
TrackIME: Enhanced Video Point Tracking via Instance Motion Estimation

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 Tracking points in video frames is essential for understanding video content. How-
2 ever, the task is fundamentally hindered by the computation demands for brute-force
3 correspondence matching across the frames. As the current models down-sample
4 the frame resolutions to mitigate this challenge, they fall short in accurately repre-
5 senting point trajectories due to information truncation. Instead, we address the
6 challenge by pruning the search space for point tracking and let the model process
7 only the important regions of the frames without down-sampling. Our first key
8 idea is to identify the object instance and its trajectory over the frames, then prune
9 the regions of the frame that do not contain the instance. Concretely, to estimate
10 the instance’s trajectory, we track a group of points on the instance and aggregate
11 their motion trajectories. Furthermore, to deal with the occlusions in complex
12 scenes, we propose to compensate for the occluded points while tracking. To this
13 end, we introduce a unified framework that jointly performs point tracking and
14 segmentation, providing synergistic effects between the two tasks. For example,
15 the segmentation results enable a tracking model to avoid the occluded points
16 referring to the instance mask, and conversely, the improved tracking results can
17 help to produce more accurate segmentation masks. Our framework can be easily
18 incorporated with various tracking models, and we demonstrate its efficacy for
19 enhanced point tracking throughout extensive experiments. For example, on the
20 recent TAP-Vid benchmark, our framework consistently improves all baselines,
21 e.g., up to 13.5% improvement on the average Jaccard metric.

22

1 Introduction

23 Obtaining accurate point trajectories over the video frames is crucial for understanding complex
24 dynamics in video data, a necessity for advanced spatial-temporal tasks like action recognition [2],
25 novel-view rendering [3], video frame prediction/interpolation [4], and video depth estimation [5].
26 Recently, video point tracking task [6, 7, 8, 9, 10] has witnessed rapid progress, which aims to predict
27 the trajectory and visibility¹ of a given query point, proving long-term trajectories robust to partial
28 occlusions of objects in real video scenes.

29 Despite their success, we find current point tracking models are fundamentally challenged by an
30 excessive computation demand since the task requires brute-force comparisons over every spatial
31 location in every frame in a given video. As a result, to meet the computation constraints, the models
32 down-sample their tracking resolutions, sacrificing detailed visual features, which eventually leads to
33 sub-optimal tracking accuracy and triggers tracking failures on intricate object parts. In this regard,
34 we pursue the direction of pruning the excessive search space for point tracking, so that models can

¹The confidence whether the trajectory is visible in each frame; *i.e.*, the point is not out-of-frame and not occluded by different objects.

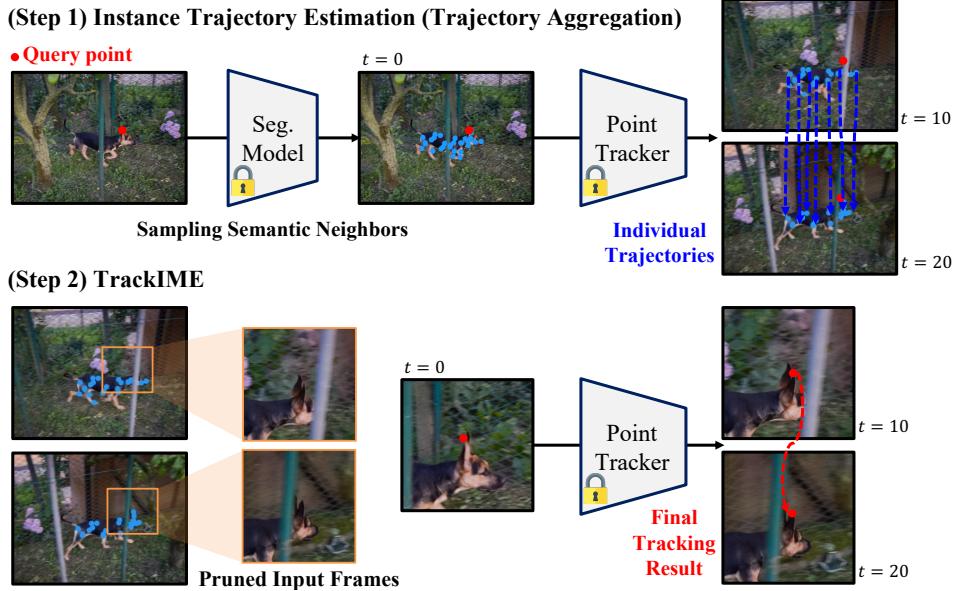


Figure 1: **The workflow of TrackIME.** Our framework enhances point tracking by pruning the search space, along the instance trajectory in video frames. To estimate the instance trajectory, our framework utilizes the point tracking results for a group of points (blue lines) on top of the object instance predicted by segmentation model (*e.g.*, SAM [1]) and aggregate their individual trajectories.

- 35 avoid the down-sampling and focus only on important regions maintaining detailed visual features,
 36 *e.g.*, the object instance masks the query point lies in.
- 37 In this paper, we introduce TrackIME: Enhanced Video Point Tracking via Instance Motion Estimation
 38 that focuses on the region occupied by the object instance that the queried point lies in and guides
 39 point tracking models to prune the video frames along the instance’s motion trajectory. Here, to obtain
 40 the instance trajectory, we first produce the instance mask for a given query point by utilizing the
 41 recent segmentation foundation models, *e.g.*, segment anything (SAM) [1], where these foundation
 42 models show strong generalization performance to different objects/scenes and we find resulting
 43 instance masks in quality are readily available. Then, given the instance mask, we sample a set of
 44 points and aggregate their tracking results as the estimate of the instance trajectory.²
- 45 Furthermore, to deal with the occlusions in complex video scenes, we propose a unified framework
 46 that jointly performs the point tracking and video segmentation, where it re-samples the occluded
 47 points by referring to the instance mask. We note that our framework provides synergistic effects
 48 for both tasks, *i.e.*, the point tracking results assisted by the segmentation can conversely bolster
 49 the quality of segmentation. Consequently, although our primary focus is on the advances in point
 50 tracking, our method can also demonstrate improved segmentation results than the baselines (see
 51 Section 4.2 for details).
- 52 Through the experiments on the TAP-Vid point tracking benchmark [11], we demonstrate the
 53 effectiveness of TrackIME by incorporating it with different point tracking models such as TAPIR [6].
 54 For example, in the DAVIS scenes [12] evaluating the point tracking for dynamic objects, our method
 55 achieved up to 13.5% relative improvement ($57.5 \rightarrow 65.3$ with TAPIR) in terms of the average Jaccard
 56 (AJ) metric. Moreover, as our framework allows pruning non-instance regions for point tracking
 57 models, the efficacy of our method stands out even more when evaluated on more harsh standards,
 58 *e.g.*, the 1-pixel error threshold, where the conventional metrics allow up to 16-pixel errors when
 59 judging the prediction to be correct.

²Intuitively, the instance as a group of points moves together even if a fine-grained motion of individual points may differ. Hence, we track multiple points on the same instance, and then aggregate their trajectories as the instance motion, which we eventually utilize for pruning video frames.

60 **2 Method**

61 In this section, we describe the detailed procedure of our TrackIME framework and its application to
 62 video point tracking. Specifically, in Section 2.1, we describe the formulation for instance trajectory
 63 estimation, which is based on the video point tracking of the query points found by the foundation
 64 segmentation model [1]. Next, in Section 2.2, we present the detailed formulation of TrackIME given
 65 the instance trajectory, which prunes unimportant regions in the video frame and achieves boosted
 66 point tracking performances.

67 As for the data notations, we denote vectors with N elements as bold letters $\mathbf{x} := [\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_N]$,
 68 tensors with N arrays as bold capital letters $\mathbf{X} := [\mathbf{X}_1; \mathbf{X}_2; \dots; \mathbf{X}_N]$, where the subscripts represent
 69 the indexed scalars or arrays. Otherwise, every non-bold symbol is scalar. We also introduce the
 70 superscripts, *e.g.*, $\mathbf{x}^{(q)}$, when denoting there is special semantics for a data, such as the query point.

71 Finally, when making binary classifications based on probability (or normalized confidence) values,
 72 we use threshold 0.5; nevertheless, the values are hyperparameters and can be altered in practice.

73 **2.1 Instance Trajectory Estimation**

74 In this section, we provide the definition of the instance trajectory and procedures to obtain it, such as
 75 sampling a group of points on the instance, trajectory aggregation, and the point re-sampling modules.

76 **Video point tracking.** Let $\mathbf{I} \in \mathbb{R}^{L \times H \times W \times 3}$ be the tensor of video frames, where L denotes the time
 77 duration and $(H \times W)$ denotes the image size, and let $\mathbf{p}^{(q)} \in \mathbb{R}^2$ be the spatial coordinates of the
 78 query point. Typically, we consider the query in the initial frame hence we do not denote the time
 79 index of the query point for clarity. Given the video \mathbf{I} and the query point $\mathbf{p}^{(q)}$, we consider a point
 80 tracking model Tracker that predicts the query trajectory $\mathbf{T}^{(q)} \in \mathbb{R}^{L \times 2}$ and the probability of being
 81 visible $\mathbf{o}^{(q)} \in (0, 1)^L$ over the entire set of frames,

$$(\mathbf{T}^{(q)}, \mathbf{o}^{(q)}) := \text{Tracker}(\mathbf{p}^{(q)}, \mathbf{I}). \quad (1)$$

82 Here, one might utilize Equation (1) as the simplest representation of the instance motion trajectory.
 83 However, modeling the instance motion solely using the query point has critical shortcomings. For
 84 example, when the instance is partially occluded by other objects, the trajectory of the query point
 85 may no longer exist (see Section 4.1 for the ablation study). To address this challenge, we propose to
 86 sample additional tracking points automatically. Specifically, our idea is to identify the instance mask
 87 of the query point so that extra query points can be added from the mask.

88 **Sampling points on the instance.** Let $\mathbf{M}_0 \in (0, 1)^{H \times W}$ denote the segmentation mask that
 89 represents the object instance associated with the query point $\mathbf{p}^{(q)}$. Given this mask, we sample a
 90 group of points on the instance,

$$\mathcal{N}(\mathbf{p}^{(q)}) := \{\mathbf{p}^{(n_0)}, \dots, \mathbf{p}^{(n_S)}\}, \quad (2)$$

91 which we refer to it as the *semantic neighbors* of $\mathbf{p}^{(q)}$. We note that S is the number of sampled
 92 points, where the query point is also counted as its semantic neighbor, *i.e.*, $\mathbf{p}^{(n_0)} \equiv \mathbf{p}^{(q)}$. For each
 93 semantic neighbor point, we employ Tracker in Equation (1) to produce its trajectory and visibility,
 94 $(\mathbf{T}^{(n_i)}, \mathbf{o}^{(n_i)}) := \text{Tracker}(\mathbf{p}^{(n_i)}, \mathbf{I})$,³ and pass it to the trajectory aggregation module.

95 **Trajectory aggregation.** We produce an instance motion trajectory by aggregating the tracking
 96 results of the semantic neighbors. Specifically, we consider the velocity, $\Delta \mathbf{T}_t^{(n_i)} := \mathbf{T}_t^{(n_i)} - \mathbf{T}_{t-1}^{(n_i)}$,
 97 and calculate the weighted average:

$$\Delta \bar{\mathbf{T}}_t^{(q)} := \sum_{(\mathbf{o}_t^{(n_i)} \geq 0.5)} \frac{\mathbf{o}_t^{(n_i)} \cdot \Delta \mathbf{T}_t^{(n_i)}}{\sum_{(\mathbf{o}_t^{(n_j)} \geq 0.5)} \mathbf{o}_t^{(n_j)}}. \quad (3)$$

98 In Equation (3), we note that velocities are aggregated only if the points are classified visible
 99 ($\mathbf{o}_t^{(n_i)} \geq 0.5$), and the visibility acts as the aggregation weight. Finally, we accumulate the aggregated

³In practice, we batch-process a set of multiple points simultaneously.

100 velocity starting from $\bar{\mathbf{T}}_0^{(q)} := \mathbf{p}^{(q)}$, to obtain the instance motion trajectory,

$$\bar{\mathbf{T}}_t^{(q)} := \bar{\mathbf{T}}_{t-1}^{(q)} + \Delta \bar{\mathbf{T}}_t^{(q)}. \quad (4)$$

101 **Instance mask.** In order to identify the instance mask, we employ the recent foundation segmentation
 102 model, *e.g.*, Segment Anything Model (SAM) [1], and prompt the model with the query point $\mathbf{p}^{(q)}$, to
 103 produce the pixel-wise confidence representing the object instance indicated by the query point. We
 104 denote this function as Seg ,

$$\mathbf{M}_0 := \text{Seg}(\mathbf{p}^{(q)}, \mathbf{I}_0) \in (0, 1)^{H \times W}. \quad (5)$$

105 Given the mask \mathbf{M}_0 , we employ a weighted sampling for the semantic neighbors. Specifically, we
 106 encode the sampling weights with the distance transform (DT) [13, 14] to the mask's region with
 107 positive classifications,

$$\mathbf{W}_0 := \text{DT}(\mathbf{1}[\mathbf{M}_0 \geq 0.5]). \quad (6)$$

108 In this way, the points near the mask's contour are preferred, which we find efficiently represent the
 109 object instance because the contour is approximately linearly proportional to the mask's radius.

110 **Point re-sampling for robustness to occlusion.** Occlusions are common in real video frames, due to
 111 dynamic objects and the camera's motion. In the extreme case, Equation (3) can become degenerate
 112 when all semantic neighbors are invisible in future frames $t > 0$. Therefore, maintaining a sufficient
 113 number of visible tracking points is crucial, and we tackle this issue by re-sampling occluded points
 114 from the instance mask jointly predicted while point tracking.

115 In a nutshell, whenever we find a certain semantic neighbor point $\mathbf{p}^{(n_i)}$ becomes invisible at time
 116 t' and does not show up again ($\mathbf{o}_{t'}^{(n_i)} < 0.5$ for $t \geq t'$), we query the segmentation model with the
 117 tracking results of other semantic neighbors to obtain a new mask to re-sample the occluded point:

$$\mathbf{M}_{t'}^{(n_j)} := \text{Seg}(\mathbf{T}_{t'}^{(n_j)}, \mathbf{I}_{t'}) \in (0, 1)^{H \times W}. \quad (7)$$

118 However, they could also have been affected by occlusions (*e.g.*, when $\mathbf{o}_{t'}^{(n_j)}$ is close to the threshold
 119 0.5), or by the severe errors in the trajectory $\mathbf{T}_{t'}^{(n_j)}$ due to sub-optimal tracking performance of
 120 Equation (1). Hence, predicting segmentation with these points in a naive way can lead to erroneous
 121 masks being produced.

122 To address this problem, our key idea is to aggregate the group of segmentation masks. Specifically, we
 123 collect individual masks by Equation (7), then apply a weighted average of the positive classifications,

$$\bar{\mathbf{M}}_{t'} := \sum_{\left(\mathbf{o}_{t'}^{(n_i)} \geq 0.5\right)} \frac{\mathbf{o}_{t'}^{(n_i)} \cdot \mathbf{1}[\mathbf{M}_{t'}^{(n_i)} > 0.5]}{\sum_{\left(\mathbf{o}_{t'}^{(n_j)} \geq 0.5\right)} \mathbf{o}_{t'}^{(n_j)}} \in (0, 1)^{H \times W}. \quad (8)$$

124 We find the mask produced by Equation (8) reflects the confidence of each segmentation mask, as
 125 well as the visibility of the associated point, and refer to it as the *mixture of segmentation distributions*.
 126 Based on the constructed mixture of segmentation distributions, we obtain the sampling weight in
 127 similar manner to Equation (6) as, $\mathbf{W}_{t'} := \text{DT}(\mathbf{1}[\bar{\mathbf{M}}_{t'} \geq r])$, where the threshold $r \in [0, 1]$ is set much
 128 smaller than the standard 0.5. This is because we should allow the confident partial segmentation
 129 distributions, but ignore the unconfident noise segmentation distributions.

130 Finally, we re-sample the additional points with $\mathbf{W}_{t'}$ as the sampling weight. They replace the
 131 occluded points for the instance trajectory estimation in subsequent frames $t > t'$. We execute this
 132 procedure during the tracking, which ensures that a sufficient number of visible points participate
 133 in Equations (3) and (4). For example, we set it to be the same as the number of initial semantic
 134 neighbors S .

135 2.2 TrackIME: Enhanced Video Point Tracking via Instance Motion Estimation

136 In this section, we describe our enhanced point tracking, which prunes the search spaces in frames
 137 and produces more accurate tracking results.

138 **Search space pruning.** Given the instance trajectory in Equation (4), we now aim to utilize it for
 139 pruning the search space. Specifically, we prune unimportant non-instance regions, by sampling each
 140 frame around the $(H_0 \times W_0)$ regions centered at the aggregated trajectory,

$$\mathbf{I}^{(q)} := \text{Prune}(\mathbf{I}, \bar{\mathbf{T}}^{(q)}, H_0, W_0) \in \mathbb{R}^{L \times H_0 \times W_0 \times 3}. \quad (9)$$

141 We note that the sizes $(H_0 \times W_0)$ are set to be close to the down-sampling resolution considered by
 142 a tracking model (*e.g.*, (256×256) for TAPIR [6]) so that the information loss is minimized.

143 Given the frames with pruned search spaces, we execute Tracker again to produce the enhanced
 144 tracking outputs. Also, for convenience, we abstract the entire process of the instance trajectory
 145 estimation (Section 2.1), the pruning (Equation (9)) and the tracking into a function TrackerHD,

$$\begin{aligned} (\mathbf{T}^{(\text{HD})}, \mathbf{o}^{(\text{HD})}) &:= \text{TrackerHD}(\mathbf{p}^{(q)}, \mathbf{I}, H_0, W_0) \\ &:= \text{Tracker}(\mathbf{p}^{(q)}, \mathbf{I}^{(q)}). \end{aligned} \quad (10)$$

146 We note that the feature resolutions inside the tracking model are not modified, therefore the computa-
 147 tional complexity does not increase.

148 **Progressive inference.** To achieve a further boost in the tracking performance, we can additionally
 149 use a progressive inference structure. Formally, we consider a collection of K different TrackerHD
 150 models equipped with different pruning sizes (H_k, W_k) :

$$\left[\mathbf{T}_1^{(\text{HD})}; \dots; \mathbf{T}_K^{(\text{HD})} \right] \text{ and } \left[\mathbf{o}_1^{(\text{HD})}; \dots; \mathbf{o}_K^{(\text{HD})} \right], \quad (11)$$

151 where $\mathbf{T}_k^{(\text{HD})} \in \mathbb{R}^{L \times 2}$ and $\mathbf{o}_k^{(\text{HD})} \in (0, 1)^L$ denotes the outputs of the k -th TrackerHD model.

152 This progressive structure can boost the tracking performance in two ways. The first is utilizing a past
 153 k -th TrackerHD as the tracking model that estimates the instance trajectory for the next $(k + 1)$ -th
 154 TrackerHD. In this way, the pruning is guided by a more accurate trajectory estimate. The second
 155 is that these K tracking results can be aggregated to produce the final trajectory. Specifically, we
 156 aggregate based on the visibility, in a similar manner to Equations (3) and (8):

$$\mathbf{T}^{(\text{Final})} := \sum_{k=1}^K \frac{\mathbf{o}_k^{(\text{HD})} \odot \mathbf{T}_k^{(\text{HD})}}{\sum_{l=1}^K \mathbf{o}_l^{(\text{HD})}}, \quad (12)$$

157 where \odot indicates the element-wise product. This aggregation allows processing multiple scales in
 158 visual features, which can enhance the generalization performance of vision models [15, 16]. We
 159 note that the visibility predictions are averaged over the K predictions.

160 3 Related Work

161 **Optical Flow.** Optical flow deals with the dense computation of instantaneous motion patterns
 162 between two given video frames. Starting with the pioneering work of applying neural networks for
 163 motion estimation [17, 18], the seminal works such as DCFNet [19], PWC-Net [20] and RAFT [21]
 164 introduced the concept of dense correspondence matching between pairs of image patches. Despite
 165 their success, the optical flow’s inherent limitations incapable of modeling trajectories and occlusions
 166 triggered the recent progress in the point tracking methods.

167 **Point Tracking.** In essence, point tracking attempts to find the long-term point correspondences
 168 over the entire video frames, and model the occlusions and trajectories. The current models in this
 169 domain, such as PIPs [22], TAPNet [11], TAPIR [6], CoTracker [7], and OmniMotion [8] has led
 170 rapid progress, with advanced neural architectures [6, 7] or test-time optimizations [8]. However,
 171 they are fundamentally hindered by the excessive search space for correspondence matching over the
 172 entire frames. Our focus is to address this issue by pruning the search space, where our method can
 173 be readily incorporated with these baselines.

174 **Instance Segmentation.** Recently, the important advancement within image segmentation has
 175 been the introduction of segment anything (SAM) [1]. SAM is specifically designed to perform
 176 image segmentation by general point prompts and exhibits an impressive capacity for class-agnostic
 177 segmentation. Specifically, in the context of point tracking, SAM serves as a valuable resource by
 178 generating segmentation masks for the object instance indicated by the query point. We also note the

Table 1: **The evaluation of point tracking performance for dynamic objects.** We benchmark the quality of point tracking in DAVIS [12] videos with the point annotations provided by TAP-Net [11]. We note that TrackIME is incorporated with TAPIR point tracker [6].

Method	First Query					Strided Query				
	J ₁	AJ	δ ₁ ^x	δ _{avg} ^x	OA	J ₁	AJ	δ ₁ ^x	δ _{avg} ^x	OA
TAPNet [11]	20.7	51.6	30.1	63.8	79.8	25.3	56.5	36.3	68.2	82.6
PIPS2 [9]	19.6	46.6	35.8	69.4	80.3	6.9	52.8	14.2	65.8	83.5
TAPIR [6]	23.0	57.5	34.3	70.5	84.4	28.1	62.8	41.0	75.1	87.7
CoTracker [7]	28.3	60.8	43.5	76.1	86.0	34.9	64.3	50.9	78.9	89.1
OmniMotion [8]	21.5	52.6	39.1	68.1	85.4	30.1	55.6	45.1	70.3	88.9
TrackIME	35.4	65.3	48.2	78.6	86.5	41.9	69.3	55.0	81.4	89.0

179 line of zero-shot video segmentation [23, 24, 25, 26, 27, 28, 29, 30]. Specifically, the recent SAM-PT
180 [30] focuses on bolstering video segmentation based on point tracking, which is fundamentally
181 different from our work; our primary goal is obtaining better point tracking, while that for SAM-PT
182 is for better segmentation. Nevertheless, our method provides synergistic effects for both tasks, and
183 even outperforms SAM-PT for segmentation tasks (see Table 4).

184 4 Experiments

185 In this section, we demonstrate the effectiveness of the proposed TrackIME on point tracking tasks
186 and the downstream video object segmentation. First, to verify the efficacy of the instance motion
187 trajectory estimation and our search space pruning technique for point tracking, we experiment
188 with video scenes that capture dynamic objects. Next, we verify the universality of our method
189 to different point tracking models and find whether it can provide general performance improve-
190 ments when incorporated into the five recent baselines, *e.g.*, TAPNet [11], PIPS2 [9], CoTracker[7],
191 OmniMotion[8] and TAPIR [6]. Finally, we further verify the efficacy of the enhanced point tracking
192 results by TrackIME in the downstream video object segmentation. Specifically, we compare the
193 zero-shot video segmentation performances with the recent SAM-PT [30] baseline which utilizes the
194 point trajectories as the inputs, as well as the conventional baselines that input the semantic classes
195 [27, 28, 29].

196 **Common implementation details.** We note that TrackIME is mainly incorporated with TAPIR
197 point tracker [6] (as it empirically performs best) unless specified otherwise, and we subject it to all
198 experiments including the point tracking and other downstream tasks. For the segmentation model,
199 we utilize the ViT-H variant of Segment Anything (SAM) [1] for the segmentation function. To
200 prepare video frames, we always adjust the resolutions of raw video data to 1080p (1080 pixels in
201 the shorter frame edges), then apply data pre-processing required by individual baseline models. For
202 example, we resize the 1080p frames to 256×256 for TAPIR [6] baseline, following the default
203 setting provided by the official open-source repository. When experimenting TrackIME, we choose
204 the hyperparameters for each baseline, *e.g.*, progressive inference steps $K = 2$, and the pruning sizes
205 $H_0 = W_0 = 960$ and $H_1 = W_1 = 384$ when incorporated with TAPIR [6]. We provide further
206 implementation details in Appendix A.

207 4.1 Point Tracking

208 **Baselines.** We compare our method to the recent baselines OmniMotion[8], CoTracker [7], TAPIR
209 [6], PIPS2 [9], and TAPNet [11]. We utilize the official checkpoints provided by the official project
210 pages and reproduce all experimental results under our common experimental set-up, except for
211 OmniMotion [8] which does not provide checkpoints. Instead, we reproduced the training of
212 OmniMotion models to obtain the experimental results. We use $S = 31$ semantic neighbors to
213 incorporate our framework with the baselines.

214 **Datasets.** We evaluate these models on three different datasets, DAVIS [12], Kinetics [32], and RGB-
215 Stacking [31], each representing different characteristics. For example, DAVIS contains 30 videos
216 specifically curated to evaluate the tracking performance under large variances in the appearance and

Table 2: **Universality of TrackIME with different point tracking models.** We incorporate recent point tracking model baselines [6, 7, 8, 9, 11] with our method, and benchmark its performance on DAVIS [12], RGBStacking [31], and Kinetics [32]. †: the underlined results are obtained with subsets of RGBStacking and Kinetics datasets due to a large optimization cost for the OmniMotion [8].

Method	DAVIS		RGBStacking		Kinetics	
	AJ	δ_{avg}^x	AJ	δ_{avg}^x	AJ	δ_{avg}^x
TAPNet [11]	51.6	63.8	56.5	79.0	49.3	60.7
+ TrackIME	57.9	72.4	66.9	80.0	51.0	63.6
PIPS2 [9]	46.6	69.4	52.3	74.9	-	-
+ TrackIME	50.3	74.0	52.8	75.8	-	-
CoTracker [7]	60.8	76.1	64.1	78.0	47.7	63.7
+ TrackIME	64.5	79.2	68.2	82.1	48.1	63.8
OmniMotion† [8]	52.6	68.1	<u>71.2</u>	<u>81.1</u>	<u>51.0</u>	<u>64.3</u>
+ TrackIME	54.1	69.3	71.9	81.9	51.2	64.6
TAPIR [6]	57.5	70.5	66.3	80.6	50.2	62.3
+ TrackIME	65.3	78.6	66.6	81.8	51.4	65.8

217 motion of object entities. Its two variants, DAVIS-F (First) and DAVIS-S (Strided) differ in how the
 218 query points are given to the models: DAVIS-F queries the model only once in the first frame, while
 219 DAVIS-S queries the model in strides of five frames. Because DAVIS-F requires long-term tracking,
 220 it is generally a more difficult setting. Kinetics contains 1,144 web videos collected from YouTube
 221 that represent realistic noisy characteristics of the video in the wild, such as sudden scene changes.
 222 RGB Stacking is a synthetically rendered dataset representing 50 different moves by a robotic arm.
 223 For all datasets, we refer to the point tracking annotations provided by TAP-Vid [6] and utilize them
 224 as the ground truth for evaluation.

225 **Metric.** To measure the quality of point tracking, we consider point tracking accuracy considered
 226 following TAP-Vid [11], such as the δ -average accuracy (δ_{avg}^x) and the average Jaccard (AJ). The
 227 average metrics are based on the δ -n accuracy (δ_n^x) which indicates the proportion of correct trajectory
 228 sequence as judged by whether they are within the n-pixel error threshold around the ground truth.
 229 In addition, the Jaccard-n (J_n) judges a trajectory sequence to be correct only if the visibility
 230 prediction is also correct. Given these definitions, the average metrics are calculated by averaging
 231 $n \in \{1, 2, 4, 8, 16\}$. To evaluate the fine-grained tracking performance in a harsh error threshold, we
 232 also report δ -1 accuracy (δ_1^x) and Jaccard-1 (J_1). For Table 1, we also discuss the occlusion accuracy
 233 (OA), the proportion of correct visibility sequence given the ground truth.

234 **Effectiveness on point tracking in dynamic objects.** We first present the point tracking scenarios
 235 with dynamic objects. Specifically, we experiment with the DAVIS video scenes [12], which is curated
 236 for evaluating instance motion estimation tasks. As shown in Table 1, we find our method achieves
 237 the best point tracking accuracy surpassing all baselines, *e.g.*, up-to 7.4% relative improvements
 238 in average Jaccard, *i.e.*, 60.8 AJ (CoTracker [7]) vs. 65.3 AJ (TrackIME) when evaluated with the
 239 DAVIS-F (denoted First Query in Table 1). We also measure the occlusion accuracies (OA) and find
 240 a relatively incremental improvement than other metrics. Intuitively, there is a trade-off between
 241 modeling the occlusions among different objects and the search space pruning for one instance, as
 242 the pruning removes information from other instances. Nevertheless, our method is beneficial for
 243 detecting occlusion in fine-grained object parts, and we recommend searching for optimal pruning
 244 parameters that fit a user’s purpose. Finally, we discuss the efficacy of TrackIME under the harsh δ_1^x
 245 and J_1 metrics, where the conventional metrics allows up to 16-pixel errors and takes the average
 246 when judging whether the prediction is correct. For example, the improvement can be even larger,
 247 *e.g.*, up to relative 25.1%, *i.e.*, 28.3 J_1 (CoTracker [7]) vs. 35.4 J_1 (TrackIME) when evaluated with
 248 DAVIS-F. We highlight these benefits of TrackIME allowed by pruning the search space.

249 **Universality to different point tracking models.** We validate the universality of our method when
 250 plugged into the state-of-the-art baselines by evaluating the average tracking accuracy (AJ and δ_n^x)
 251 of the vanilla models and the variants incorporated with our method in Table 2 on DAVIS (First)
 252 [12], RGBStacking [31], and Kinetics [32] datasets. As a result, we observe that our method can

Table 3: **Ablation study of the components in our model.** We ablate the effect of search space pruning (Pruning), trajectory aggregation (Aggregation), and the progressive inference (Progressive) modules for point tracking. We evaluate the tracking benchmark in DAVIS scenes [11, 12].

Pruning	Aggregation	Progressive	J_1	AJ	δ_1^x	δ_{avg}^x
✗	✗	✗	23.0	57.5	34.3	70.5
✓	✗	✗	28.2	62.5	41.1	75.3
✓	✗	✓	28.3	62.6	41.2	75.6
✓	✓	✗	34.0	62.9	48.0	77.0
✓	✓	✓	35.4	65.3	48.2	78.6

Table 4: **Zero-shot video object segmentation performance in DAVIS benchmark.** We consider two set of zero-shot baselines, those utilizing the set of classes [23, 24, 25, 26, 27, 28, 29] and the baseline utilizing a set of query points [30] in a similar manner to our TrackIME. †: we produced the results for TrackIME and SAM-PT [30] under the common set-up, such as the number of tracking points, segmentation function (HQ-SAM [33]), and the same mask formatting for the benchmark.

Method	Input	DAVIS-2017-val			DAVIS-2017-test-dev		
		(J&F) _m	J _m	F _m	(J&F) _m	J _m	F _m
PDB [23]	class	55.1	53.2	57.0	40.4	37.7	43.0
RVOS [24]	class	41.2	36.8	45.7	22.5	17.7	27.3
AGS [25]	class	57.5	55.5	59.5	45.6	42.1	49.0
MAST [26]	class	65.5	63.3	67.6	-	-	-
Propose-Reduce [27]	class	70.4	67.0	73.8	-	-	-
UnOVSOT [28]	class	67.9	66.4	69.3	58.0	54.0	62.0
EntitySeg [29]	class	73.4	70.4	76.4	62.1	-	-
SAM-PT† [30]	points	78.8	76.3	81.3	65.3	62.3	68.3
TrackIME†	points	79.6	76.4	82.8	65.9	62.5	69.4

provide consistent and significant performance improvements in all the baselines, *e.g.*, 13.6% relative improvements (*i.e.*, $57.5 \rightarrow 65.3$ AJ) in TAPIR [6] when evaluated on the DAVIS. Since the model variant incorporated with TAPIR demonstrates the best performance, we chose it as our main model and subjected it to other studies. We note that the experiments for OmniMotion [8] have been conducted in 16 subsets for RGBStacking and Kinetics, and $K = 1$ progressive inference, due to its heavy optimization costs, *e.g.*, approximately 13 gpu-hours for processing one scene. We also note that PIPS2 [9] in Kinetics [12] is unavailable, as its memory requirement for processing Kinetics exceeds our system’s capacity.

Ablation study. We perform an ablation study to understand how each component affects the point trajectory accuracy in Table 3. Specifically, we consider the search space pruning, the trajectory aggregation, and progressive inference modules as the subjects for the ablation.

First of all, we reveal the pure efficacy of our pruning method, separate from the effect of segmentation prior. Notably, when the trajectory aggregation module is removed (the first 2 rows in Table 3), we observe the pruning solely based on the query point’s trajectory provides the most significant effect (*e.g.*, $23.0 \rightarrow 28.2$ in J_1). This validates our key motivation for pruning the search space, which provides superior results even if SAM [1] is not employed.

Next, we discuss the effect of employing SAM [1] by enabling the trajectory aggregation. As expected, aggregating the trajectories for a group of points found in the segmentation mask provides another comparable gain (*e.g.*, $28.2 \rightarrow 34.0$ in J_1), which validates that the aggregation improves the quality of instance trajectory estimation. We additionally note that the progressive inference boosts the performance, (*e.g.*, $34.0 \rightarrow 35.4$ in J_1) when combined with the trajectory aggregation, otherwise the gain is lesser (*e.g.*, $23.2 \rightarrow 28.3$ in J_1). As the progressive inference refers to the estimated instance trajectory, the estimation quality is essential for this module.

We note that further ablation study is available in Appendix D, *e.g.*, the number of semantic neighbors, progressive inference steps, or the pruning sizes.

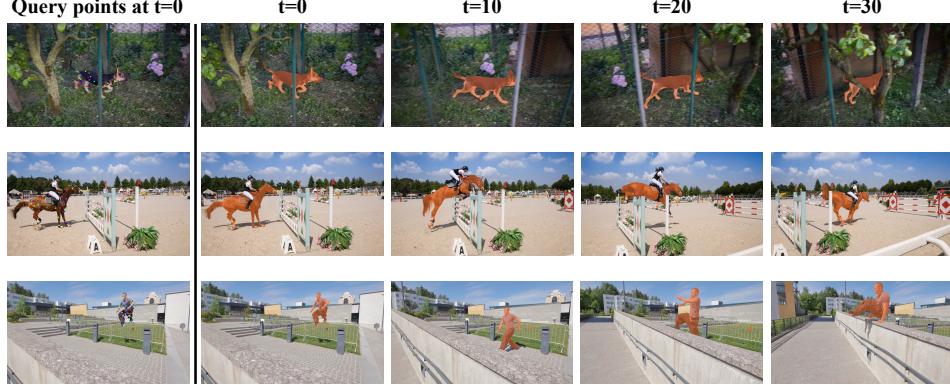


Figure 2: **Demonstration of the video instance segmentation results by our TrackIME framework.** Given the query points in the reference frame, our framework can produce the video instance segmentation masks at quality by performing the weighted aggregation of the mask associated each query point, based on the visibility values.

278 4.2 Video Object Segmentation

279 In this section, we validate the efficacy of TrackIME by performing the zero-shot video segmentation.
 280 We also provide the visualization results for selected scenes from DAVIS [12] in Figure 2.

281 **Baselines.** We experiment with zero-shot video object segmentation to check the efficacy
 282 of TrackIME for improving segmentation. Specifically, we consider the class-guided baselines
 283 for unsupervised video segmentation tasks, *e.g.*, EntitySeg [29]. In addition, we consider the SAM-PT
 284 [30] baseline which also proposes to take point tracking for producing segmentation. To consider the
 285 equivalent experimental set-ups for SAM-PT [30] and TrackIME, we incorporate the models with
 286 HQ-SAM [33] variant for the segmentation, 16 points from the initial frame’s mask, and employ the
 287 iterative refinement technique [33] to produce the video segmentation results.

288 **Datasets.** We evaluate our model on the DAVIS-2017 [12] video segmentation. In particular, we use
 289 the validation and the test-dev sets for the zero-shot benchmark. Both sets contain 30 non-overlapping
 290 scenes with single or multiple objects.

291 **Metrics.** To measure the quality of video instance segmentation, we consider the standard metrics in
 292 baselines: the mean Jaccard (J_m); the mean F-measure (F_m); and the average ($J\&F$) $_m$.

293 **Effectiveness on zero-shot video object segmentation.** In Table 4, we first confirm that the
 294 point tracking provides useful guidance for video segmentation, observing that both SAM-PT [30]
 295 and TrackIME demonstrates significant improvement over the conventional class-prompted baselines.
 296 More importantly, as our framework brings synergistic improvements for both point tracking and
 297 segmentation tasks, we find TrackIME achieves even larger improvement, *e.g.*, 78.8 vs. 79.6 ($J\&F$) $_m$
 298 in the validation set of DAVIS-2017 [12].

299 5 Conclusion

300 In this work, we introduce TrackIME, a novel approach for point tracking to overcome the fundamen-
 301 tal challenge of computation demands in existing models. Specifically, we reduce the search space
 302 by identifying the instance trajectory and pruning the video frames along it. To obtain the instance
 303 trajectory, we aggregate the motion for a group of points on the segmentation masks. To this end, we
 304 propose a unified framework that jointly performs point tracking and segmentation, with the tech-
 305 niques to ensure robustness to occlusion in complex video scenes. TrackIME demonstrates consistent
 306 and significant impacts by bolstering existing point tracking baselines. The joint framework also
 307 reveals the synergistic effects, which also demonstrates the improvements in the video segmentation
 308 task. Overall, our work highlights the effectiveness of considering instance motion trajectory and
 309 jointly solving the tracking and segmentation, and we believe our work could inspire researchers to
 310 consider a new direction to further leverage it in the future.

311 **References**

- 312 [1] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
313 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint*
314 *arXiv:2304.02643*, 2023.
- 315 [2] Haoqi Fan, Yanghao Li, Bo Xiong, Wan-Yen Lo, and Christoph Feichtenhofer. Pyslowfast. <https://github.com/facebookresearch/slowfast>, 2020.
- 316 [3] Zhengqi Li, Qianqian Wang, Forrester Cole, Richard Tucker, and Noah Snavely. Dynibar: Neural dynamic
317 image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
318 *Recognition (CVPR)*, pages 4273–4284, June 2023.
- 320 [4] Jiaben Chen and Huaizu Jiang. Sportsslomo: A new benchmark and baselines for human-centric video
321 frame interpolation. *arXiv preprint arXiv:2308.16876*, 2023.
- 322 [5] Zhoutong Zhang, Forrester Cole, Richard Tucker, William T Freeman, and Tali Dekel. Consistent depth of
323 moving objects in video. *ACM Transactions on Graphics (TOG)*, 40(4):1–12, 2021.
- 324 [6] Carl Doersch, Yi Yang, Mel Vecerik, Dilara Gokay, Ankush Gupta, Yusuf Aytar, Joao Carreira, and Andrew
325 Zisserman. Tapir: Tracking any point with per-frame initialization and temporal refinement. *arXiv preprint*
326 *arXiv:2306.08637*, 2023. URL <https://github.com/google-deepmind/tapnet>.
- 327 [7] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian
328 Rupprecht. Cotracker: It is better to track together. *arXiv preprint arXiv:2307.07635*, 2023. URL
329 <https://github.com/facebookresearch/co-tracker>.
- 330 [8] Qianqian Wang, Yen-Yu Chang, Ruojin Cai, Zhengqi Li, Bharath Hariharan, Aleksander Holynski, and
331 Noah Snavely. Tracking everything everywhere all at once. *arXiv preprint arXiv:2306.05422*, 2023.
- 332 [9] Yang Zheng, Adam W Harley, Bokui Shen, Gordon Wetzstein, and Leonidas J Guibas. Pointodyssey: A
333 large-scale synthetic dataset for long-term point tracking. *arXiv preprint arXiv:2307.15055*, 2023. URL
334 <https://github.com/aharley/pips2>.
- 335 [10] Yuxi Xiao, Qianqian Wang, Shangzhan Zhang, Nan Xue, Sida Peng, Yujun Shen, and Xiaowei Zhou.
336 Spatialtracker: Tracking any 2d pixels in 3d space. In *Proceedings of the IEEE/CVF Conference on*
337 *Computer Vision and Pattern Recognition (CVPR)*, 2024.
- 338 [11] Carl Doersch, Ankush Gupta, Larisa Markeeva, Adrià Recasens, Lucas Smaira, Yusuf Aytar, João
339 Carreira, Andrew Zisserman, and Yi Yang. Tap-vid: A benchmark for tracking any point in a
340 video. *Advances in Neural Information Processing Systems*, 35:13610–13626, 2022. URL <https://github.com/google-deepmind/tapnet>.
- 342 [12] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung, and Luc
343 Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*,
344 2017.
- 345 [13] Christina Karam, Kenjiro Sugimoto, and Keigo Hirakawa. Fast convolutional distance transform. *IEEE*
346 *Signal Processing Letters*, 26(6):853–857, 2019.
- 347 [14] Duc Duy Pham, Gurbandy Dovletov, and Josef Pauli. A differentiable convolutional distance transform
348 layer for improved image segmentation. In *Pattern Recognition: 42nd DAGM German Conference,*
349 *DAGM GCPR 2020, Tübingen, Germany, September 28–October 1, 2020, Proceedings 42*, pages 432–444.
350 Springer, 2021.
- 351 [15] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature
352 pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and*
353 *pattern recognition*, pages 2117–2125, 2017.
- 354 [16] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and
355 Alexander C Berg. Ssd: Single shot multibox detector. In *Computer Vision–ECCV 2016: 14th European*
356 *Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, pages 21–37.
357 Springer, 2016.
- 358 [17] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick
359 Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional
360 networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2758–2766,
361 2015.

- 362 [18] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox.
 363 Flownet 2.0: Evolution of optical flow estimation with deep networks. In *Proceedings of the IEEE*
 364 *conference on computer vision and pattern recognition*, pages 2462–2470, 2017.
- 365 [19] Jia Xu, René Ranftl, and Vladlen Koltun. Accurate optical flow via direct cost volume processing. In
 366 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1289–1297,
 367 2017.
- 368 [20] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid,
 369 warping, and cost volume. In *Proceedings of the IEEE conference on computer vision and pattern*
 370 *recognition*, pages 8934–8943, 2018.
- 371 [21] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer*
 372 *Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II*
 373 *16*, pages 402–419. Springer, 2020.
- 374 [22] Adam W Harley, Zhaoyuan Fang, and Katerina Fragkiadaki. Particle video revisited: Tracking through
 375 occlusions using point trajectories. In *European Conference on Computer Vision*, pages 59–75. Springer,
 376 2022.
- 377 [23] Hongmei Song, Wenguan Wang, Sanyuan Zhao, Jianbing Shen, and Kin-Man Lam. Pyramid dilated deeper
 378 convlstm for video salient object detection. In *Proceedings of the European conference on computer vision*
 379 (*ECCV*), pages 715–731, 2018.
- 380 [24] Carles Ventura, Miriam Bellver, Andreu Girbau, Amaia Salvador, Ferran Marques, and Xavier Giro-i Nieto.
 381 Rvos: End-to-end recurrent network for video object segmentation. In *Proceedings of the IEEE/CVF*
 382 *conference on computer vision and pattern recognition*, pages 5277–5286, 2019.
- 383 [25] Wenguan Wang, Hongmei Song, Shuyang Zhao, Jianbing Shen, Sanyuan Zhao, Steven CH Hoi, and Haibin
 384 Ling. Learning unsupervised video object segmentation through visual attention. In *Proceedings of the*
 385 *IEEE/CVF conference on computer vision and pattern recognition*, pages 3064–3074, 2019.
- 386 [26] Zihang Lai, Erika Lu, and Weidi Xie. Mast: A memory-augmented self-supervised tracker. In *Proceedings*
 387 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6479–6488, 2020.
- 388 [27] Huaijia Lin, Ruizheng Wu, Shu Liu, Jiangbo Lu, and Jiaya Jia. Video instance segmentation with a
 389 propose-reduce paradigm. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
 390 pages 1739–1748, 2021.
- 391 [28] Jonathon Luiten, Idil Esen Zulfikar, and Bastian Leibe. Unovost: Unsupervised offline video object
 392 segmentation and tracking. In *Proceedings of the IEEE/CVF winter conference on applications of*
 393 *computer vision*, pages 2000–2009, 2020.
- 394 [29] Lu Qi, Jason Kuen, Weidong Guo, Tiancheng Shen, Jiuxiang Gu, Wenbo Li, Jiaya Jia, Zhe Lin, and
 395 Ming-Hsuan Yang. Fine-grained entity segmentation. *arXiv preprint arXiv:2211.05776*, 2022.
- 396 [30] Frano Rajič, Lei Ke, Yu-Wing Tai, Chi-Keung Tang, Martin Danelljan, and Fisher Yu. Segment anything
 397 meets point tracking. *arXiv preprint arXiv:2307.01197*, 2023.
- 398 [31] Alex X. Lee, Coline Devin, Yuxiang Zhou, Thomas Lampe, Konstantinos Bousmalis, Jost Tobias
 399 Springenberg, Arunkumar Byravan, Abbas Abdolmaleki, Nimrod Gileadi, David Khosid, Claudio Fan-
 400 tacci, Jose Enrique Chen, Akhil Raju, Rae Jeong, Michael Neunert, Antoine Laurens, Stefano Sal-
 401 iceti, Federico Casarini, Martin Riedmiller, Raia Hadsell, and Francesco Nori. Beyond pick-and-place:
 402 Tackling robotic stacking of diverse shapes. In *Conference on Robot Learning (CoRL)*, 2021. URL
 403 <https://openreview.net/forum?id=U0Q8CrtBJxJ>.
- 404 [32] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset.
 405 In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308,
 406 2017.
- 407 [33] Lei Ke, Mingqiao Ye, Martin Danelljan, Yifan Liu, Yu-Wing Tai, Chi-Keung Tang, and Fisher Yu. Segment
 408 anything in high quality. *arXiv preprint arXiv:2306.01567*, 2023.
- 409 [34] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin,
 410 George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: compos-
 411 able transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.

- 412 [35] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
413 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf,
414 Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit
415 Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-
416 performance deep learning library. In *Advances in Neural Information Processing Systems 32*,
417 pages 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- 418
- 419 [36] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J Fleet,
420 Dan Gnanapragasam, Florian Golemo, Charles Herrmann, et al. Kubric: A scalable dataset generator. In
421 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3749–3761,
422 2022.

423 Appendix

424 A Experimental details for point tracking

425 In this section, we present detailed experimental setups considered by our experiments in Section 4.

426 **Baselines.** We consider 5 different baseline point tracking models, TAPNet [11], PIPS2 [9], CoTracker
 427 [7], OmniMotion [8], and TAPIR [6]. We experiment with the checkpoint provided in the official
 428 open-source repository hosted by their authors, following the default hyperparameters in each
 429 model, *e.g.*, for the input dimensions, in TAPNet [11] and TAPIR [6] consider a square-shaped
 430 (256×256) dimension, while PIPS2 and CoTracker do rectangular-shaped dimensions, (512×896)
 431 and (384×512), respectively. We also note that the backend library of TAPIR and TAPNet is ported
 432 from JAX [34] to PyTorch [35] in our experiments, which provides subtle enhancements in the
 433 tracking accuracy, *e.g.*, AJ 56.2 [6] \rightarrow 57.5 (Table 2) in DAVIS-F.

Table 5: **The pruned resolutions in our method for each baseline point tracking model.** We report the specific values for H_0 , H_1 , W_0 , W_1 when TAPNet[11], PIPS2[9], CoTracker[7], OmniMotion [8], or TAPIR[6] is used as the baseline.

Baseline Model	H_0	H_1	W_0	W_1
TAPNet [11]	960	384	960	384
PIPS2 [9]	960	512	1680	896
CoTracker [7]	960	384	1280	512
OmniMotion [8]	960	—	960	—
TAPIR [6]	960	384	960	384

434 **Hyperparameters for point tracking.** Unless otherwise specified, we always choose the number of
 435 semantic neighbors $S = 31$, and the progressive inference steps $K = 2$ for TAPNet [11], PIPS2 [9],
 436 CoTracker [7], and TAPIR [6]. For OmniMotion [8], we set the progressive inference step $K = 1$.
 437 To incorporate our framework with the baselines, we select different pruning sizes to meet the shape
 438 requirements of a specific model (*e.g.*, TAPIR [6] needs a shape in multiples of 8 to be compatible
 439 with its convolution layers). For example, if our method is plugged into TAPIR [6] and a video
 440 with the 1080p resolution, *e.g.*, (1080×1920), we set the pruning resolutions $H_0 = W_0 = 960$ and
 441 $H_1 = W_1 = 384$. For clarity, we present the resolutions for all baseline models in Table 5.

442 **Datasets.** We evaluate the baselines and TrackIME in three different datasets from the TAP-Vid
 443 benchmark [11]: DAVIS [12]; Kinetics [32]; and RGBStacking [8]. The sizes of the raw samples
 444 can vary, *e.g.*, from 256 to 2160 in their shorter sides, hence we process the frames by resizing
 445 the shorter sides to 1080 with the aspect ratio fixed. As a result, the video frame resolutions are
 446 typically (1080×1920) for DAVIS [12] and Kinetics [32]. We note that RGB-Stacking is originally
 447 in (256×256), but we do bilinear up-sampling to (1080×1080) for simplicity.

448 **Experimental environment.** Every baseline model and internal module in TrackIME (*e.g.*, Segment
 449 Anything [1]) is implemented in PyTorch 2.1 [35] compiled for CUDA 11.8, which we run on an
 450 NVIDIA RTX 3090 GPU. In default, we experiment with the float32 numerical precision; however,
 451 in case of out-of-memory errors (*e.g.*, RTX 3090’s 24 GiB VRAM cannot handle hundreds of frames),
 452 we employ the bfloat16 precision to fit such samples into the limited memory.

453 B Backgrounds

454 In this section, we describe technical details behind the limitations in the current point tracking
 455 models.

456 In the common canonical design of recent model architectures for Equation (1), *e.g.*, our baselines:
 457 TAPNet [11], CoTracker [7], TAPIR [6], etc., the key component is the cost volume [21], which
 458 represents the likelihood of the query point’s spatial-temporal location over the entire video frames.
 459 In principle, predicting this cost map requires a brute-force search over every spatial-temporal
 460 location, which is often computationally infeasible on the raw video dimensions, *e.g.*, 1080p. To
 461 mitigate this problem, current models first down-sample the raw video into a lower spatial resolutions,

Table 6: **FLOP counts by each module in TrackIME.** We report the FLOP counts for point tracking given 64 video frames, during the instance motion stage, and the high-fidelity tracking with $K = 2$ progressive steps.

TrackIME Modules	FLOPs
Instance trajectory estimation (stage 1)	1355G
- Segmentation	533G
- Instance Tracking (32 points; $S = 31$)	822G
Progressive inference (stage 2)	1434G
- $k = 0$ (32 points; $S = 31$)	822G
- $k = 1$ (1 point)	612G
Total	2789G

462 *e.g.*, (256×256) in TAPIR [6]. While the reduced resolution enables models to process the entire
463 video frames for tracking, the lost information during the resolution reduction induces quantization
464 noises into the cost volume. Recent baselines, including the state-of-the-art [6], employ refinement
465 techniques to mitigate these noises.⁴ Nevertheless, the lost detail in the visual feature after the
466 down-sampling still hinder representing high-frequency patterns, and the model can suffer from
467 tracking failure modes.

468 In this regard, our method pursues the direction of pruning the excessive search space for point
469 tracking, so that models can avoid the down-sampling and focus only on important regions maintaining
470 detailed visual features.

471 C Computational costs for point tracking

472 In this section, we study the computational costs and efficiency of TrackIME by examining the FLOP
473 (floating-point operations) counts for performing the point tracking.

474 **FLOP count of TrackIME.** To check the exact cost of each module in TrackIME, we report the
475 FLOPs for tracking under our default setting, *e.g.*, TAPIR [6] as the baseline, given 64 video frames.
476 Specifically, as given in Table 6, the segmentation with SAM [1] needs 533 GFLOPs, tracking 32
477 points (*e.g.*, $S = 31$ semantic neighbors plus one query point) demand 822 GFLOPs, and tracking a
478 single point demands 612 GFLOPs, respectively. As a result, the net FLOP count of TrackIME (with
479 $K = 2$ progressive steps) is 2789 GFLOPs.

480 **Computation efficiency compared to baselines.** Next, we compare the baseline TAPIR [6] with
481 various input dimensions and TrackIME, in terms of their FLOP counts versus the point tracking
482 performances, AJ (Average Jaccard), δ_{avg}^x , and OA (Occlusion Accuracy), evaluated under DAVIS-F
483 and DAVIS-S in Table 7.

484 For TAPIR, the FLOP count is mostly governed by the input dimension of a model (256×256) , *e.g.*,
485 612 GFLOPs for processing 64 video frames, and it grows quadratically as the input dimension gets
486 increased.

487 An interesting finding in Table 7 is that the baseline [6] cannot benefit from the larger input dimensions
488 without fine-tuning. For example, we observe that the baseline’s performance only deteriorates given
489 larger inputs, as the model is only optimized for a low-resolution input frames (256×256) to meet
490 the memory constraints while training; it is non-trivial to process high-resolution inputs without
491 fine-tuning. Furthermore, even if fine-tuning is employed (*e.g.*, TAPIR Hi-Res [6]), the performance
492 gain (*e.g.*, $62.8 \rightarrow 65.7$ AJ) is not significant considering the excessive increase in FLOP counts (*e.g.*,
493 612 \rightarrow 8257 GFLOPs), and the occlusion accuracy (OA) can even get worse (*e.g.*, $88.3 \rightarrow 86.7$).

494 These results further demonstrate the merits of employing TrackIME for point tracking, which can
495 enable point tracking models to process the frames in a computationally efficient manner, even without
496 fine-tuning, and provide consistent performance gains. For example, comparing TrackIME (ours) vs.

⁴We refer the readers to literature for the refinement mechanisms [6, 7].

Table 7: **The comparison of the FLOP counts of the TAPIR [6] models and TrackIME.** We report the FLOP counts to process 64 video frames by TAPIR with the input dimensions (256×256) (default), (512×512), and (768×768), TAPIR Hi-Res (a fine-tuned model for (1080×1080)) and TrackIME (ours). For each model, we further report the benchmark results in terms of AJ (Average Jaccard), δ_{avg}^x , and OA (Occlusion Accuracy), evaluated under DAVIS-F and DAVIS-S. For TAPIR Hi-Res, numbers are excerpted from [6], where results for DAVIS-F are not available.

Method (Input Dim.)	FLOPs	DAVIS-F			DAVIS-S		
		AJ	δ_{avg}^x	OA	AJ	δ_{avg}^x	OA
TAPIR (256×256)	612G	57.5	70.5	85.5	62.8	75.1	88.3
TAPIR (512×512)	2429G	53.9	65.9	79.8	62.5	74.0	81.8
TAPIR (768×768)	5457G	53.3	65.5	73.2	58.3	70.2	76.6
TAPIR Hi-Res (1080×1080)	8257G	-	-	-	65.7	77.6	86.7
TrackIME (256×256)	2789G	65.3	78.6	86.5	69.3	81.4	89.0

497 TAPIR Hi-Res [6] gives: **2789G** vs. 8257G (FLOPs); **69.3** vs. 65.7 (AJ); **81.4** vs. 77.6 (δ_{avg}^x); and
498 **89.0** vs. 86.7 (OA), in DAVIS-S, respectively.

499 D Ablation study

500 In this section, we ablate the choice of hyperparameters in our enhanced point tracking, namely
501 the pruning sizes (H_0, W_0) without the progressive fusion (*i.e.*, $K = 1$) and our default setting
502 in TrackIME ($K = 2$), and the number of sampling semantic neighbors S for estimating the instance
503 trajectory.

504 In Table 8, we find that smaller pruning sizes tend to introduce positive effects in the fine-grained
505 metrics (*e.g.*, 1- and 2-pixel error thresholds), but also trade off the average-scale metrics (*e.g.*, AJ
506 and δ_{avg}^x). These results are expected, as the pruning size gets smaller, the amount of down-sampling
507 reduces and more detailed visual features would be preserved, but at the same time, the chance of
508 erroneous pruning increases where the true location of the query point is lost.

509 We note that the progressive fusion ($K = 2$) in our method can mitigate the trade-off in pruning by
510 considering multiple scales, *e.g.*, $H_0 = 960$ and $H_1 = 384$, providing additional performance gains.

511 Next, in Table 9, we ablate the effect of the choice for the number of semantic neighbors $S + 1$
512 (including the query point), halving down its value starting from $(S + 1) = 128$ to $(S + 1) = 2$. As
513 a result, we find that all of the choices can provide satisfactory performance in general, although
514 there exist mild trade-offs between the 1- and 2-pixel scale metrics and the average scale metrics. As
515 one of our goal is on achieving the optimal pixel-scale performance in point tracking, we empirically
516 choose $(S + 1) = 32$, which reveals the best 1-pixel scale metrics.

517 E Additional experiments and visualizations

518 In this section, we provide the additional experiment and visualizations with TrackIME.

519 **Fine-tuning SAM with pseudo labels from video segmentation by TrackIME.** In Table 10, we
520 experiment the fine-tuning of SAM [1] and HQ-SAM [33] models with the pseudo labels produced
521 by TrackIME. Specifically, we utilize the masks produced by TrackIME for the training set of
522 DAVIS-2017 [12] as the pseudo label for training SAM [1] and HQ-SAM [33] and measure the mean
523 intersection-over-union (mIoU) in the validation set of the DAVIS scenes. For the optimization, we
524 utilize the efficient tuning architecture for SAM [33] with a batch size of 4, using the pseudo labels
525 and a learning rate of 1e-3 with a decay factor of 0.1 after 10 epochs, over a total of 12 epochs.

Table 8: **Ablation study of the pruning size in our framework.** We ablate the pruning size considered in TrackIME. For the evaluation, we calculate both pixel-scale and average-scale metrics under the DAVIS-F dataset [11].

Pruning Size (K)	J ₁	δ_1^x	J ₂	δ_2^x	AJ	δ_{avg}^x
1080 ($K = 1$)	28.1	41.0	52.3	66.0	62.5	75.2
960 ($K = 1$)	29.7	42.4	52.9	66.3	63.1	75.6
768 ($K = 1$)	31.9	44.6	55.7	68.1	64.0	76.4
512 ($K = 1$)	34.6	47.6	56.5	68.9	63.9	76.9
384 ($K = 1$)	35.3	47.3	56.5	68.3	62.3	76.1
$960 \rightarrow 384$ ($K = 2$)	35.4	48.2	57.7	70.1	65.3	78.6

Table 9: **Ablation study of the effect of the number of semantic neighbors in our method.** We ablate the number of semantic neighbors considered in our method. For the evaluation, we calculate both pixel-scale and average-scale metrics under the DAVIS-F dataset [11].

$S + 1$	J ₁	δ_1^x	J ₂	δ_2^x	AJ	δ_{avg}^x
128	35.1	47.7	57.1	69.4	64.8	78.2
64	35.0	47.5	57.2	69.8	64.8	78.3
32	35.4	48.2	57.7	70.1	65.3	78.6
16	35.2	47.9	57.3	69.9	64.8	78.5
8	35.1	47.8	57.4	70.2	64.9	78.5
4	35.0	47.6	57.5	69.8	64.9	78.3
2	35.1	47.9	57.2	69.7	64.7	78.1

526 Both the models demonstrate approximately the same mIoUs before the fine-tuning, where the gains
527 are observed after tuning with the zero-shot segmentation results provided by TrackIME, *e.g.*, 80.00
528 → 83.02 in SAM [1]. This results suggest the potential application of TrackIME for obtaining
529 pseudo-labels from unlabeled video data, which can be used for training segmentation models.

Table 10: **Fine-tuning SAM and HQ-SAM with pseudo labels.** We measure the mean intersection-over-union (mIoU) in the validation set of the DAVIS scenes.

Model	Before tuning	After tuning
SAM [1]	80.00	83.02
HQ-SAM [33]	79.86	83.08

530 **Visualization of the progressive inference.** We additionally visualize the progressive inference
531 structure in Appendix E. Specifically, we incorporated TrackIME with TAPIR [6] and apply the
532 progressive pruning sizes of (960×960) and (384×384) . As depicted by Appendix E, the latest
533 progressive step is well focused around the query point, *e.g.*, the dog’s ear, so that the search space
534 for point tracking is effectively pruned.



Figure 3: Demonstration of the progressive inference by TrackIME framework.

535 **F Limitation**

536 **F.1 Limitation and Future Works**

537 TrackIME relies on the pre-trained models for point tracking, often trained with synthetic datasets,
538 such as Kubric [36] and PointOdyssey [9], while the segmentation models are primarily trained on
539 the real images [1]. An interesting future direction is to integrate the TrackIME with training on the
540 real video scenes. As this could include the development of point tracking algorithms capable of
541 generalization to diverse intricate objects or, alternatively, optimizing the segmentation models for
542 better video scene understanding. This approach could further improve the accuracy and applicability
543 of TrackIME in various real-world scenarios.

544 **F.2 Potential Negative Societal Impact**

545 While point tracking by TrackIME can be beneficial for various video understanding applications,
546 such as novel-view synthesis, depth estimation, and action recognition, the emergence of unexpected
547 behavior within TrackIME can lead to misrepresentations of the real video data. For those applications
548 that require extremely accurate models for safety-related judgements, such as depth estimation for
549 autonomous driving, the unexpected behaviors must be carefully managed. To ensure the reliability
550 of systems using point tracking predictions, we recommend to conduct thorough investigations and
551 implement robust mitigation strategies to minimize potential risks, thereby increasing the overall
552 safety and effectiveness of these applications.

553 **NeurIPS Paper Checklist**

554 **1. Claims**

555 Question: Do the main claims made in the abstract and introduction accurately reflect the
556 paper's contributions and scope?

557 Answer: [Yes]

558 Justification: All claims in the introduction and abstract accurately reflect the contribution
559 and scope, which are then verified in Section 4.

560 Guidelines:

- 561 • The answer NA means that the abstract and introduction do not include the claims
562 made in the paper.
- 563 • The abstract and/or introduction should clearly state the claims made, including the
564 contributions made in the paper and important assumptions and limitations. A No or
565 NA answer to this question will not be perceived well by the reviewers.
- 566 • The claims made should match theoretical and experimental results, and reflect how
567 much the results can be expected to generalize to other settings.
- 568 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
569 are not attained by the paper.

570 **2. Limitations**

571 Question: Does the paper discuss the limitations of the work performed by the authors?

572 Answer: [Yes]

573 Justification: Appendix F discusses it. The trade-offs regarding the hyperparameter selection
574 is discussed in Section 4.

575 Guidelines:

- 576 • The answer NA means that the paper has no limitation while the answer No means that
577 the paper has limitations, but those are not discussed in the paper.
- 578 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 579 • The paper should point out any strong assumptions and how robust the results are to
580 violations of these assumptions (e.g., independence assumptions, noiseless settings,
581 model well-specification, asymptotic approximations only holding locally). The authors
582 should reflect on how these assumptions might be violated in practice and what the
583 implications would be.
- 584 • The authors should reflect on the scope of the claims made, e.g., if the approach was
585 only tested on a few datasets or with a few runs. In general, empirical results often
586 depend on implicit assumptions, which should be articulated.
- 587 • The authors should reflect on the factors that influence the performance of the approach.
588 For example, a facial recognition algorithm may perform poorly when image resolution
589 is low or images are taken in low lighting. Or a speech-to-text system might not be
590 used reliably to provide closed captions for online lectures because it fails to handle
591 technical jargon.
- 592 • The authors should discuss the computational efficiency of the proposed algorithms
593 and how they scale with dataset size.
- 594 • If applicable, the authors should discuss possible limitations of their approach to
595 address problems of privacy and fairness.
- 596 • While the authors might fear that complete honesty about limitations might be used by
597 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
598 limitations that aren't acknowledged in the paper. The authors should use their best
599 judgment and recognize that individual actions in favor of transparency play an impor-
600 tant role in developing norms that preserve the integrity of the community. Reviewers
601 will be specifically instructed to not penalize honesty concerning limitations.

602 **3. Theory Assumptions and Proofs**

603 Question: For each theoretical result, does the paper provide the full set of assumptions and
604 a complete (and correct) proof?

605 Answer: [NA]

606 Justification: We do not have a theory in this paper.

607 Guidelines:

- 608 • The answer NA means that the paper does not include theoretical results.
- 609 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 610 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 611 • The proofs can either appear in the main paper or the supplemental material, but if
- 612 • they appear in the supplemental material, the authors are encouraged to provide a short
- 613 • proof sketch to provide intuition.
- 614 • Inversely, any informal proof provided in the core of the paper should be complemented
- 615 • by formal proofs provided in appendix or supplemental material.
- 616 • Theorems and Lemmas that the proof relies upon should be properly referenced.

617 4. Experimental Result Reproducibility

619 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
620 perimental results of the paper to the extent that it affects the main claims and/or conclusions
621 of the paper (regardless of whether the code and data are provided or not)?

622 Answer: [Yes]

623 Justification: We have included the implementation details of TrackIME in Section 4
624 and Appendix A.

625 Guidelines:

- 626 • The answer NA means that the paper does not include experiments.
- 627 • If the paper includes experiments, a No answer to this question will not be perceived
- 628 • well by the reviewers: Making the paper reproducible is important, regardless of
- 629 • whether the code and data are provided or not.
- 630 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 631 • to make their results reproducible or verifiable.
- 632 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 633 • For example, if the contribution is a novel architecture, describing the architecture fully
- 634 • might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 635 • be necessary to either make it possible for others to replicate the model with the same
- 636 • dataset, or provide access to the model. In general, releasing code and data is often
- 637 • one good way to accomplish this, but reproducibility can also be provided via detailed
- 638 • instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 639 • of a large language model), releasing of a model checkpoint, or other means that are
- 640 • appropriate to the research performed.
- 641 • While NeurIPS does not require releasing code, the conference does require all submis-
- 642 • sions to provide some reasonable avenue for reproducibility, which may depend on the
- 643 • nature of the contribution. For example
 - 644 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
 - 645 • to reproduce that algorithm.
 - 646 (b) If the contribution is primarily a new model architecture, the paper should describe
 - 647 • the architecture clearly and fully.
 - 648 (c) If the contribution is a new model (e.g., a large language model), then there should
 - 649 • either be a way to access this model for reproducing the results or a way to reproduce
 - 650 • the model (e.g., with an open-source dataset or instructions for how to construct
 - 651 • the dataset).
 - 652 (d) We recognize that reproducibility may be tricky in some cases, in which case
 - 653 • authors are welcome to describe the particular way they provide for reproducibility.
 - 654 • In the case of closed-source models, it may be that access to the model is limited in
 - 655 • some way (e.g., to registered users), but it should be possible for other researchers
 - 656 • to have some path to reproducing or verifying the results.

657 5. Open access to data and code

658 Question: Does the paper provide open access to the data and code, with sufficient instruc-
659 tions to faithfully reproduce the main experimental results, as described in supplemental
660 material?

661 Answer: [No]

662 Justification: We will make a decision for this after the acceptance.

663 Guidelines:

- 664 • The answer NA means that paper does not include experiments requiring code.
- 665 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 666 • While we encourage the release of code and data, we understand that this might not be
667 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
668 including code, unless this is central to the contribution (e.g., for a new open-source
669 benchmark).
- 670 • The instructions should contain the exact command and environment needed to run to
671 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 672 • The authors should provide instructions on data access and preparation, including how
673 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 674 • The authors should provide scripts to reproduce all experimental results for the new
675 proposed method and baselines. If only a subset of experiments are reproducible, they
676 should state which ones are omitted from the script and why.
- 677 • At submission time, to preserve anonymity, the authors should release anonymized
678 versions (if applicable).
- 679 • Providing as much information as possible in supplemental material (appended to the
680 paper) is recommended, but including URLs to data and code is permitted.

683 6. Experimental Setting/Details

684 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
685 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
686 results?

687 Answer: [Yes]

688 Justification: We provide the detail of the training/evaluation setup, dataset, and hyperpa-
689 rameters in Section 4 and Appendix A.

690 Guidelines:

- 691 • The answer NA means that the paper does not include experiments.
- 692 • The experimental setting should be presented in the core of the paper to a level of detail
693 that is necessary to appreciate the results and make sense of them.
- 694 • The full details can be provided either with the code, in appendix, or as supplemental
695 material.

696 7. Experiment Statistical Significance

697 Question: Does the paper report error bars suitably and correctly defined or other appropriate
698 information about the statistical significance of the experiments?

699 Answer: [No]

700 Justification: All experiments are conducted with the same and commonly used random
701 seed.

702 Guidelines:

- 703 • The answer NA means that the paper does not include experiments.
- 704 • The authors should answer “Yes” if the results are accompanied by error bars, confi-
705 dence intervals, or statistical significance tests, at least for the experiments that support
706 the main claims of the paper.
- 707 • The factors of variability that the error bars are capturing should be clearly stated (for
708 example, train/test split, initialization, random drawing of some parameter, or overall
709 run with given experimental conditions).

- 710 • The method for calculating the error bars should be explained (closed form formula,
 711 call to a library function, bootstrap, etc.)
 712 • The assumptions made should be given (e.g., Normally distributed errors).
 713 • It should be clear whether the error bar is the standard deviation or the standard error
 714 of the mean.
 715 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 716 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 717 of Normality of errors is not verified.
 718 • For asymmetric distributions, the authors should be careful not to show in tables or
 719 figures symmetric error bars that would yield results that are out of range (e.g. negative
 720 error rates).
 721 • If error bars are reported in tables or plots, The authors should explain in the text how
 722 they were calculated and reference the corresponding figures or tables in the text.

723 **8. Experiments Compute Resources**

724 Question: For each experiment, does the paper provide sufficient information on the com-
 725 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 726 the experiments?

727 Answer: [Yes]

728 Justification: We provide the computational costs in Appendix C.

729 Guidelines:

- 730 • The answer NA means that the paper does not include experiments.
- 731 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 732 or cloud provider, including relevant memory and storage.
- 733 • The paper should provide the amount of compute required for each of the individual
 734 experimental runs as well as estimate the total compute.
- 735 • The paper should disclose whether the full research project required more compute
 736 than the experiments reported in the paper (e.g., preliminary or failed experiments that
 737 didn't make it into the paper).

738 **9. Code Of Ethics**

739 Question: Does the research conducted in the paper conform, in every respect, with the
 740 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

741 Answer: [Yes]

742 Justification: We do not have any ethical concerns regarding the paper.

743 Guidelines:

- 744 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 745 • If the authors answer No, they should explain the special circumstances that require a
 746 deviation from the Code of Ethics.
- 747 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 748 eration due to laws or regulations in their jurisdiction).

749 **10. Broader Impacts**

750 Question: Does the paper discuss both potential positive societal impacts and negative
 751 societal impacts of the work performed?

752 Answer: [Yes]

753 Justification: We discussed the societal impact in Appendix F

754 Guidelines:

- 755 • The answer NA means that there is no societal impact of the work performed.
- 756 • If the authors answer NA or No, they should explain why their work has no societal
 757 impact or why the paper does not address societal impact.
- 758 • Examples of negative societal impacts include potential malicious or unintended uses
 759 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
 760 (e.g., deployment of technologies that could make decisions that unfairly impact specific
 761 groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our method does not introduce risks for misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We have cited all papers and datasets used in the paper.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- 814 • If this information is not available online, the authors are encouraged to reach out to
815 the asset's creators.

816 **13. New Assets**

817 Question: Are new assets introduced in the paper well documented and is the documentation
818 provided alongside the assets?

819 Answer: [No]

820 Justification: We will release the Pytorch implementation of TrackIME after the acceptance.

821 Guidelines:

- 822 • The answer NA means that the paper does not release new assets.
823 • Researchers should communicate the details of the dataset/code/model as part of their
824 submissions via structured templates. This includes details about training, license,
825 limitations, etc.
826 • The paper should discuss whether and how consent was obtained from people whose
827 asset is used.
828 • At submission time, remember to anonymize your assets (if applicable). You can either
829 create an anonymized URL or include an anonymized zip file.

830 **14. Crowdsourcing and Research with Human Subjects**

831 Question: For crowdsourcing experiments and research with human subjects, does the paper
832 include the full text of instructions given to participants and screenshots, if applicable, as
833 well as details about compensation (if any)?

834 Answer: [NA]

835 Justification: We use existing benchmark datasets and do not have any crowdsourcing
836 datasets or experiments in the paper.

837 Guidelines:

- 838 • The answer NA means that the paper does not involve crowdsourcing nor research with
839 human subjects.
840 • Including this information in the supplemental material is fine, but if the main contribu-
841 tion of the paper involves human subjects, then as much detail as possible should be
842 included in the main paper.
843 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
844 or other labor should be paid at least the minimum wage in the country of the data
845 collector.

846 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
847 Subjects**

848 Question: Does the paper describe potential risks incurred by study participants, whether
849 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
850 approvals (or an equivalent approval/review based on the requirements of your country or
851 institution) were obtained?

852 Answer: [NA]

853 Justification: We do not have human subject in the research.

854 Guidelines:

- 855 • The answer NA means that the paper does not involve crowdsourcing nor research with
856 human subjects.
857 • Depending on the country in which research is conducted, IRB approval (or equivalent)
858 may be required for any human subjects research. If you obtained IRB approval, you
859 should clearly state this in the paper.
860 • We recognize that the procedures for this may vary significantly between institutions
861 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
862 guidelines for their institution.
863 • For initial submissions, do not include any information that would break anonymity (if
864 applicable), such as the institution conducting the review.