

ProTracker: Probabilistic Integration for Robust and Accurate Point Tracking

Supplementary Material

001 1. Video Results

002 Please refer to our Supplementary Webpage for the corre-
003 sponding videos of images illustrated in the paper and more
004 results on different data.

005 2. Implementation Details

006 2.1. Flow Preparation

007 We utilize the RAFT [5] model, as adopted by MFT [4], as
008 our flow estimation model. It takes two images, I_j and I_i ,
009 captured at different times as input and outputs the flow map
010 \mathcal{F}_{ji} , occlusion map \mathcal{O}_{ji} , and uncertainty map \mathcal{U}_{ji} . Since
011 the uncertainty of an estimation can also be interpreted as its
012 variance, the initial flow prediction from frame j to frame i
013 at location p is computed as follows:

$$014 (\mathbf{f}_{ji}, \mu_{ji}) = \begin{cases} (S(\mathcal{F}_{ji}, p), S(\mathcal{U}_{ji}, p)) & \text{if } (\mathcal{O}_{ji}, p) > \rho \\ \text{None} & \text{otherwise} \end{cases} \quad (1)$$

015 where $\rho = 0.1$, and $S(\text{Target}, p)$ indicates sampling the
016 target at location p . The initial flow predictions are then
017 set as input to flow integration. Following MFT, we adopt
018 $\{\infty, 1, 2, 4, 8, 16, 32\}$ as the time intervals, meaning that
019 each frame's prediction is computed by combining results
020 from these previous frames.

021 2.2. Filter Thresholds

022 During the dual filtering stage, we apply different thresh-
023 olds to predictions from flow and long-term keypoints. For
024 long-term keypoints, a prediction is marked as invalid if the
025 cosine similarity to the query point on the geometry-aware
026 feature falls below 0.5. For predictions from flow, however,
027 we use a threshold of 0.3 instead. These distinct thresh-
028 olds allow flow to track featureless areas while ensuring
029 that long-term keypoints do not drift into visually similar
030 regions.

031 2.3. Outlier Removal in Integration

032 Even after dual filtering, occasional erroneous predictions
033 may persist. To ensure a more stable integration pro-
034 cess and minimize the impact of these incorrect predic-
035 tions, we discard rough predictions that deviate signif-
036 icantly from others. Denoting the input predictions as
037 $\{(\mathbf{f}_1, \mu_1), (\mathbf{f}_2, \mu_2), \dots, (\mathbf{f}_N, \mu_N)\}$, we remove outliers us-
038 ing the following criterion:

$$039 (\mathbf{f}_i, \mu_i) = \begin{cases} (\mathbf{f}_i, \mu_i) & \text{if } \text{Norm}(\mathbf{f}_i - \mathbf{f}_T) < \rho_{dist} \\ \text{None} & \text{otherwise} \end{cases} \quad (2)$$

where $\rho_{dist} = 10$ and $\text{Norm}(\mathbf{x})$ measures the magnitude
of a vector.

040 3. Dense Inference

041 As discussed in Sec.3.1 in the main paper, we utilize a
042 geometry-aware feature extractor and a video mask gen-
043 erator for the dual-filter stage. While optical flow and
044 geometry-aware features can be computed densely, generat-
045 ing masks for each pixel is both time-intensive and memory-
046 intensive. To address this, we adopt an iterative approach to
047 efficiently generate a set of masks that collectively cover all
048 pixels, as described below:

049 Algorithm 1 Dense Mask Generation Algorithm

```
050 1: Initialize: Set all pixels as unassigned.  
051 2: while there exist unassigned pixels do  
052 3:   Select the first unassigned pixel  $p$ .  
053 4:   Generate a new mask  $M$  starting from pixel  $p$ .  
054 5:   for each pixel  $q$  in  $M$  do  
055 6:     if  $q$  is unassigned then  
056 7:       Assign  $q$  to mask  $M$ .  
057 8:     end if  
058 9:   end for  
059 10: end while  
060 11: Output: All pixels assigned to corresponding masks.
```

061 Subsequently, the dual filter can be applied to each pixel
062 based on its corresponding mask.

063 4. Training and Inference Speed

064 The total time consumed for our method includes the time
065 for keypoint extraction, mask generation, geometry-aware
066 feature extraction and probabilistic integration. During key-
067 point extraction, we follow DINO-Tracker [6] to train a
068 delta-DINO model and a heatmap refiner, which takes about
069 1 hour for an 80-frame video on a single RTX 4090 GPU.
070 We refer to DINO-Tracker [6] for more details. How-
071 ever, our method skips the time-consuming occlusion pre-
072 diction and directly uses points with high cosine similarity
073 as keypoints, which saves much time. The mask genera-
074 tion and geometry-aware feature extraction together takes
075 about 2 minutes and the probabilistic integration takes about
076 1 minute for the same video.

077 In total, during the inference stage, the time spent track-
078 ing 3,000 points on a single object in an 80-frame video
079 is about 3 minutes, which is about 20x faster than DINO-
080 Tracker [6]. For dense inference, an additional 4 minutes

071 may be required due to the increased number of masks gen-
072 erated, but our method remains more than 30x faster than
073 DINO-Tracker [6].

074 5. More Qualitative Results

075 To further illustrate our methods’ robustness. We conduct
076 experiments on more challenging cases and show the qual-
077 iative results.

078 Some of the previous methods rely on computing a
079 heatmap between the query point and the target frame.
080 However, the per-frame heatmap lacks temporal-awareness
081 and may confuse different objects. We address this issue
082 by leveraging the mask and combining the heatmap with
083 optical flow. As illustrated in Fig. 1 and Fig. 2, by com-
084 paring the results of our method with DINO-Tracker [6] and
085 TAPIR [1], we show that although our method also relies
086 on per-frame heatmap to extract keypoints,our method has
087 strong temporal-awareness and is able to tell between simi-
088 lar objects.

089 To further demonstrate the robustness of our method,
090 we conduct experiments on extended videos from TAP-
091 Vid-DAVIS, simulating high frame-rate videos by repeat-
092 ing each frame three times, as illustrated in Fig. 3 and
093 Fig. 4. In contrast to typical sliding-window or flow-based
094 trackers (such as TAPTR [3], SpatialTracker [7] and Co-
095 Tracker [2]), which tend to accumulate errors and drift over
096 time, our integration of long-term key points with short-
097 term optical flow enables continuous, drift-free tracking of
098 the same point through occlusions.

099 References

- 100 [1] Carl Doersch, Yi Yang, Mel Vecerik, Dilara Gokay, Ankush
101 Gupta, Yusuf Aytar, Joao Carreira, and Andrew Zisserman.
102 Tapir: Tracking any point with per-frame initialization and
103 temporal refinement. In *Proceedings of the IEEE/CVF Inter-*
104 *national Conference on Computer Vision*, pages 10061–
105 10072, 2023. [2](#)
- 106 [2] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia
107 Neverova, Andrea Vedaldi, and Christian Rupprecht. Co-
108 tracker: It is better to track together. *arXiv preprint*
109 *arXiv:2307.07635*, 2023. [2](#)
- 110 [3] Hongyang Li, Hao Zhang, Shilong Liu, Zhaoyang Zeng,
111 Tianhe Ren, Feng Li, and Lei Zhang. Taptr: Tracking any
112 point with transformers as detection. In *Proceedings of the*
113 *IEEE/CVF European Conference on Computer Vision*, 2024.
114 [2](#)
- 115 [4] Michal Neoral, Jonáš Šerých, and Jiří Matas. Mft: Long-
116 term tracking of every pixel. In *Proceedings of the IEEE/CVF*
117 *Winter Conference on Applications of Computer Vision*, pages
118 6837–6847, 2024. [1](#)
- 119 [5] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field
120 transforms for optical flow. In *Computer Vision–ECCV 2020:*
121 *16th European Conference, Glasgow, UK, August 23–28,*

- 122 2020, *Proceedings, Part II 16*, pages 402–419. Springer, 2020.
123 [1](#)
- 124 [6] Narek Tumanyan, Assaf Singer, Shai Bagon, and Tali Dekel.
125 Dino-tracker: Taming dino for self-supervised point tracking
126 in a single video, 2024. [1](#), [2](#)
- 127 [7] Yuxi Xiao, Qianqian Wang, Shangzhan Zhang, Nan Xue, Sida
128 Peng, Yujun Shen, and Xiaowei Zhou. Spatialtracker: Track-
129 ing any 2d pixels in 3d space. In *Proceedings of the IEEE/CVF*
130 *Conference on Computer Vision and Pattern Recognition*,
131 pages 20406–20417, 2024. [2](#)



Figure 1. Results of tracking a single object. While DINO-Tracker may mispredict parts onto similar objects and TAPIR can be disrupted by similar patterns, our method avoids these errors.

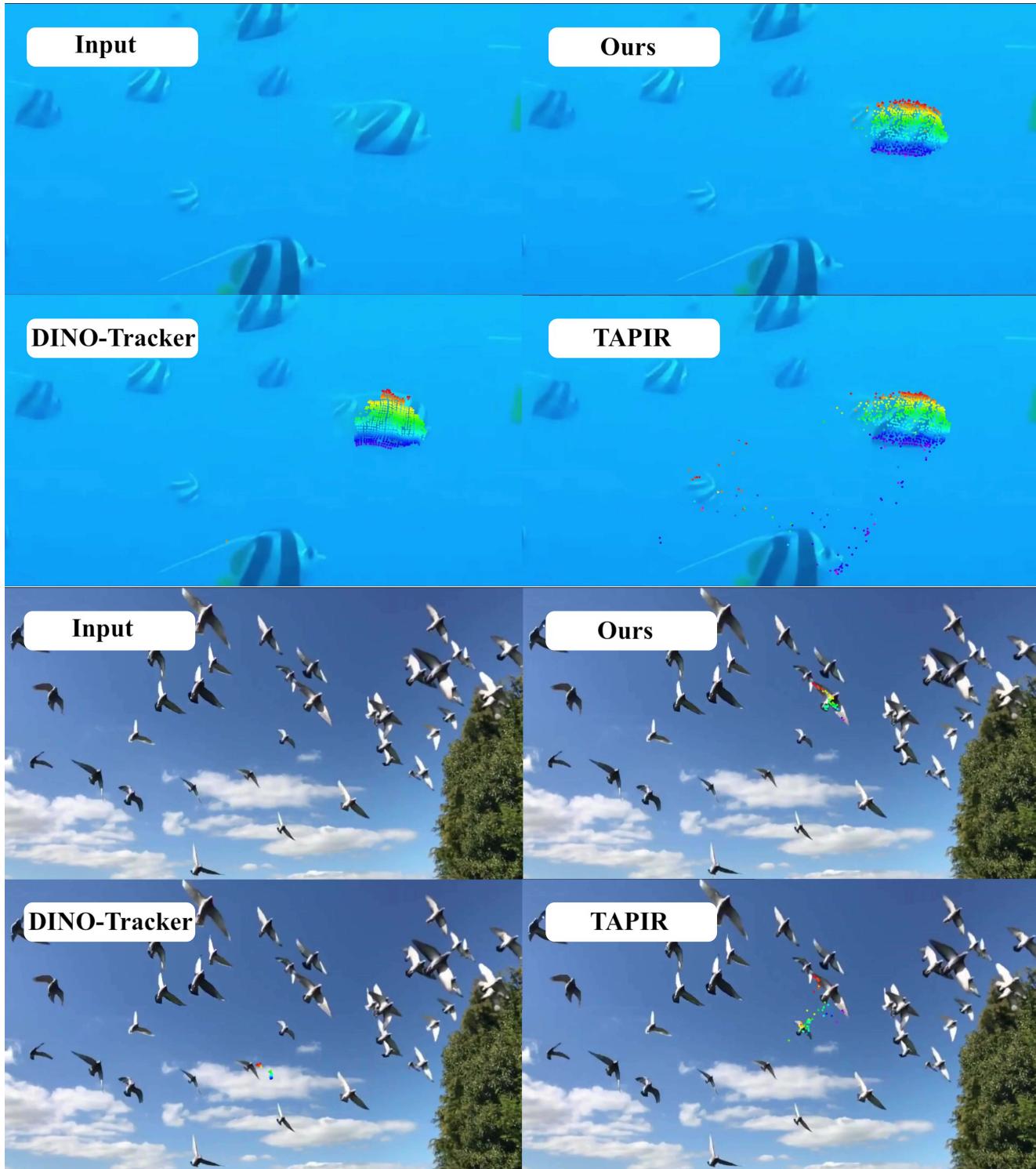


Figure 2. Results of tracking a single object. While DINO-Tracker may lose some parts and TAPIR can be disrupted by multiple similar patterns, our method avoids these errors.

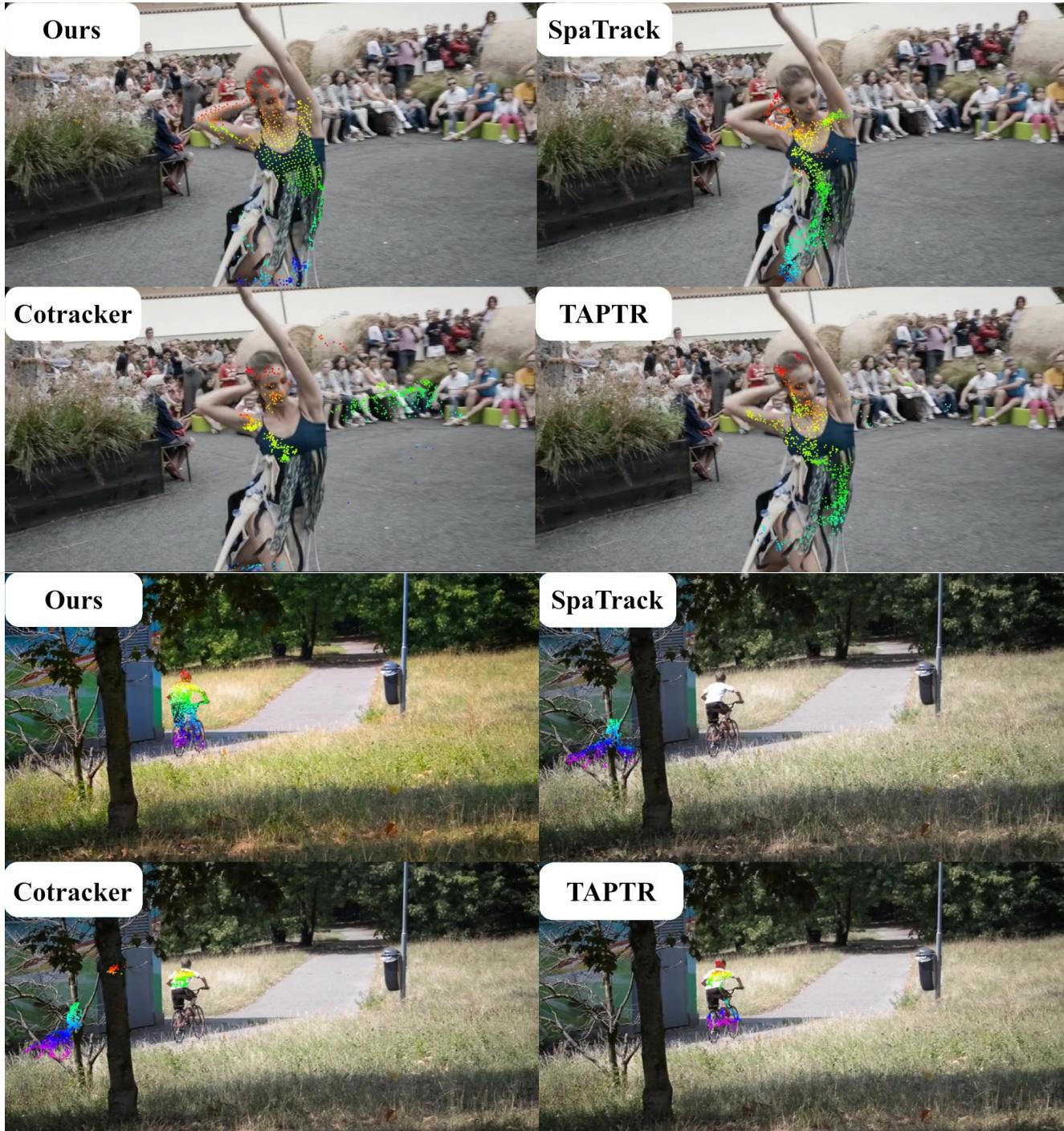


Figure 3. Results of tracking at a higher frame rate. Sliding window based methods can easily lose track after occlusion and drift due to accumulating errors, while ours exhibit robustness.



Figure 4. Results of tracking at a higher frame rate. Sliding window based methods can mispredict points to other regions during occlusion (e.g. the gun and rope in *shooting* and the wrong person in *india*), while ours exhibit robustness.