTrack-EViT <sub>384</sub>	70.1	79.9	76.4	<b>52.1</b>	<b>62.9</b>	<b>59.2</b>	<b>84.3</b>	<b>89.0</b>	<b>83.3</b>	76.4	86.4	<b>75.3</b>
<b>UTPTrack-O<sub>384</sub></b> (ours)	<b>70.7</b>	<b>80.6</b>	<b>77.1</b>	49.4	59.8	55.8	84.0	88.7	82.7	<u>76.5</u>	<u>86.9</u>	74.5

Table 10. SOTA comparison on unified tracker across RGB-Based tracking.

Method	LaSOT			LaSOT <sub>ext</sub>			TrackingNet			GOT-10k		
	AUC	P <sub>Norm</sub>	P	AUC	P <sub>Norm</sub>	P	AUC	P <sub>Norm</sub>	P	AO	SR <sub>0.5</sub>	SR <sub>0.75</sub>
<i>Trackers with Resolution 224</i>												
SUTrack <sub>224</sub> (Baseline)	73.7	83.8	80.7	53.2	64.5	61.6	85.9	90.4	85.7	77.5	86.7	78.1
SUTrack-CE <sub>224</sub>	<b>72.8</b>	<b>82.9</b>	<b>80.2</b>	<u>52.9</u>	<u>63.7</u>	<u>60.4</u>	<b>85.3</b>	<b>90.0</b>	<b>84.8</b>	<b>77.6</b>	<b>87.6</b>	<b>78.0</b>
SUTrack-ToMe <sub>224</sub>	71.4	81.2	78.3	52.0	62.9	59.3	<u>85.2</u>	<b>90.0</b>	84.6	<u>77.5</u>	<u>87.4</u>	77.5
SUTrack-EViT <sub>224</sub>	71.3	81.9	78.3	51.3	62.7	58.5	84.5	<u>89.5</u>	83.5	76.4	86.5	75.7
SUTrack-DynamicViT <sub>224</sub>	57.4	65.5	59.3	40.8	49.7	46.6	65.6	69.1	60.2	71.6	81.2	64.7
<b>UTPTrack-S<sub>224</sub></b> (ours)	<u>72.6</u>	<u>82.8</u>	<u>79.8</u>	<b>53.6</b>	<b>64.8</b>	<b>61.4</b>	<u>85.2</u>	<b>90.0</b>	<u>84.7</u>	77.3	87.1	<u>77.7</u>
<i>Trackers with Higher Resolution 384</i>												
SUTrack <sub>384</sub> (Baseline)	73.9	83.2	81.2	52.8	63.3	60.3	86.8	90.9	87.3	79.1	87.6	80.2
SUTrack-CE <sub>384</sub>	73.0	82.3	80.2	<b>54.1</b>	<b>64.7</b>	<b>61.9</b>	86.2	90.6	86.6	<b>79.7</b>	<b>89.0</b>	<u>80.1</u>
SUTrack-ToMe <sub>384</sub>	<u>74.2</u>	<u>83.8</u>	<b>81.9</b>	53.0	63.3	60.5	<u>86.3</u>	<u>90.7</u>	<u>86.8</u>	79.1	88.1	79.1
SUTrack-EViT <sub>384</sub>	73.7	83.4	81.6	53.0	63.7	61.0	86.1	90.6	86.4	78.4	87.9	78.1
SUTrack-DynamicViT <sub>384</sub>	63.0	72.1	63.0	44.8	53.4	47.3	73.5	80.5	66.3	69.5	79.0	61.8
<b>UTPTrack-S<sub>384</sub></b> (ours)	<b>74.3</b>	<b>83.9</b>	<u>81.8</u>	<u>53.6</u>	<u>64.3</u>	<u>61.5</u>	<b>86.4</b>	<b>90.9</b>	<b>86.9</b>	<u>79.3</u>	<u>88.3</u>	<b>80.3</b>

### C. Detailed Overall Performance Comparisons

In Sec. 4.2, we validate the effectiveness of UTPTrack under two tracking paradigms: RGB-based tracking and Unified tracking, which includes RGB-Depth (RGB-D), RGB-Thermal (RGB-T), RGB-Event (RGB-E), and RGB-Language (RGB-Lang) tasks. Using OTrack and SUTrack as representative base models, we compare UTPTrack against state-of-the-art token compression methods, including CE, ToMe, EViT, and DynamicViT, across varying input resolutions. We report comprehensive results for Overall Performance, with detailed tabular data and analysis provided in this section.

**RGB-based Tracking.** We evaluate UTPTrack on four large-scale RGB-based tracking benchmarks, covering both long-term (LaSOT, LaSOT<sub>ext</sub>) and short-term (TrackingNet, GOT-10k) scenarios. The results are summarized in Tab. 9 and Tab. 10. For the RGB-based variant, UTPTrack-O maintains the accuracy of base model across all four benchmarks and both resolutions. At 256 resolution, UTPTrack-O<sub>256</sub> closely matches base model on LaSOT, LaSOT<sub>ext</sub>, and TrackingNet while slightly improving AO on GOT-10k, and remains competitive with other compression methods. At 384 resolution, UTPTrack-O<sub>384</sub> similarly matches the performance of base model across all benchmarks. In con-

Table 11. SOTA comparison on unified tracker across RGB-Depth, RGB-Thermal, RGB-Event, and RGB-Language tracking.

Method	VOT-RGBD22			LasHeR		RGBT234		VisEvent		TNL2K		OTB99	
	EAO	Acc	Rob	AUC	P	MSR	MPR	AUC	P	AUC	P	AUC	P
<i>Trackers with Resolution 224</i>													
SUTrack <sub>224</sub> (Baseline)	75.5	82.5	91.3	59.9	74.8	70.0	92.1	63.3	80.7	67.8	73.8	68.4	91.1
SUTrack-CE <sub>224</sub>	<b>76.7</b>	<b>82.7</b>	<b>92.4</b>	<b>59.6</b>	<b>74.0</b>	<u>70.1</u>	<u>92.4</u>	<b>62.8</b>	<u>79.9</u>	<u>65.7</u>	<u>71.1</u>	<u>70.8</u>	<u>92.7</u>
SUTrack-ToMe <sub>224</sub>	76.1	<b>82.7</b>	91.8	<u>59.2</u>	73.7	69.4	91.8	61.7	78.9	<b>65.9</b>	<b>71.4</b>	70.5	91.6
SUTrack-EViT <sub>224</sub>	75.0	82.1	90.9	58.6	73.0	69.2	91.6	62.0	79.4	65.0	70.0	<u>70.8</u>	91.9
SUTrack-DynamicViT <sub>224</sub>	62.8	74.2	84.4	46.1	57.5	63.2	85.1	50.9	71.8	55.6	55.0	65.8	86.8
<b>UTPTrack-S<sub>224</sub>(ours)</b>	<u>76.4</u>	<b>82.7</b>	<u>92.1</u>	<u>59.2</u>	<u>73.9</u>	<b>70.4</b>	<b>92.9</b>	<u>62.6</u>	<b>80.0</b>	65.6	71.0	<b>71.9</b>	<b>93.7</b>
<i>Trackers with Higher Resolution 384</i>													
SUTrack <sub>384</sub> (Baseline)	76.0	83.2	91.5	60.6	75.3	69.2	92.4	63.0	79.3	68.2	74.7	69.2	91.0
SUTrack-CE <sub>384</sub>	<b>77.9</b>	83.6	<b>92.8</b>	<b>60.4</b>	<b>75.0</b>	<b>69.7</b>	<u>92.0</u>	<b>62.9</b>	<u>79.5</u>	66.2	72.6	70.0	90.6
SUTrack-ToMe <sub>384</sub>	<u>76.8</u>	<b>83.8</b>	91.5	<u>59.7</u>	74.1	<u>69.4</u>	<u>92.0</u>	62.6	79.3	<u>66.4</u>	<u>72.8</u>	70.8	91.8
SUTrack-EViT <sub>384</sub>	77.1	83.2	<u>92.5</u>	58.8	73.2	68.2	90.8	62.1	79.2	66.2	72.7	<b>72.5</b>	<b>94.6</b>
SUTrack-DynamicViT <sub>384</sub>	71.2	81.0	87.8	51.6	63.5	56.7	82.2	56.4	73.4	58.9	58.0	68.0	86.5
<b>UTPTrack-S<sub>384</sub>(ours)</b>	<u>76.8</u>	<b>83.8</b>	<u>91.6</u>	<b>60.4</b>	<u>74.8</u>	<b>69.7</b>	<b>92.2</b>	<u>62.8</u>	<b>79.7</b>	<b>66.6</b>	<b>72.9</b>	<u>72.2</u>	<u>93.7</u>

trast, the unified variant, UTPTrack-S, consistently delivers strong results across all benchmarks. UTPTrack-S<sub>224</sub> improves the AUC on LaSOT<sub>ext</sub> and TrackingNet over base model, while UTPTrack-S<sub>384</sub> achieves the best or second-best scores on all four datasets. These results demonstrate that our pruning strategy effectively preserves accuracy across different resolutions and tracking paradigms.

**RGB-D/T/E Tracking.** As shown in Tab.11, UTPTrack-S<sub>224</sub> and UTPTrack-S<sub>384</sub> consistently rank among the top two across all benchmarks, including VOT-RGBD22 (RGB-D), LasHeR and RGBT234 (RGB-T), and VisEvent (RGB-E). This demonstrates the robustness and generalizability of our pruning strategy across diverse modalities and dataset scales, while maintaining strong performance and high efficiency.

**RGB-Lang Tracking.** As shown in Tab. 11, UTPTrack-S significantly outperforms the baseline model on the short-term tracking dataset OTB99 across different resolutions. It achieves gains of 3.5% at 224 resolution and 2.0% at 384 resolution. On TNL2K, however, the performance of UTPTrack-S decreases relative to the baseline model, likely due to compression-induced performance loss. Even so, UTPTrack-S still delivers competitive performance on TNL2K, and at 384 resolution it achieves the best results among all compared methods. Such degradation across all compressed models may be attributed to the limited precision of the initial language descriptions in later stages of long videos, especially in long-term sequences, which reduces the reliability of language-guided cues.

## D. Detaild Controlled-Budget Performance Comparisons

In Sec. 4.2, we conduct controlled-budget experiments under fixed compression ratios to ensure fair comparisons for RGB-based and unified tracking, with comprehensive results presented in this section.

**RGB-based Tracking.** For RGB-based tracking, we compare UTPTrack with other methods on four widely used RGB-based benchmarks. As shown in Tab. 12, UTPTrack-O consistently outperforms all other methods across all pruning rates. Notably, with 18.8% of the visual tokens pruned, UTPTrack-O<sub>256</sub> outperforms the baseline model by 0.2%, suggesting that removing redundant tokens can even lead to performance gains. Similarly, in Tab. 16, UTPTrack-S surpasses all other compression methods across all pruning rates, and at a retained ratio of 35.4% visual tokens, it achieves a 1.0% improvement over the second-best method.

**RGB-D/T/E Tracking.** As shown in Tab. 17, UTPTrack outperforms the baseline on both the RGB-D and RGB-T benchmarks across different pruning ratios. On the VOT-RGBD22 dataset, it surpasses the baseline by 1.3% when pruning 25.6% visual tokens. On the RGB-E benchmark, UTPTrack consistently ranks among the top two methods. These results demonstrate that UTPTrack generalizes effectively across diverse multimodal tracking scenarios.

**RGB-Lang Tracking.** As shown in Tab. 17, UTPTrack maintains competitive performance across different compression rates on the RGB-Language benchmark. At all compression rates, UTPTrack outperforms the baseline on

Table 12. Performance comparisons under different vision token compression settings for **RGB-based tracking**. The baseline uses 384 tokens. The first row shows raw scores, and the second row reports percentages relative to the upper bound.

Method	LaSOT			LaSOT <sub>ext</sub>			TrackingNet			GOT-10k			Avg Perf.(%)
	AUC	P <sub>Norm</sub>	P	AUC	P <sub>Norm</sub>	P	AUC	P <sub>Norm</sub>	P	AO	SR <sub>0.5</sub>	SR <sub>0.75</sub>	
<i>Upper Bound, 384 Tokens (100%)</i>													
OSTrack <sub>256</sub> (Baseline)	67.9	76.9	73.3	47.6	57.5	53.6	84.5	89.0	83.0	74.9	84.2	73.0	100%
<i>Retain 87.2% Tokens (≈ 335 tokens, ↓ 18.8%)</i>													
OSTrack-CE <sub>256</sub>	67.1	75.8	72.0	46.9	56.7	52.4	84.1	88.8	83.0	74.9	85.0	73.5	<u>99.3%</u> (↓ 0.7%)
OSTrack-ToMe <sub>256</sub>	68.0	76.9	72.9	46.4	56.3	51.5	83.8	88.7	82.7	74.3	84.2	73.0	99.0%(↓ 1.0%)
OSTrack-EViT <sub>256</sub>	66.8	75.7	71.7	46.7	56.3	52.3	84.1	88.9	82.8	75.0	85.2	73.9	99.0%(↓ 1.0%)
UTPTrack-O <sub>256</sub>	67.1	75.7	72.2	48.3	58.3	54.4	84.0	88.7	83.0	75.6	85.8	74.7	<b>100.2%</b> (↑ 0.2%)
<i>Retain 75.5% Tokens (≈ 290 tokens, ↓ 24.5%)</i>													
OSTrack-CE <sub>256</sub>	67.3	76.2	72.6	47.3	57.2	53.1	84.2	89.0	83.2	74.2	84.1	72.8	99.3%(↓ 0.7%)
OSTrack-ToMe <sub>256</sub>	68.2	77.1	73.2	46.5	56.0	52.2	84.2	89.0	83.2	74.9	85.1	73.6	<u>99.4%</u> (↓ 0.6%)
OSTrack-EViT <sub>256</sub>	67.8	76.8	73.4	47.1	56.6	52.8	83.8	88.5	82.5	72.9	82.8	70.4	98.8%(↓ 1.2%)
UTPTrack-O <sub>256</sub>	67.6	76.3	72.9	47.3	57.2	53.0	83.8	88.4	82.7	74.8	85.1	73.5	<b>99.5%</b> (↓ 0.5%)
<i>Retain 65.6% Tokens (≈ 252 tokens, ↓ 34.4%)</i>													
OSTrack-CE <sub>256</sub>	67.2	76.1	72.5	46.5	56.3	51.8	83.3	88.2	82.0	72.8	82.9	70.9	<u>98.1%</u> (↓ 1.9%)
OSTrack-ToMe <sub>256</sub>	66.3	75.7	71.1	46.0	55.9	50.9	82.5	87.5	80.2	71.0	81.3	67.9	96.7%(↓ 3.3%)
OSTrack-EViT <sub>256</sub>	68.1	77.2	73.2	46.1	55.8	51.2	83.0	87.9	81.7	72.8	83.3	69.7	<u>98.1%</u> (↓ 1.9%)
UTPTrack-O <sub>256</sub>	68.2	77.3	74.0	46.2	56.1	51.8	84.0	88.9	83.1	74.9	85.0	73.3	<b>99.2%</b> (↓ 0.8%)

the short-term tracking dataset OTB99. With 25.6% visual tokens pruned, UTPTrack surpasses the baseline by 3.3%, and even at a pruning rate of 64.6%, it still leads the baseline by 2.7%. On TNL2K, UTPTrack consistently ranks in the top two compared to other compression methods.

## E. Additional Ablation Study

### E.1. The Effect of Pruning Location for RGB-based Tracking

To analyze the impact of pruning locations, we systematically vary the pruning positions of both the CE and DTE modules across different transformer layers, with detailed results summarized in Tab. 13. Based on empirical results and a meticulous consideration of the performance-efficiency trade-off, we select the pruning configuration

(#3) that prunes CE at layers [3, 6, 9], and DTE at layers [4, 7, 10]. This arrangement achieves an optimal balance, maintaining competitive tracking performance while substantially reducing the computational overhead.

Table 13. Ablation Study on CTEM Location for **RGB-based Trackers** with a 12-layer ViT backbone.

#	CE Location	DTE Location	LaSOT	LaSOT <sub>ext</sub>	TrackingNet	GOT-10k	Avg Vis Tok	Cm Vis Tok	MACs (G)	Avg Perf. (%)
1	-	-	67.9	47.6	84.5	74.9	384.0	384	34.5	100.0
2	[3,6,9]	[3,6,9]	67.6	47.0	83.8	74.5	267.8	176	25.1	99.3
3	[3,6,9]	[4,7,10]	68.4	46.6	83.8	75.3	271.2	176	25.4	<b>99.6</b>
4	[3,6,9]	[2,5,8]	67.4	45.7	83.8	75.4	264.3	176	24.8	98.8
5	[2,5,8]	[3,6,9]	68.0	46.4	83.9	75.7	253.8	176	23.9	<u>99.5</u>
6	[4,7,10]	[3,6,9]	67.5	46.5	83.9	73.4	281.7	176	26.3	98.6

## E.2. The Effect of Spatial Priors in Token Type-Aware Strategy

TTA does not replace attention-based importance estimation, but stabilizes static template pruning, which is sensitive due to limited tokens and foreground-background imbalance. The bounding-box prior uses only initialization information. As shown in Table 14, attention-based pruning without priors consistently outperforms random pruning, while adding spatial priors yields modest but consistent gains. These results indicate that TTA complements attention-based pruning by improving stability under high compression ratios.

Table 14. Effect of spatial priors in attention-guided pruning.

Method	LaSOT	LaSOT <sub>ext</sub>	TrackingNet	GOT-10k
Attn-based Pruning w/o Priors	67.2	46.5	<b>84.3</b>	74.4
Attn-based Pruning w/ Priors	<b>67.4</b>	<b>47.3</b>	84.0	<b>75.2</b>
Random Pruning w/ Priors	67.0	46.8	83.9	74.2

## E.3. The Effect of Token Type-Aware Strategy for Unified Tracking

Table 15 compares the effects of three foreground bonus strategies on ST pruning for unified trackers, corresponding to the Full, Soft and All variants introduced in Sec. 3.2.1. Among them, the Soft bonus delivers the best overall performance across the five benchmarks, achieving an average accuracy of 99.8%, outperforming the Full and All strategies by approximately 0.5% and 0.7%, respectively. Under a strict pruning ratio, the Full bonus rewards only patches that lie entirely within the target box, making it overly conservative and prone to underestimating foreground tokens near object boundaries. Conversely, the All bonus assigns the maximum reward as long as a patch merely intersects the box, which is excessively permissive and easily introduces background noise. The Soft bonus, computed as the mean mask value within each patch, provides a more fine-grained estimate of foreground coverage and yields smoother boundary transitions, making it a more reliable continuous signal for guiding pruning.

Table 15. Ablation Study on bonus for Unified Trackers.

#	Bonus	LaSOT	VOT-RGBD22	LasHeR	VisEvent	OTB99	Avg Perf. (%)
1	Full	71.7	75.7	59.5	61.7	69.7	99.3
2	Soft	72.4	75.9	59.4	61.9	70.7	<b>99.8</b>
3	All	72.2	75.6	59.5	61.2	69.4	99.1

## F. Visualization

We present visualizations of the UTPTrack pruning process for RGB-based, RGB-Lang, RGB-D, RGB-T, and RGB-E tracking, as shown in Figs. 5 to 9. The top-left shows the search region image, with the static template on the bottom-left and the dynamic template on the bottom-right. Masked areas indicate pruned tokens. The six stages illustrate a progressive schedule: stages 1, 3, and 5 progressively prune tokens from the search region, while stages 2, 4, and 6 further prune tokens from both the static and dynamic templates.

Table 16. Performance comparisons under different vision token compression configurations across **unified tracking** on RGB-based tracking benchmark. The vanilla number of visual tokens is 294. The first line of each method is the raw performance of benchmarks, and the second line is the proportion relative to the upper limit. The average performance listed is calculated across all 10 benchmarks.

Method	LaSOT			LaSOT <sub>ext</sub>			TrackingNet			GOT-10k			Avg Perf.(%)
	AUC	P <sub>Norm</sub>	P	AUC	P <sub>Norm</sub>	P	AUC	P <sub>Norm</sub>	P	AO	SR <sub>0.5</sub>	SR <sub>0.75</sub>	
<i>Upper Bound, 294 Tokens (100%)</i>													
SUTrack <sub>224</sub> (Baseline)	73.7	83.8	80.7	53.2	64.5	61.6	85.9	90.4	85.7	77.5	86.7	78.1	100%
	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	
<i>Retain 71.4% Tokens (<math>\approx 218</math> tokens, <math>\downarrow 25.6\%</math>)</i>													
SUTrack-CE <sub>224</sub>	72.4	82.4	79.2	52.8	63.6	60.2	85.1	89.7	84.4	77.4	87.0	78.0	<u>99.5%</u> ( $\downarrow 0.5\%$ )
	98.3%	98.4%	98.1%	99.1%	98.7%	97.7%	99.0%	99.2%	98.5%	99.9%	100.3%	99.9%	
SUTrack-ToMe <sub>224</sub>	72.1	82.0	78.9	52.6	63.7	59.7	85.4	90.0	84.7	76.9	86.5	77.2	99.4%( $\downarrow 0.6\%$ )
	97.8%	97.8%	97.7%	98.8%	98.8%	96.9%	99.4%	99.6%	98.9%	99.2%	99.8%	98.8%	
SUTrack-EViT <sub>224</sub>	71.6	81.9	79.0	52.3	63.2	59.5	84.9	89.7	84.1	76.9	87.2	76.8	99.0%( $\downarrow 1.0\%$ )
	97.1%	97.7%	97.8%	98.3%	98.1%	96.6%	98.8%	99.2%	98.2%	99.2%	100.6%	98.3%	
SUTrack-DyViT <sub>224</sub>	68.8	81.6	71.1	49.7	63.0	55.3	80.2	86.4	74.9	74.6	87.1	72.1	96.5%( $\downarrow 3.5\%$ )
	93.4%	97.4%	88.0%	93.4%	97.8%	89.8%	93.4%	95.6%	87.4%	96.3%	100.5%	92.3%	
UTPTrack-S <sub>224</sub>	72.1	82.2	78.9	52.9	63.8	60.1	85.1	89.7	84.5	77.7	87.4	78.2	<b>99.8%</b> ( $\downarrow 0.2\%$ )
	97.9%	98.1%	97.7%	99.4%	98.9%	97.5%	99.0%	99.3%	98.6%	100.3%	100.8%	100.1%	
<i>Retain 52.0% Tokens (<math>\approx 153</math> tokens, <math>\downarrow 48.0\%</math>)</i>													
SUTrack-CE <sub>224</sub>	72.0	81.8	78.9	52.8	63.9	60.4	85.3	90.0	84.7	77.1	86.8	77.7	<u>99.2%</u> ( $\downarrow 0.8\%$ )
	97.7%	97.6%	97.7%	99.2%	99.1%	97.9%	99.3%	99.6%	98.8%	99.5%	100.1%	99.5%	
SUTrack-ToMe <sub>224</sub>	71.4	81.2	78.3	52.0	62.9	59.3	85.2	90.0	84.6	77.5	87.4	77.5	99.0%( $\downarrow 1.0\%$ )
	96.9%	96.9%	96.9%	97.8%	97.6%	96.2%	99.1%	99.6%	98.7%	100.0%	100.8%	99.2%	
SUTrack-EViT <sub>224</sub>	71.2	81.6	78.4	51.9	61.1	58.9	84.7	89.7	84.0	76.9	87.2	76.9	98.5%( $\downarrow 1.5\%$ )
	96.6%	97.4%	97.1%	97.5%	94.8%	95.6%	98.6%	99.3%	98.0%	99.2%	100.6%	98.5%	
SUTrack-DyViT <sub>224</sub>	53.2	60.6	55.3	36.5	44.1	41.7	58.2	62.4	52.2	52.2	60.0	34.8	78.8%( $\downarrow 21.2\%$ )
	72.1%	72.3%	68.5%	68.5%	68.4%	67.6%	67.7%	69.0%	61.0%	67.4%	69.2%	44.6%	
UTPTrack-S <sub>224</sub>	72.5	82.5	79.6	53.1	64.0	60.5	85.5	90.1	84.9	77.1	86.8	77.5	<b>99.5%</b> ( $\downarrow 0.5\%$ )
	98.4%	98.5%	98.6%	99.7%	99.3%	98.2%	99.5%	99.7%	99.1%	99.5%	100.1%	99.2%	
<i>Retain 35.4% Tokens (<math>\approx 104</math> tokens, <math>\downarrow 64.6\%</math>)</i>													
SUTrack-CE <sub>224</sub>	70.5	80.5	77.2	51.8	63.0	59.3	84.8	89.5	83.9	75.9	85.9	75.3	<u>98.3%</u> ( $\downarrow 1.7\%$ )
	95.7%	96.1%	95.6%	97.3%	97.7%	96.2%	98.7%	99.1%	97.9%	97.9%	99.1%	96.4%	
SUTrack-ToMe <sub>224</sub>	67.9	73.9	70.7	48.6	55.9	53.8	80.0	83.1	75.5	71.6	81.2	64.7	92.5%( $\downarrow 7.5\%$ )
	92.1%	88.1%	87.6%	91.3%	86.6%	87.2%	93.1%	92.0%	88.1%	92.4%	93.7%	82.8%	
SUTrack-EViT <sub>224</sub>	69.6	80.0	76.5	50.6	61.8	57.6	84.3	89.2	83.2	74.0	84.3	72.4	96.7%( $\downarrow 3.3\%$ )
	94.5%	95.4%	94.7%	95.1%	95.9%	93.4%	98.1%	98.8%	97.0%	95.5%	97.2%	92.7%	
SUTrack-DyViT <sub>224</sub>	4.4	3.6	3.6	2.9	5.8	2.5	8.8	6.7	5.7	7.4	7.1	0.8	14.7%( $\downarrow 85.3\%$ )
	5.9%	4.2%	4.4%	5.4%	9.0%	4.1%	10.3%	7.4%	6.7%	9.5%	8.2%	1.0%	
UTPTrack-S <sub>224</sub>	72.3	82.3	79.3	52.7	63.7	60.0	85.2	89.7	84.4	77.5	87.1	77.9	<b>99.3%</b> ( $\downarrow 0.7\%$ )
	98.2%	98.2%	98.2%	99.0%	98.8%	97.4%	99.1%	99.3%	98.5%	100.0%	100.5%	99.7%	

Table 17. Performance comparisons under different vision token compression configurations across **unified tracking** on RGB-D/T/E/Lang tracking benchmark. The vanilla number of visual tokens is 294. The first line of each method is the raw performance of benchmarks, and the second line is the proportion relative to the upper limit. The average performance listed is calculated across all 10 benchmarks.

Method	VOT-RGBD22			LasHeR		RGBT234		VisEvent		TNL2K		OTB99		Avg Perf.(%)
	EAO	Acc	Rob	AUC	P	MSR	MPR	AUC	P	AUC	P	AUC	P	
<i>Upper Bound, 294 Tokens (100%)</i>														
SUTrack <sub>224</sub>	75.5 100%	82.5 100%	91.3 100%	59.9 100%	74.8 100%	70.0 100%	92.1 100%	63.3 100%	80.7 100%	67.8 100%	73.8 100%	68.4 100%	91.1 100%	100%
<i>Retain 71.4% Tokens (<math>\approx 218</math> tokens, <math>\downarrow 25.6\%</math>)</i>														
SUTrack-CE <sub>224</sub>	75.8 100.4%	82.8 100.4%	91.4 100.1%	59.9 100.0%	74.7 99.9%	70.0 100.0%	92.5 100.4%	61.8 97.6%	79.2 98.1%	66.0 97.3%	71.5 96.9%	70.5 103.0%	92.6 101.7%	<b>99.5%</b> ( $\downarrow 0.5\%$ )
SUTrack-ToMe <sub>224</sub>	76.0 100.7%	82.5 100.0%	91.9 100.7%	60.0 100.2%	74.6 99.7%	69.0 98.6%	92.0 99.9%	62.3 98.4%	79.6 98.6%	66.2 97.6%	71.8 97.3%	70.4 102.9%	91.3 100.3%	99.4% ( $\downarrow 0.6\%$ )
SUTrack-EViT <sub>224</sub>	75.7 100.3%	82.6 100.1%	91.4 100.1%	59.0 98.5%	73.2 97.9%	70.0 100.0%	92.0 99.9%	62.0 97.9%	79.1 98.0%	65.7 96.9%	71.1 96.4%	70.7 103.4%	91.4 100.4%	99.0% ( $\downarrow 1.0\%$ )
SUTrack-DyViT <sub>224</sub>	74.4 98.5%	81.5 98.8%	91.2 99.9%	57.6 96.2%	71.8 96.0%	68.7 98.1%	90.8 98.6%	61.3 96.8%	79.4 98.4%	64.0 94.4%	65.5 88.9%	71.5 104.5%	93.6 102.8%	96.5% ( $\downarrow 3.5\%$ )
<b>UTPTrack<sub>224</sub></b>	76.5 101.3%	82.7 100.4%	92.1 100.9%	60.1 100.3%	74.9 100.1%	70.2 100.3%	93.1 101.1%	62.0 97.9%	79.2 98.1%	66.4 97.9%	72.0 97.6%	70.7 103.3%	93.1 102.2%	<b>99.8%</b> ( $\downarrow 0.2\%$ )
<i>Retain 52.0% Tokens (<math>\approx 153</math> tokens, <math>\downarrow 48.0\%</math>)</i>														
SUTrack-CE <sub>224</sub>	76.1 100.8%	82.8 100.4%	91.8 100.5%	58.4 97.5%	72.7 97.2%	70.4 100.6%	93.0 101.0%	61.5 97.2%	78.6 97.4%	66.1 97.5%	71.7 97.2%	70.5 103.0%	92.6 101.7%	<b>99.2%</b> ( $\downarrow 0.8\%$ )
SUTrack-ToMe <sub>224</sub>	76.1 100.8%	82.7 100.2%	91.8 100.5%	59.2 98.8%	73.7 98.5%	69.4 99.1%	91.8 99.7%	61.7 97.5%	78.9 97.8%	65.9 97.1%	71.4 96.9%	70.5 103.1%	91.6 100.6%	99.0% ( $\downarrow 0.5\%$ )
SUTrack-EViT <sub>224</sub>	75.2 99.6%	82.1 99.5%	91.4 100.1%	57.8 96.5%	71.9 96.1%	69.9 99.9%	92.3 100.2%	62.3 98.4%	79.8 98.9%	65.1 96.0%	70.3 95.4%	70.5 103.1%	91.7 100.7%	98.5% ( $\downarrow 1.5\%$ )
SUTrack-DyViT <sub>224</sub>	63.0 83.4%	74.3 90.1%	84.8 92.9%	50.6 84.5%	63.4 84.8%	63.3 90.4%	86.3 93.7%	49.6 78.4%	67.2 83.3%	54.3 80.0%	54.4 73.7%	65.3 95.5%	86.0 94.4%	78.8% ( $\downarrow 21.2\%$ )
<b>UTPTrack<sub>224</sub></b>	75.7 100.3%	82.8 100.4%	91.4 100.1%	59.6 99.5%	74.4 99.5%	70.7 101.0%	93.6 101.6%	61.7 97.5%	78.8 97.6%	65.9 97.2%	71.5 96.9%	70.1 102.5%	91.2 100.1%	<b>99.5%</b> ( $\downarrow 0.5\%$ )
<i>Retain 35.4% Tokens (<math>\approx 104</math> tokens, <math>\downarrow 64.6\%</math>)</i>														
SUTrack-CE <sub>224</sub>	75.5 100.0%	82.2 99.6%	91.8 100.5%	57.8 96.5%	71.9 96.1%	69.3 99.0%	92.2 100.1%	61.4 97.0%	78.7 97.5%	65.4 96.3%	70.5 95.5%	71.7 104.8%	94.2 103.5%	<b>98.3%</b> ( $\downarrow 1.7\%$ )
SUTrack-ToMe <sub>224</sub>	70.9 93.9%	78.6 95.3%	90.2 98.8%	53.8 89.8%	66.6 89.0%	63.6 90.9%	87.9 95.4%	58.0 91.6%	77.4 95.9%	61.7 91.0%	63.6 86.3%	67.3 98.4%	87.0 95.6%	92.5% ( $\downarrow 7.5\%$ )
SUTrack-EViT <sub>224</sub>	75.0 99.3%	81.9 99.3%	91.1 99.8%	56.4 94.2%	70.3 94.0%	67.8 96.9%	90.2 97.9%	60.8 96.1%	78.4 97.1%	64.4 95.0%	69.3 94.0%	70.4 102.9%	91.5 100.5%	96.7% ( $\downarrow 3.3\%$ )
SUTrack-DyViT <sub>224</sub>	13.7 18.1%	49.7 60.2%	23.6 25.8%	10.3 17.2%	15.5 20.7%	17.9 25.6%	32.4 35.2%	7.3 11.5%	14.8 18.3%	15.9 23.4%	13.5 18.3%	13.6 19.9%	20.8 22.9%	14.7% ( $\downarrow 85.3\%$ )
<b>UTPTrack<sub>224</sub></b>	76.1 100.8%	83.0 100.6%	91.5 100.2%	58.6 97.8%	72.9 97.5%	70.4 100.6%	92.9 100.9%	61.5 97.2%	79.1 98.0%	66.0 97.3%	71.4 96.9%	70.3 102.7%	91.9 100.9%	<b>99.3%</b> ( $\downarrow 0.7\%$ )

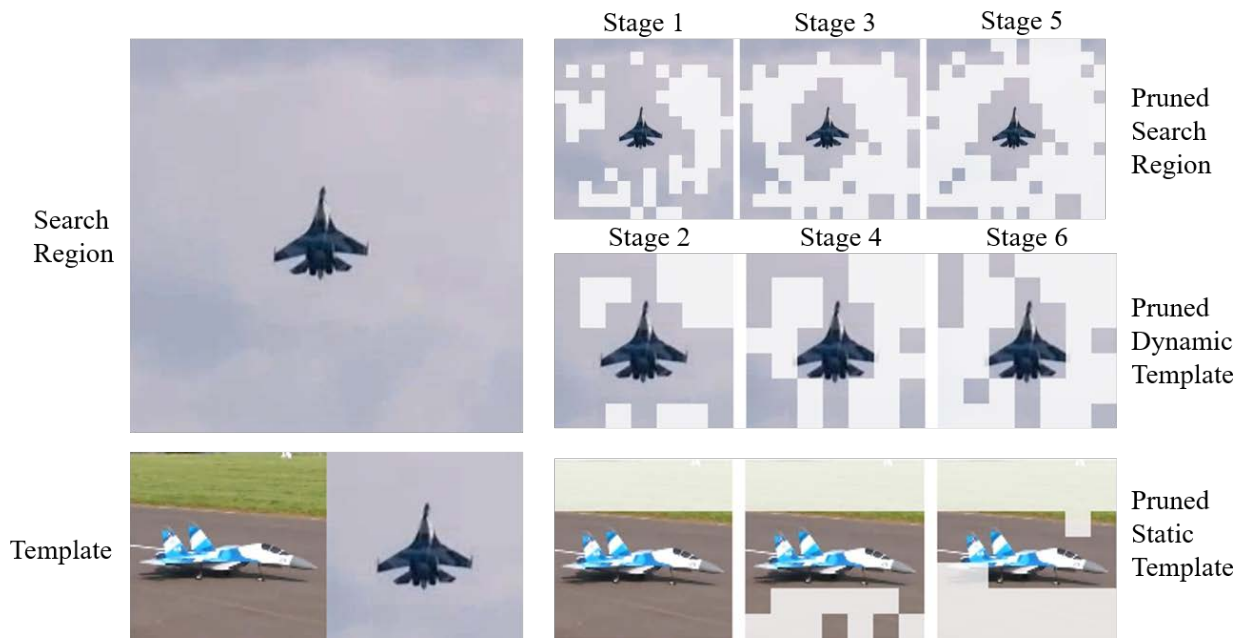


Figure 5. Visualization of the UTPTrack Pruning Process for RGB-based Tracking.

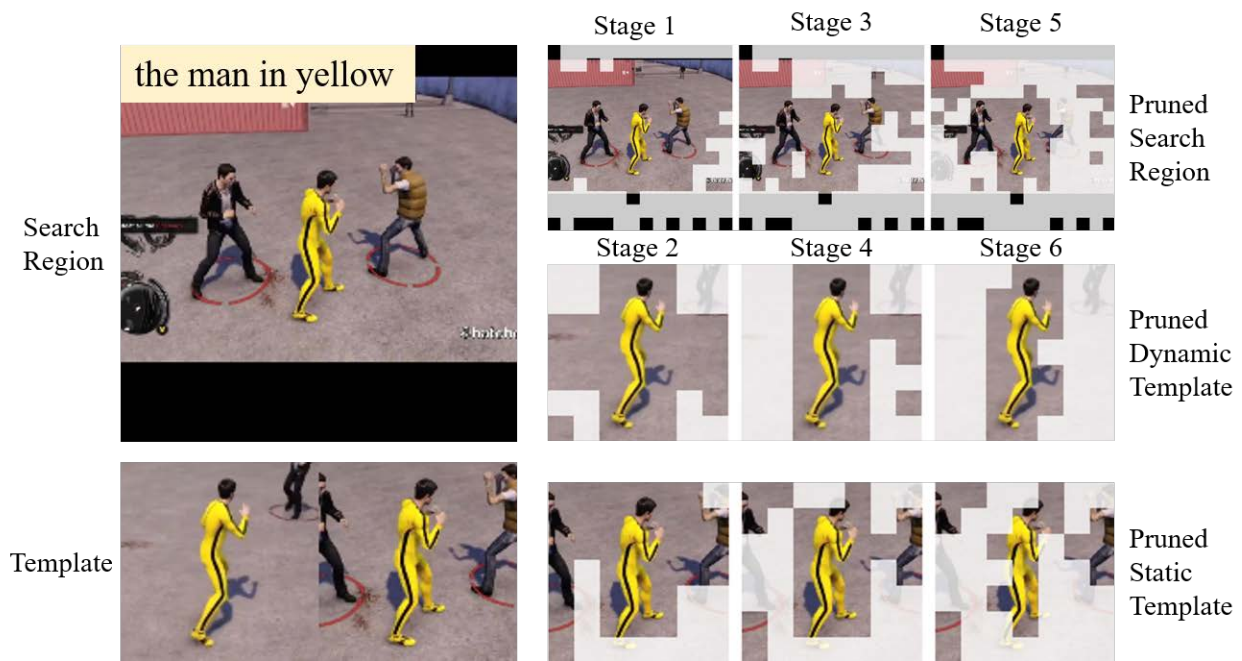


Figure 6. Visualization of the UTPTrack Pruning Process for RGB-Language Tracking.

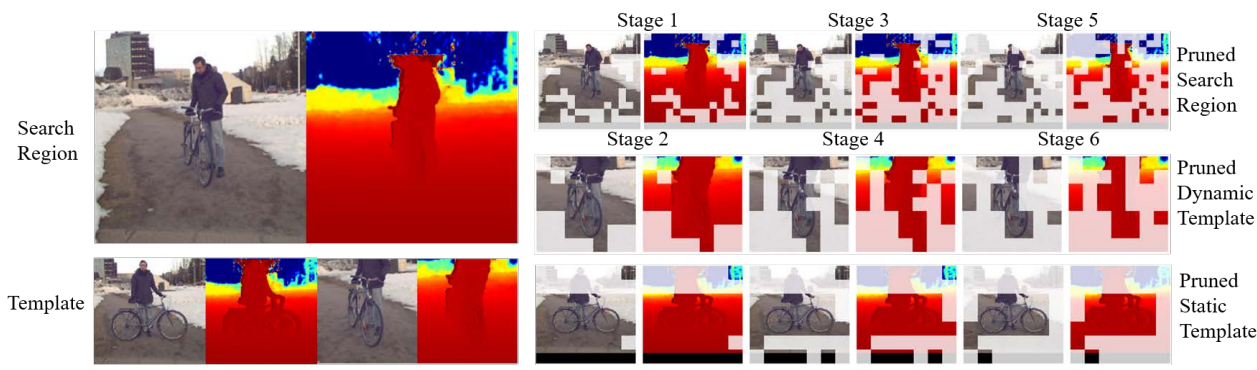


Figure 7. Visualization of the UTPTrack Pruning Process for RGB-Depth Tracking.

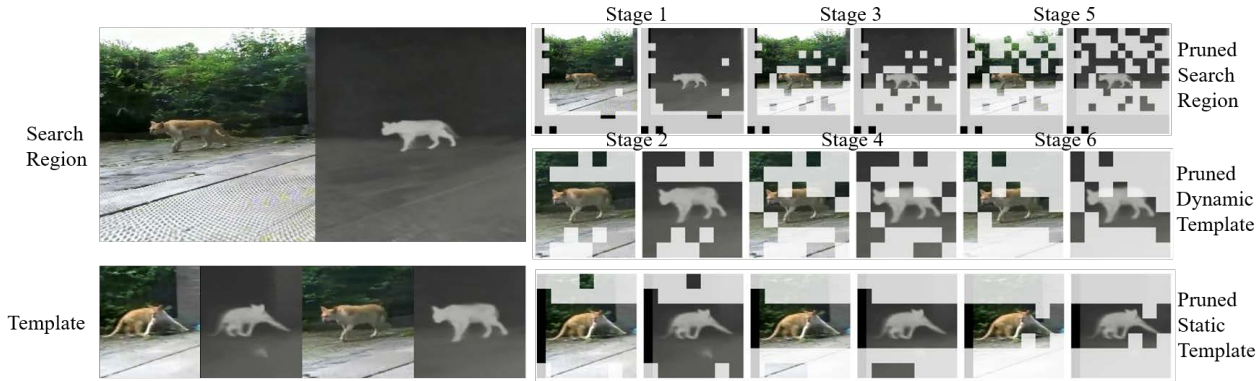


Figure 8. Visualization of the UTPTrack Pruning Process for RGB-Thermal Tracking.

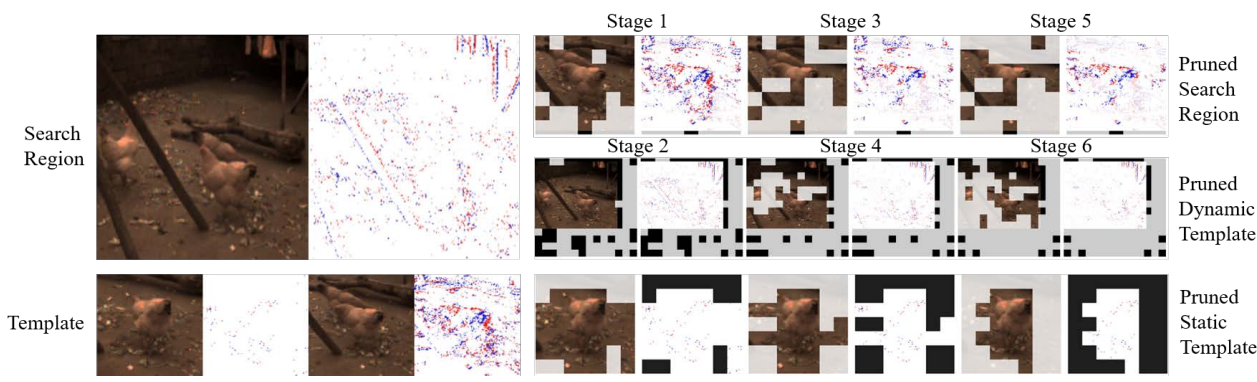


Figure 9. Visualization of the UTPTrack Pruning Process for RGB-Event Tracking.