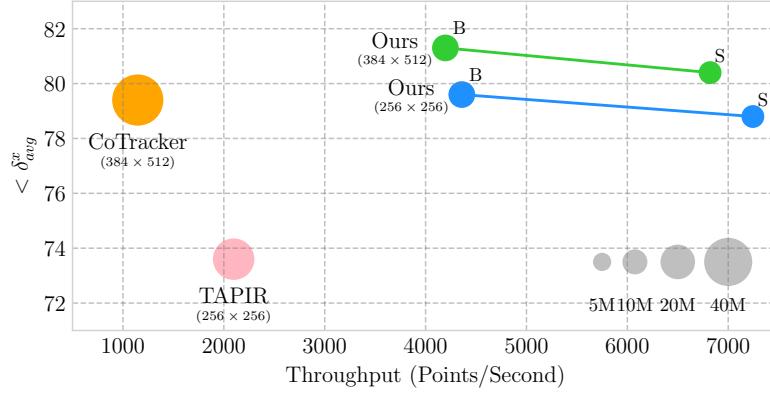


# Local All-Pair Correspondence for Point Tracking

Seokju Cho<sup>1</sup>, Jiahui Huang<sup>2</sup>, Jisu Nam<sup>1</sup>, Honggyu An<sup>1</sup>,  
Seungryong Kim<sup>1,†</sup>, and Joon-Young Lee<sup>2,†</sup>

<sup>1</sup> Korea University <sup>2</sup> Adobe Research



**Fig. 1: Evaluating LocoTrack against state-of-the-art methods.**

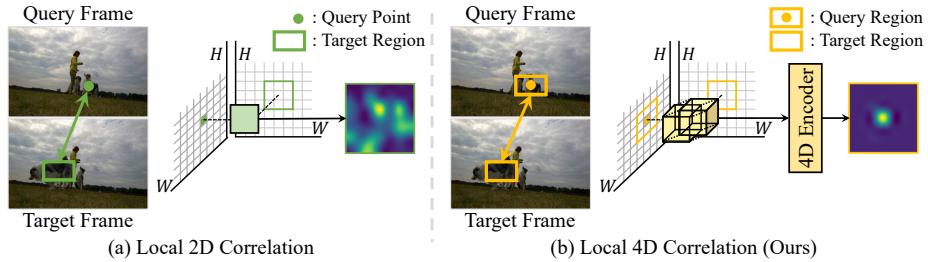
We compare our LocoTrack against other SOTA methods [12,24] in terms of model size (circle size), position accuracy (y-axis), and throughput (x-axis). LocoTrack shows exceptionally high precision and efficiency.

**Abstract.** We introduce LocoTrack, a highly accurate and efficient model designed for the task of tracking any point (TAP) across video sequences. Previous approaches in this task often rely on local 2D correlation maps to establish correspondences from a point in the query image to a local region in the target image, which often struggle with homogeneous regions or repetitive features, leading to matching ambiguities. LocoTrack overcomes this challenge with a novel approach that utilizes all-pair correspondences across regions, *i.e.*, local 4D correlation, to establish precise correspondences, with bidirectional correspondence and matching smoothness significantly enhancing robustness against ambiguities. We also incorporate a lightweight correlation encoder to enhance computational efficiency, and a compact Transformer architecture to integrate long-term temporal information. LocoTrack achieves unmatched accuracy on all TAP-Vid benchmarks and operates at a speed almost 6× faster than the current state-of-the-art.

## 1 Introduction

Finding corresponding points across different views of a scene, a process known as point correspondence [1, 31, 68], is one of fundamental problems in computer

<sup>†</sup>Co-corresponding authors.



**Fig. 2: Illustration of our core component.** Our local all-pair formulation, achieved with local 4D correlation, demonstrates robustness against matching ambiguity. This contrasts with previous works [12, 15, 24, 63] that rely on point-to-region correspondences, achieved with local 2D correlation, which are susceptible to the ambiguity.

vision, which has a variety of applications such as 3D reconstruction [34, 49], autonomous driving [21, 40], and pose estimation [47–49]. Recently, the emerging point tracking task [11, 15] addresses the point correspondence across a video. Given an input video and a query point on a physical surface, the task aims to find the corresponding position of the query point for every target frame along with its visibility status. This task demands a sophisticated understanding of motion over time and a robust capability for matching points accurately.

Recent methods in this task often rely on constructing a 2D local correlation map [12, 15, 24, 63], comparing the deep features of a query point with a local region of the target frame to predict the corresponding positions. However, this approach encounters substantial difficulties in precisely identifying positions within homogeneous areas, regions with repetitive patterns, or differentiating among co-occurring objects [46, 57, 68]. To resolve matching ambiguities that arise in these challenging scenarios, establishing effective correspondence between frames is crucial. Existing works attempt to resolve these ambiguities by considering the temporal context [12, 15, 24, 69], however, in cases of severe occlusion or complex scenes, challenges often persist.

In this work, we aim to alleviate the problem with better spatial context which is lacking in local 2D correlations. We revisit dense correspondence methods [6, 7, 28, 44, 58, 59], as they demonstrate robustness against matching ambiguity by leveraging rich spatial context. Dense correspondence establishes a corresponding point for every point in an image. To achieve this, these methods often calculate similarities for every pair of points across two images, resulting in a 4D correlation volume [6, 37, 44, 57, 58]. This high-dimensional tensor provides dense bidirectional correspondence, offering matching priors that 2D correlation does not, such as dense matching smoothness from one image to another and vice versa. For example, 4D correlation can provide the constraint that the correspondence of one point to another image is spatially coherent with the correspondences of its neighboring points [46]. However, incorporating the advantages of dense correspondence, which stem from the use of 4D correlation, into point tracking poses significant challenges. Not only does it introduce a

substantial computational burden but the high-dimensionality of the correlation also necessitates a dedicated design for proper processing [6, 35, 46].

We solve the problem by formulating point tracking as a local all-pair correspondence problem, contrary to predominant point-to-region correspondence methods [12, 15, 24, 63], as illustrated in Fig. 2. We construct a local 4D correlation that finds all-pair matches between the local region around a query point and a corresponding local region on the target frame. With this formulation, our framework gains the ability to resolve matching ambiguities, provided by 4D correlation, while maintaining efficiency due to a constrained search range. The local 4D correlation is then processed by a lightweight correlation encoder carefully designed to handle high-dimensional correlation volume. This encoder decomposes the processing into two branches of 2D convolution layers and produces a compact correlation embedding. We then use a Transformer [24] to integrate temporal context into the embeddings. The Transformer’s global receptive field facilitates effective modeling of long-range dependencies despite its compact architecture. Our experiments demonstrate that stack of 3 Transformer layers is sufficient to significantly outperform state-of-the-arts [12, 24]. Additionally, we found that using relative position bias [42, 43, 50] allows the Transformer to process sequences of variable length. This enables our model to handle long videos without the need for a hand-designed chaining process [15, 24].

Our model, dubbed LocoTrack, outperforms the recent state-of-the-art model while maintaining an extremely lightweight architecture, as illustrated in Fig. 1. Specifically, our small model variant achieves a +2.5 AJ increase in the TAP-Vid-DAVIS dataset compared to Cotracker [24] and offers 6 $\times$  faster inference speed. Additionally, it surpasses TAPIR [12] by +5.6 AJ with 3.5 $\times$  faster inference in the same dataset. Our larger variant, while still faster than competing state-of-the-art models [12, 24], demonstrates even further performance gains.

In summary, LocoTrack is a highly efficient and accurate model for point tracking. Its core components include a novel local all-pair correspondence formulation, leveraging dense correspondence to improve robustness against matching ambiguity, a lightweight correlation encoder that ensures computational efficiency, and a Transformer for incorporating temporal information over variable context lengths.

## 2 Related Work

**Point correspondence.** The aim of point correspondence, which is also known as sparse feature matching [10, 13, 31, 48], is to identify corresponding points across images within a set of detected points. This is often achieved by matching a hand-designed descriptors [1, 31] or, more recently, learnable deep features [10, 22, 32, 52]. They are also applicable to videos [40], as the task primarily targets image pairs with large baselines, which is similar to the case with video frames. These approaches filter out noisy correspondences using geometric constraints [16, 49, 56] or their learnable counterparts [22, 48, 68]. However, they often struggle with objects that exhibit deformation [67]. Also, they primarily target

the correspondence of geometrically salient points (*i.e.*, detected points) rather than any arbitrary point.

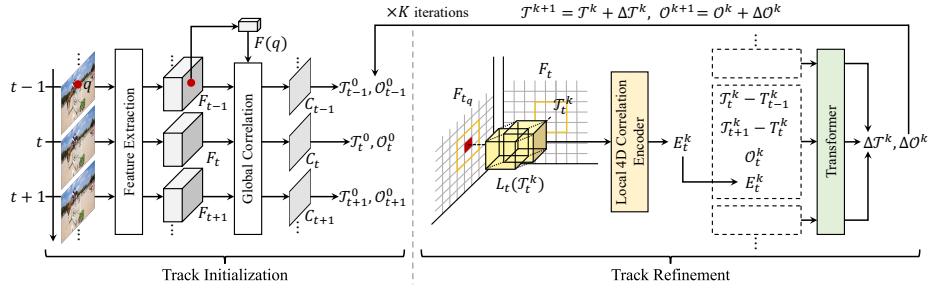
**Long-range point correspondence in video.** Recent methods [2, 11, 12, 15, 24, 63, 69] finds point correspondence in a video, aiming to find a track for a query point over a long sequence of video. They capture a long-range temporal context with MLP-Mixer [2, 15], 1D convolution [12, 69], or Transformer [24]. However, they either leverage a constrained length of sequence within a local temporal window and use sliding window inference to process videos longer than the fixed window size [2, 15, 24], or they necessitate a series of convolution layers to expand the temporal receptive field [12, 69]. Recent Cotracker [24] leverage spatial context by aggregating supporting tracks with self-attention. However, this approach requires tracking additional query points, which introduces significant computational overhead. Notably, Context-PIPs [2] constructs a correlation map across sparse points around the query and the target region. However, this sparsity may limit the model’s ability to fully leverage the matching prior that all-pair correlation can provide, such as matching smoothness.

**Dense correspondence.** Dense correspondence [28] aims to establish pixel-wise correspondence between a pair of images. Conventional methods [6, 18, 33, 37, 45, 54, 58–60] often leverage a 4-dimensional correlation volume, which computes pairwise cosine similarity between localized deep feature descriptors from two images, as the 4D correlation provides a mean for disambiguate the matching process. Traditionally, bidirectional matches from 4D correlation are filtered to remove spurious matches using techniques such as the second nearest neighbor ratio test [31] or the mutual nearest neighbor constraint. Recent methods instead learn patterns within the correlation map to disambiguate matches. DGC-Net [33] and GLU-Net [58] proposed a coarse-to-fine architecture leveraging global 4D correlation followed by local 2D correlation. CATs [6, 7] propose a transformer-based architecture to aggregate the global 4D correlation. GoCor [57], NCNet [45], and RAFT [54] developed an efficient framework using local 4D correlation to learn spatial priors in both image pairs, addressing matching ambiguities.

The use of 4D correlation extends beyond dense correspondence. It has been widely applied in fields such as video object segmentation [5, 39], few-shot semantic segmentation [19, 35], and few-shot classification [23]. However, its application in point tracking remains underexplored. Instead, several attempts have been made to integrate the strengths of off-the-shelf dense correspondence model [54] into point tracking. These include chaining dense correspondences [15, 54], which has limitations in recovering from occlusion, or directly finding correspondences with distant frames [36, 38, 64], which is computationally expensive.

### 3 Method

In this work, we integrate the effectiveness of a 4D correlation volume into our point tracking pipeline. Compared to the widely used 2D correlation [11, 12, 15,



**Fig. 3: Overall architecture of LocoTrack.** Our model comprises two stages: track initialization and track refinement. The track initialization stage determines a rough position by conducting feature matching with global correlation. The track refinement stage iteratively refines the track by processing the local 4D correlation.

[24], 4D correlation offers two distinct characteristics that provide valuable information for filtering out noisy correspondences, leading to more robust tracking:

- **Bidirectional correspondence:** 4D correlation provides bidirectional correspondences, which can be used to verify matches and reduce ambiguity [31]. This prior is often leveraged by checking for mutual consensus [46] or by employing a ratio test [31].
- **Smooth matching:** A 4D correlation volume is constructed using dense all-pair correlations, which can be leveraged to enforce matching smoothness and improve matching consistency across neighboring points [46, 57, 58].

We aim to leverage these benefits of the 4D correlation volume while maintaining efficient computation. We achieve this by restricting the search space to a local neighborhood when constructing the 4D correlation volume. Along with the use of local 4D correlation, we also propose a recipe to benefit from the global receptive field of Transformers for long-range temporal modeling. This enables our model to capture long-range context within a few (even only with 3) stacks of transformer layers, resulting in a compact architecture.

Our method, dubbed LocoTrack, takes as input a query point  $q = (x_q, y_q, t_q) \in \mathbb{R}^3$  and a video  $\mathcal{V} = \{\mathcal{I}_t\}_{t=0}^{t=T-1}$ , where  $T$  indicates the number of frames and  $\mathcal{I}_t \in \mathbb{R}^{H \times W \times 3}$  represents the  $t$ -th frame. We assume query point can be given in the arbitrary time step. Our goal is to produce a track  $\mathcal{T} = \{\mathcal{T}_t\}_{t=0}^{t=T-1}$ , where  $\mathcal{T}_t \in \mathbb{R}^2$ , and associated occlusion probabilities  $\mathcal{O} = \{\mathcal{O}_t\}_{t=0}^{t=T-1}$ , where  $\mathcal{O}_t \in [0, 1]$ . Following previous works [12, 24], our method predicts the track in two stage approach: an initialization stage followed by a refinement stage, each detailed in the follows, as illustrated in Fig. 3.

### 3.1 Stage I: Track Initialization

To estimate the initial track of a given query point, we conduct feature matching that constructs a global similarity map between features derived from the query point and the target frame’s feature map, and choose the positions with the



**Fig. 4: Visualization of correspondence.** We visualize the correspondences established between the query and target regions. Our refined 4D correlation (e) demonstrates a clear reduction in matching ambiguity and yields better correspondences compared to the noisy results produced by 2D correlation (d). This improvement aligns closely with the ground truth (c).

highest scores as the initial track. This similarity map, often referred to as a correlation map, provides a strong signal for accurately initializing the track’s positions. We use a global correlation map for the initialization stage, which calculates the similarity for every pixel in each frame.

Specifically, we use hierarchical feature maps derived from the feature backbone [17]. Given a set of pyramidal feature maps  $\{F_t^l\}_{t=1}^{T-1} = \mathcal{E}(\mathcal{V})$ , where  $\mathcal{E}(\cdot)$  represents the feature extractor and  $F_t^l$  indicates a level  $l \in \{0, \dots, L-1\}$  feature map in frame  $t$ , we sample a query feature vector  $F^l(q)$  at position  $q$  from  $F^l$  using linear interpolation for each level  $l$ . The global correlation map is calculated as  $C_t^l = \frac{F_t^l \cdot F^l(q)}{\|F_t^l\|_2 \|F^l(q)\|_2} \in \mathbb{R}^{H^l \times W^l}$ , where  $H^l$  and  $W^l$  denote the height and width of the feature map at the  $l$ -th level, respectively. The correlation maps obtained from multiple levels are resized to the largest feature map size and concatenated as  $C_t \in \mathbb{R}^{H^0 \times W^0 \times L}$ . The concatenated maps are processed as follows to generate the initial track and occlusion probabilities:

$$\begin{aligned} \mathcal{T}_t^0 &= \text{Softargmax}(\text{Conv2D}(C_t); \tau), \\ \mathcal{O}_t^0 &= \text{Linear}([\text{Maxpool}(C_t); \text{Avgpool}(C_t)]), \end{aligned} \quad (1)$$

where  $\text{Conv2D} : \mathbb{R}^{H \times W \times L} \rightarrow \mathbb{R}^{H \times W}$  is a single-layered 2D convolution layer,  $\text{Softargmax} : \mathbb{R}^{H \times W} \rightarrow \mathbb{R}^2$  is a differentiable argmax function with a Gaussian kernel [27] that provides the 2D position of the maximum value,  $\tau$  is a temperature parameter,  $[ \cdot ]$  indicates concatenation, and  $\text{Linear} : \mathbb{R}^{2L} \rightarrow \mathbb{R}$  is a linear projection. Similar to CBAM [65], we apply global max and average pooling followed by a linear projection to calculate initial occlusion probabilities.

### 3.2 Stage II: Track Refinement

We found that the initial track  $\mathcal{T}^0$  and  $\mathcal{O}^0$  often exhibit severe jittering, arising from the matching ambiguity from the noisy correlation map. We iteratively refine the noise in the initial tracks  $\mathcal{T}^0$  and  $\mathcal{O}^0$ . For each iteration, we estimate the residuals  $\Delta\mathcal{T}^k$  and  $\Delta\mathcal{O}^k$ , which are then applied to the tracks as  $\mathcal{T}^{k+1} := \mathcal{T}^k + \Delta\mathcal{T}^k$  and  $\mathcal{O}^{k+1} := \mathcal{O}^k + \Delta\mathcal{O}^k$ . During the refining process, the matching noise can be rectified in two ways: 1) by establishing locally dense correspondences with local 4D correlation, and 2) through temporal modeling with a Transformer [62].

**Local 4D correlation.** The 2D correlation  $\mathcal{C}_t$  often exhibits limitations when dealing with repetitive patterns or homogeneous regions as exemplified in Fig. 4. Inspired by dense correspondence literatures, we utilize 4D correlation to provide richer information for refining tracks compared to 2D correlation. The 4D correlation  $\mathbf{C}^{4D} \in \mathbb{R}^{H \times W \times H \times W}$ , which computes every pairwise similarity, can be formally defined as follows:

$$\mathbf{C}_t^{4D}(i, j) = \frac{F_t(i) \cdot F_{t_q}(j)}{\|F_t(i)\|_2 \|F_{t_q}(j)\|_2}, \quad (2)$$

where  $F_{t_q}$  is the feature map from the frame in which the query point is located, and  $i$  and  $j$  specify the locations within the feature map. However, since a global 4D correlation volume with the shape of  $H \times W \times H \times W$  becomes computationally intractable, we employ a local 4D correlation  $\mathbf{L} \in \mathbb{R}^{h_p \times w_p \times h_q \times w_q}$ , where  $(h_p, w_p, h_q, w_q)$  denotes spatial resolution of local correlation. We define the correlation as follows:

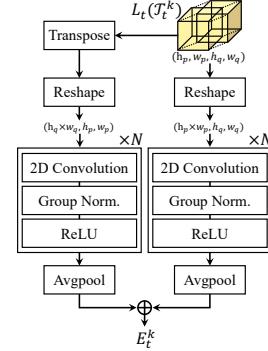
$$\begin{aligned} \mathcal{N}(p, r) &= \{p + \delta \mid \delta \in \mathbb{Z}^2, \|\delta\|_\infty \leq r\}, \\ \mathbf{L}_t(i, j; p) &= \frac{F_t(i) \cdot F_{t_q}(j)}{\|F_t(i)\|_2 \|F_{t_q}(j)\|_2}, \quad i \in \mathcal{N}(p; r_p), \quad j \in \mathcal{N}(q; r_q), \end{aligned} \quad (3)$$

where  $r_p$  and  $r_q$  are the radii of the regions around points  $p$  and  $q$ , respectively, resulting in  $h_p = w_p = 2r_p + 1$  and  $h_q = w_q = 2r_q + 1$ . The correlation then serves as a cue for refining the track  $\mathcal{T}^k$ . To achieve this, we calculate the set of local correlations around the intermediate predicted position, denoted as  $\{\mathbf{L}_t(\mathcal{T}_t^k)\}_{t=0}^{T-1}$  with abuse of notation.

**Local 4D correlation encoder.** We then process the local 4D correlation volume to disambiguate matching ambiguities, leveraging the smoothness of both the query and target dimensions of correlations. Note that the obtained 4D correlation is a high-dimensional tensor, posing an additional challenge for its correct processing. In this regard, we introduce an efficient encoding strategy that decomposes the processing of the correlation. We process the 4D correlation in two symmetrical branches as shown in Fig. 5. One branch spatially processes the dimensions of the query, treating the flattened target dimensions as a channel dimension. The other branch, on the other hand, considers the query dimensions as channel. Each branch compresses the correlation into a single vector, which are then concatenated to form a correlation embedding  $E_t^k$ :

$$E_t^k = [\mathcal{E}_L(\mathbf{L}_t(\mathcal{T}_t^k)); \mathcal{E}_L((\mathbf{L}_t(\mathcal{T}_t^k))^T)], \quad (4)$$

where  $\mathbf{L}(i, j) = \mathbf{L}^T(j, i)$ . The convolutional encoder  $\mathcal{E}_L : \mathbb{R}^{h_p \times w_p \times h_q \times w_q} \rightarrow \mathbb{R}^{C_E}$  consists of stacks of strided 2D convolutions, group normalization [66], and ReLU



**Fig. 5: Local 4D correlation encoder.**

activations. These operations progressively reduce the correlation’s spatial dimensions, followed by a final average pooling layer for a compact representation. We obtain the correlation embedding for all feature levels  $l$ , and concatenate them to form the final embedding. For more details on the local 4D correlation encoder, please refer to the supplementary material.

**Temporal modeling with length-generalizable transformer.** The encoded correlation is then provided to the refinement model. The model refines the initial trajectory and predicts its error with respect to the ground truth,  $\Delta\mathcal{T}$  and  $\Delta\mathcal{O}$ , which requires an ability to leverage temporal context. For the temporal modelling, we explore three candidates widely used in the literature: 1D Convolution [12, 69], MLP-Mixer [15], and Transformer [24]. We consider two aspects to select the appropriate architecture: 1) Can the architecture handle arbitrary sequence lengths  $T$  at test time? 2) Can the temporal receptive field, crucial for capturing long-range context, be sufficiently large with just a few layers stacked? Based on these criteria, we choose the Transformer as our architecture because it can handle arbitrary sequence lengths, a capability the MLP-Mixer lacks. This lack would necessitate an additional test-time strategy (*e.g.*, sliding window inference [15]) to accommodate sequences longer than those used during training. Additionally, the Transformer can form a global receptive field with a single layer, unlike convolution, which requires multiple layers to achieve an expanded receptive field.

Although the Transformer can process sequences of arbitrary length at test time, we found that sinusoidal position encoding [62] degrades performance for videos with sequence lengths that differ from those used during training. Instead, we use relative position bias [42, 43, 50], which disproportionately reduces the impact of distant tokens by adjusting the bias within the Transformer’s attention map. However, relative position bias is based solely on the distance between tokens cannot distinguish their relative direction (*e.g.*, whether token A is before or after token B), which makes it only suitable for causal attention. To address this, we divide the attention head into two groups: one group encodes relative position only for tokens on the left, and the other for tokens on the right:

$$\text{Softmax}(\mathbf{q} \cdot \mathbf{k}^T + b(h)), \text{ where} \\ b(t_1, t_2; h) = \begin{cases} b_{\text{left}}(t_1, t_2; h), & h < \lfloor \frac{N_h}{2} \rfloor, \\ b_{\text{right}}(t_1, t_2; h - \lfloor N_h/2 \rfloor), & h \geq \lfloor \frac{N_h}{2} \rfloor, \end{cases} \quad (5)$$

where  $\mathbf{q}$  and  $\mathbf{k}$  denote the query and key, respectively,  $N_h$  is the number of heads, and  $h \in \{0, \dots, N_h - 1\}$  is the index of the attention head. The bias term  $b_{\text{left}}$  adjusts the attention map to ensure that each query token attends only to key tokens located to its left or within the same position, as follows:

$$b_{\text{left}}(t_1, t_2; h) = \begin{cases} -\infty, & \text{if } t_1 < t_2, \\ -s_h|t_1 - t_2|, & \text{if } t_1 \geq t_2, \end{cases} \quad (6)$$

where  $s_h \in \mathbb{R}^+$  is a scaling factor that controls the rate of bias decay as distance increases. We employ different scaling factors for each attention head, follow-

ing Press et al. [42]. The function  $b_{\text{right}}(\cdot)$  can be similarly defined. With this design choice, we found that the Transformer can generalize to videos of arbitrary length, eliminating the need for test-time hand-designed techniques such as sliding window inference [15, 24].

**Iterative update.** We stack  $N_S$  Transformer layers with the modified self-attention and feed the correlation embedding  $\{E_t^k\}_{t=0}^{t=T-1}$ , the encoded initialized track  $\mathcal{T}^k$ , and occlusion status  $\mathcal{O}^k$  to the Transformer  $\mathcal{E}_S$  to predict track updates. We found using position differences between adjacent frames improves training convergence compared to using the absolute positions. This is formally defined as:

$$\begin{aligned}\Delta \mathcal{T}^k, \Delta \mathcal{O}^k &= \mathcal{E}_S \left( \{ [\sigma(\mathcal{T}_t^k - \mathcal{T}_{t-1}^k); \sigma(\mathcal{T}_{t+1}^k - \mathcal{T}_t^k); \mathcal{O}_t^k; E_t^k] \}_{t=0}^{t=T-1} \right), \\ \mathcal{T}_{-1}^k &:= \mathcal{T}_0^k, \quad \mathcal{T}_T^k := \mathcal{T}_{T-1}^k,\end{aligned}\quad (7)$$

where  $\sigma(\cdot)$  is a sinusoidal encoding [53],  $[\cdot]$  denotes concatenation, and  $\Delta \mathcal{T}^k$  and  $\Delta \mathcal{O}^k$  are predicted updates. Sequentially, the predicted updates are applied to initial track as  $\mathcal{T}^{k+1} := \mathcal{T}^k + \Delta \mathcal{T}^k$  and  $\mathcal{O}^{k+1} := \mathcal{O}^k + \Delta \mathcal{O}^k$ . We perform  $K$  iterations, yielding the final refined track  $\mathcal{T}^K$  and  $\mathcal{O}^K$ .

## 4 Experiments

### 4.1 Implementation Details

We use JAX [3] for implementation. For training, we utilize the Panning MOVi-E dataset [12] generated with Kubric [14]. We employ the loss functions introduced in Doersch et al. [12], including the prediction of additional uncertainty estimation for both track initialization and a refinement model. We use the AdamW [30] optimizer and use  $1 \cdot 10^{-3}$  for both learning rate and weight decay. We employ a cosine learning rate scheduler with a 1000-step warmup stage [29]. Following Sun et al. [51], we apply gradient clipping with a value of 1.0. The initialization stage is first trained for 100K steps, followed by track refinement model training for an additional 300K steps. This process takes approximately 4 days on 8 NVIDIA RTX 3090 GPUs with a batch size of 1 per GPU. For each batch, we randomly sample 256 tracks. We use a  $256 \times 256$  training resolution, following the standard protocol of TAP-Vid benchmark.

Our feature backbone is ResNet18 [17] with instance normalization [61] replacing batch normalization [20]. We use three pyramidal feature maps ( $L = 3$ ) from ResNet, each with a stride of 2, 4 and 8, respectively. The temperature value for softargmax is set to  $\sigma = 20.0$ . The radii of the local correlation window are  $r_q = r_p = 3$ . We stack  $N_S = 3$  Transformer layers for  $\mathcal{E}_S$ . The number of iterations ( $K$ ) is set to 4. For the track refinement model, we propose two variants: a small model and a base model. All ablations are conducted using the base model. The hidden dimension of the Transformer is set to 256 for the small model and 384 for the base model. The number of heads is set to 4 for the small model and 6 for the base model. For more details, please refer to supplementary materials.

**Table 1: Quantitative comparison on the TAP-Vid datasets with the strided query mode.** Throughput is measured on a single Nvidia RTX 3090 GPU.

Method	Kinetics			DAVIS			RGB-Stacking			Throughput (points/sec)
	AJ	$< \delta_{avg}^x$	OA	AJ	$< \delta_{avg}^x$	OA	AJ	$< \delta_{avg}^x$	OA	
<i>Input Resolution 256×256</i>										
Kubric-VFS-Like [14]	40.5	59.0	80.0	33.1	48.5	79.4	57.9	72.6	91.9	-
TAP-Net [11]	46.6	60.9	85.0	38.4	53.1	82.3	59.9	72.8	90.4	<b>29,535.98</b>
RAFT [54]	34.5	52.5	79.7	30.0	46.3	79.6	44.0	58.6	90.4	23,405.71
TAPIR [12]	57.2	70.1	87.8	61.3	73.6	88.8	62.7	74.6	91.6	2,097.32
LocoTrack-S	<b>59.6</b>	<u>72.7</u>	<u>88.1</u>	<u>66.9</u>	<u>78.8</u>	<u>88.9</u>	<b>77.4</b>	<b>87.0</b>	<u>92.9</u>	7,244.47
LocoTrack-B	59.5	<b>73.0</b>	<b>88.5</b>	<b>67.8</b>	<b>79.6</b>	<b>89.9</b>	77.1	<u>86.9</u>	<b>93.2</b>	4,358.96
<i>Input Resolution 384×512</i>										
PIPs [15]	35.3	54.8	77.4	42.0	59.4	82.1	37.3	51.0	<b>91.6</b>	46.43
FlowTrack [8]	-	-	-	66.0	79.8	87.2	-	-	-	-
CoTracker [24]	-	-	-	65.9	79.4	<b>89.9</b>	-	-	-	1,146.79
LocoTrack-S	<u>58.7</u>	<u>72.2</u>	<u>84.5</u>	<u>68.4</u>	<u>80.4</u>	87.5	<b>71.0</b>	<b>84.4</b>	83.3	<b>6,820.57</b>
LocoTrack-B	<b>59.1</b>	<b>72.5</b>	<b>85.7</b>	<b>69.4</b>	<b>81.3</b>	<u>88.6</u>	70.8	<u>83.2</u>	84.1	4,196.36

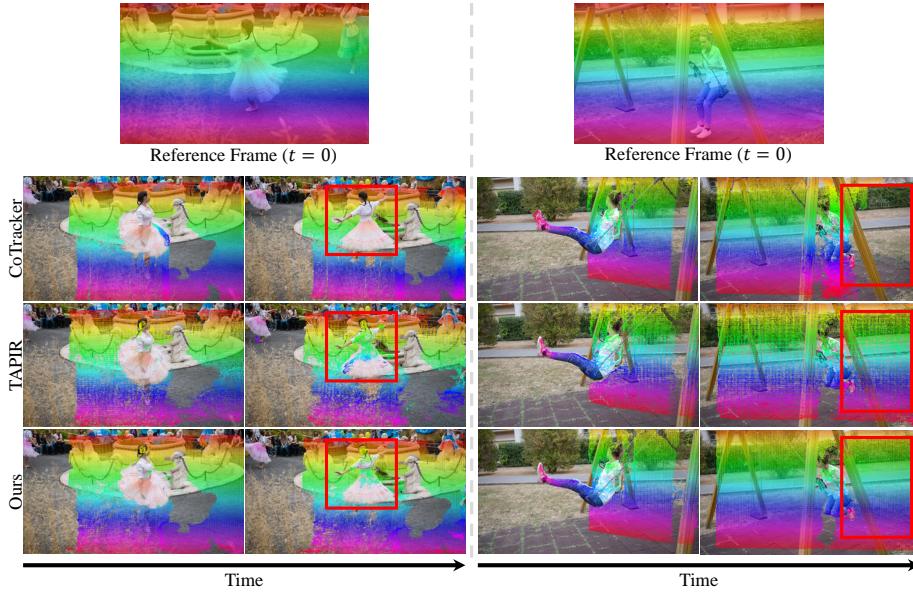
**Table 2: Quantitative comparison on the query first mode.**

Method	Kinetics-First			DAVIS-First			RoboTAP-First		
	AJ	$< \delta_{avg}^x$	OA	AJ	$< \delta_{avg}^x$	OA	AJ	$< \delta_{avg}^x$	OA
<i>Input Resolution 256×256</i>									
TAP-Net [11]	38.5	54.4	80.6	33.0	48.6	78.8	45.1	62.1	82.9
TAPIR [12]	49.6	64.2	<u>85.0</u>	56.2	70.0	<u>86.5</u>	59.6	73.4	<u>87.0</u>
LocoTrack-S	<b>52.8</b>	<u>66.5</u>	84.9	<u>62.0</u>	<u>74.3</u>	86.1	<b>62.5</b>	<u>76.0</u>	<u>87.0</u>
LocoTrack-B	<b>52.9</b>	<b>66.8</b>	<b>85.3</b>	<b>63.0</b>	<b>75.3</b>	<b>87.2</b>	<u>62.3</u>	<b>76.2</b>	<b>87.1</b>
<i>Input Resolution 384×512</i>									
CoTracker [24]	48.7	64.3	<b>86.5</b>	60.6	75.4	<b>89.3</b>	-	-	-
LocoTrack-S	<u>51.9</u>	<u>66.1</u>	81.2	<u>63.2</u>	<u>76.2</u>	84.6	-	-	-
LocoTrack-B	<b>52.3</b>	<b>66.4</b>	<u>82.1</u>	<u>64.8</u>	<b>77.4</b>	<u>86.2</u>	-	-	-

## 4.2 Evaluation Protocol

We evaluate the precision of the predicted tracks using the TAP-Vid benchmark [11] and the RoboTAP dataset [63]. For evaluation metrics, we use position accuracy ( $< \delta_{avg}^x$ ), occlusion accuracy (OA), and average Jaccard (AJ).  $< \delta_{avg}^x$  calculates position accuracy for the points visible in ground-truth. It calculates the percentage of correct points (PCK) [46], averaged over the error threshold values of 1, 2, 4, 8, and 16 pixels. OA represents the average accuracy of the binary classification results for occlusion. AJ is a metric that evaluates both position accuracy and occlusion accuracy.

Following Doersch et al. [11], we evaluate the datasets in two modes: strided query mode and first query mode. Strided query mode samples the query point along the ground-truth track at fixed intervals, sampling every 5 frames, whereas first query mode samples the query point solely from the first visible point.



**Fig. 6: Qualitative comparison of long-range tracking.** We visualize dense tracking results generated by LocoTrack and state-of-the-art methods [12, 24]. These visualizations use query points densely distributed within the initial reference frame. Our model can establish highly precise correspondences over long ranges, even in the presence of occlusions and matching challenges like homogeneous areas or deforming objects. Best viewed in color.

### 4.3 Main Results

**Quantitative comparison.** We compare our method with recent state-of-the-art approaches [8, 11, 12, 14, 15, 24, 54] in both strided query mode, with scores shown in Table 1, and first query mode, with scores shown in Table 2. To ensure a fair comparison, we categorize models based on their input resolution sizes:  $256 \times 256$  and  $384 \times 512$ . Along with performance, we also present the throughput of each model, which indicates the number of points a model can process within a second. Higher throughput implies more efficient computation.

Our small variant, LocoTrack-S, already achieves state-of-the-art performance on AJ and position accuracy across all benchmarks, surpassing both TAPIR and CoTracker by a large margin. In the DAVIS benchmark with strided query mode, we achieved a +5.6 AJ improvement compared to TAPIR and a +2.5 AJ improvement compared to CoTracker. This small variant model is not only powerful but also extremely efficient compared to recent state-of-the-art methods. Our model demonstrates  $3.5 \times$  higher throughput than TAPIR and  $6 \times$  higher than CoTracker. LocoTrack-B model shows even better performance, achieving a +0.9 AJ improvement over our small variant in DAVIS strided query mode.

However, our model often shows degradation on some datasets in  $384 \times 512$ . We attribute this degradation to the diminished effective receptive field of local correlation when resolution is increased.

**Table 3: Comparison of computation cost.** We measure the inference time with a varying number of query points and calculate the FLOPs for the feature backbone and refinement stage, along with the number of parameters. All metrics are measured using a video consisting of 24 frames on a single Nvidia RTX 3090 GPU.

Method	Inference Time (s)						Throughput (points/sec)	Backbone FLOPs (G)	FLOPs per point (G)	# of Params. (M)
	$10^0$ point	$10^1$ points	$10^2$ points	$10^3$ points	$10^4$ points	$10^5$ points				
RAFT [54]	-	-	-	-	-	-	23,405.71	325.45	-	5.3
CoTracker [24]	0.53	0.53	0.53	1.18	8.40	87.2	1,146.79	<b>624.83</b>	4.65	45.5
TAPIR [12]	<b>0.06</b>	<b>0.06</b>	0.19	0.82	4.89	47.68	2,097.32	<b>442.16</b>	5.12	29.3
LocoTrack-S	<b>0.04</b>	<b>0.04</b>	<b>0.05</b>	<b>0.17</b>	<b>1.44</b>	<b>14.23</b>	<b>7,244.47</b>	<b>442.16</b>	<b>1.08</b>	<b>8.2</b>
LocoTrack-B	<b>0.04</b>	<b>0.04</b>	<b>0.06</b>	<b>0.26</b>	<b>2.39</b>	<b>23.37</b>	4,358.96	<b>442.16</b>	<b>2.10</b>	<b>11.5</b>

**Table 4: Ablation on construction of correlation volume.**

	Local Corr. Size	Query Neighbour	DAVIS		
			AJ	$<\delta_{avg}^2$	OA
(I)	$7 \times 7$	No neighbour (2D corr.)	65.0	77.2	89.0
(II)	$9 \times 7 \times 7$	Uniform random in local region	65.7	77.8	88.9
(III)	$1 \times 9 \times 7 \times 7$	Horizontal line	66.5	78.4	89.4
(IV)	$3 \times 3 \times 7 \times 7$	Regular grid ( $r_q = 1$ )	67.2	79.1	89.5
(V)	$7 \times 7 \times 7 \times 7$	Regular grid ( $r_q = 3$ , Ours)	<b>67.8</b>	<b>79.6</b>	<b>89.9</b>

**Qualitative comparison.** The qualitative comparison is shown in Fig. 6. We visualize the results from the DAVIS [41] dataset, with the input resized to  $384 \times 512$  resolution. Note that images at their original resolution are used for visualization. Overall, our method demonstrates superior smoothness compared to TAPIR. Our predictions are spatially coherent, even over long-range tracking sequences with occlusion.

#### 4.4 Analysis and Ablation Study

**Efficiency comparison.** We compare efficiency to recent state of the arts [12, 24, 54] in Table 3. We measure inference time, throughput, FLOPs, and the number of parameters for a 24 frame video. We report inference time for a varying number of query points, increasing exponentially from  $10^0$  to  $10^5$ . To measure throughput, we calculate the average time required to add each query point. Also, we measure FLOPs for both the feature backbone and the refinement model, focusing on the incremental FLOPs per additional point.

All variants of our model demonstrate superior efficiency across all metrics. Our small variant exhibits  $4.7 \times$  lower FLOPs per point compared to TAPIR and  $4.3 \times$  lower than CoTracker. Additionally, our model boasts a compact parameter count of only 8.2M, which is  $5.5 \times$  lower than CoTracker. Remarkably, our model can process  $10^4$  points in approximately one second, implying real-time processing of  $64 \times 64$  near-dense query points for a 24 frame rate video. This underscores the practicality of our model, paving the way for real-time applications.

**Analysis on local correlation.** In Table 4, we analyze the construction of our local correlation method, focusing on how we sample neighboring points around the query points rather than target points. (I) represents the performance of local 2D correlation, a common approach in the literature [12, 15, 24]. The performance gap between (I) and (VI) demonstrates the superiority of our

**Table 5: Ablation on position encoding.**

Method	DAVIS		
	AJ	$\delta_{avg}^x$	OA
(I) Sinusoidal encoding [62]	61.9	73.9	83.5
(II) Relative position bias (Ours)	<b>67.8</b>	<b>79.6</b>	<b>89.9</b>

**Table 6: Ablation on architecture of  $\mathcal{E}_S$ .** We found that our model outperforms its counterpart while using the same number of parameters.

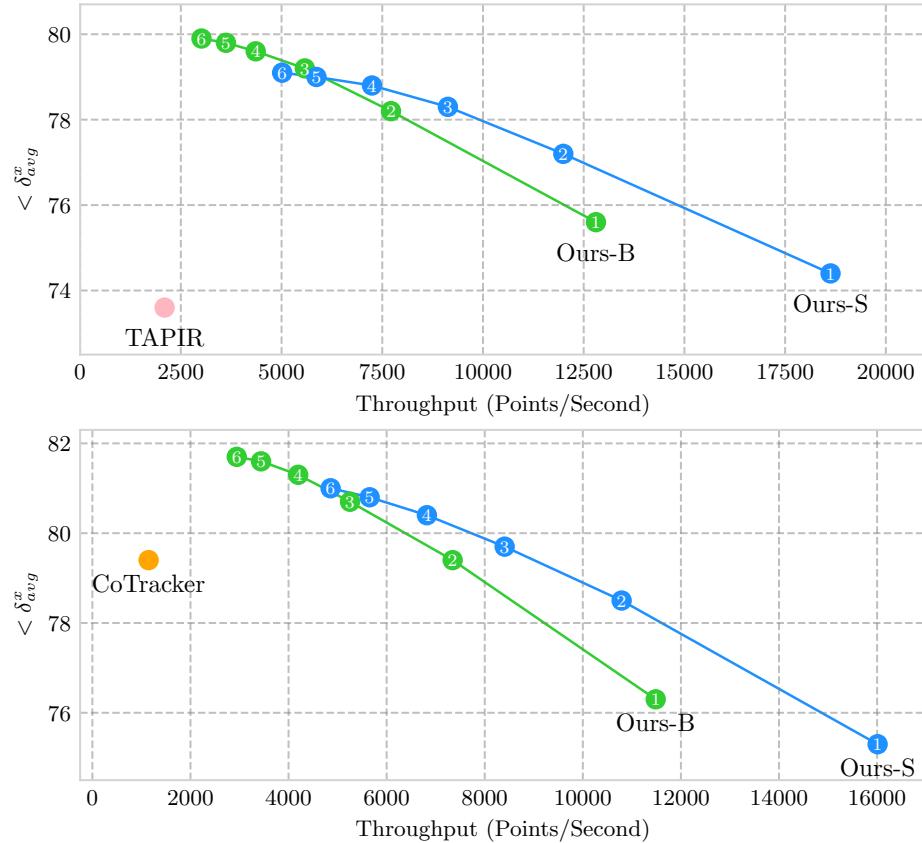
Method	# of Layers	# of Params.	DAVIS		
			AJ	$\delta_{avg}^x$	OA
(I) 1D Conv Mixer (TAPIR) [12]	3	11.5	66.1	78.0	87.5
(I) LocoTrack-B (Ours)	3	11.5	<b>67.8</b>	<b>79.6</b>	<b>89.9</b>

4D correlation approach over 2D. (II) and (III) investigate the importance of calculating dense all-pair correlations within the local region. In (II), we use randomly sampled positions for the query point’s neighbors, while (III) uses a horizontal line-shaped neighborhood. Their inferior performance compared to (IV), which samples the same number of points densely, emphasizes the value of our all-pair local 4D correlation. (IV) and (V) examine the effect of local region size. The gap between (IV) and (V), supports our choice of region size. (V) represents our final model.

**Ablation on position encoding of Transformer.** In Table 5, we ablate the effect of relative position bias. With sinusoidal encoding [62], we observe significant performance degradation during inference (I) with variable length. In contrast, relative position bias demonstrates generalization to unseen sequence lengths at inference time (II). This approach eliminates the need for hand-designed chaining processes (*i.e.*, sliding window inference [15, 24]) where window overlapping leads to computational inefficiency.

**Ablation on the architecture of refinement model.** We verify the advantages of using a Transformer architecture over a Convolution-based architecture in Table 6. Our comparison includes the architecture proposed in Doersch et al. [12], which replaces the token mixing layer of MLP-Mixer [55] with depth-wise 1D convolution. We ensure a fair comparison by matching the number of parameters and layers between the models. Our Transformer-based model achieves superior performance. We believe this difference stems from their receptive fields: Transformers can achieve a global receptive field within a single layer, while convolutions require multiple stacked layers. Although convolutions can also achieve large receptive fields with lightweight designs [4, 9], their exploration in long-range point tracking remains a promising area for future work.

**Analysis on the number of iterations.** We show the performance and throughput of our model, varying the number of iterations, in Fig. 7. We compare our model with TAPIR and CoTracker at their respective resolutions. Surprisingly, our model surpasses TAPIR even with a single iteration for both the small and base variants. With a single iteration, our small variant is about 9× faster than TAPIR. Compared to CoTracker, our model is about 8× faster at the same performance level.



**Fig. 7: Results with a varying number of refinement iterations on TAP-Vid-DAVIS.** The number in the circle denotes the number of iterations. (up) In a  $256 \times 256$  resolution, compared to TAPIR [12], LocoTrack achieves better performance in a *single iteration* while being about  $9\times$  faster. (below) In a  $384 \times 512$  resolution, compared to CoTracker [24], LocoTrack achieves comparable performance while being about  $8\times$  faster.

## 5 Conclusion

We introduce LocoTrack, an approach to the point tracking task, addressing the shortcomings of existing methods that rely solely on local 2D correlation. Our core innovation lies in a local all-pair correspondence formulation, combining the rich spatial context of 4D correlation with computational efficiency by limiting the search range. Further, a length-generalizable Transformer empowers the model to handle videos of varying lengths, eliminating the need for hand-designed processes. Our approach demonstrates superior performance and real-time inference while requiring significantly less computation compared to state-of-the-art methods.

## Acknowledgements

This research was supported by the MSIT, Korea (IITP-2024-2020-0-01819, RS-2023-00227592), Culture, Sports, and Tourism R&D Program through the Korea Creative Content Agency grant funded by the Ministry of Culture, Sports and Tourism (Research on neural watermark technology for copyright protection of generative AI 3D content, RS-2024-00348469, RS-2024-00333068) and National Research Foundation of Korea (RS-2024-00346597).

# Local All-Pair Correspondence for Point Tracking

## –Supplementary Material–

### A More Implementation Details

For generating the Panning-MOVIE dataset [12], we randomly add 10-20 static objects and 5-10 dynamic objects to each scene. The dataset comprises 10,000 videos, including a validation set of 250. For the sinusoidal position encoding function [53]  $\sigma(\cdot)$ , we use a channel size of 20 along with the original unnormalized coordinate. This results in a total of 21 channels. For all qualitative comparisons, we use LocoTrack-B model with a resolution of  $384 \times 512$ .

**Details of the evaluation benchmark.** We evaluate the precision of the predicted tracks using the TAP-Vid benchmark [11]. This benchmark comprises both real-world video datasets and synthetic video datasets. **TAP-Vid-Kinetics** includes 1,189 real-world videos from the Kinetics [25] dataset. As the videos are collected from YouTube, they often contain edits such as scene cuts, text, fade-ins or -outs, or captions. **TAP-Vid-DAVIS** comprises real-world videos from the DAVIS [41] dataset. This dataset includes 30 videos featuring various concepts of objects with deformations. **TAP-Vid-RGB-Stacking** consists of 50 synthetic videos [26]. These videos feature a robot arm stacking geometric shapes against a monotonic background, with the camera remaining static. In addition to the TAP-Vid benchmark, we also evaluate our model on the **RoboTAP** dataset [63], which comprises 265 real-world videos of robot arm manipulation.

**Table 7: Convolutional layer configurations for different model sizes.**

Model	Channel Sizes	Kernel Size	Strides
Small	(64, 128)	(5, 2)	(4, 2)
Base	(64, 128, 128)	(3, 3, 2)	(2, 2, 2)

**Detailed architecture of local 4D correlation encoder.** We stack blocks of convolutional layers, where each block consists of a 2D convolution, group normalization [66], and ReLU activation. See Table 7 for details. For the small model, we use an intermediate channel size of (64, 128) for each block. For the base model, the intermediate channel sizes are (64, 128, 128) for each block. For every instance of group normalization, we set the group size to 16.

**Details of correlation visualization.** For the correlation visualization in Fig. 3 of the main text, we train a linear layer to project the correlation embedding  $E_t^k$  into a local 2D correlation with a shape of  $7 \times 7$ . This local 2D correlation then undergoes a softargmax operation to predict the error relative to the ground truth. We begin with the pre-trained model and train the linear layer for 20,000 iterations. For clarity, we bilinearly upsample the  $7 \times 7$  correlation to  $256 \times 256$ .

## B More Qualitative Comparison

We provide more qualitative comparisons to recent state-of-the-art methods [12, 24] in Fig. 8 and Fig. 9. Our model establishes accurate correspondences in homogeneous areas and on deforming objects, demonstrating robust occlusion handling even under severe occlusion conditions.



**Fig. 8:** Additional qualitative comparison with state-of-the-art [12, 24].

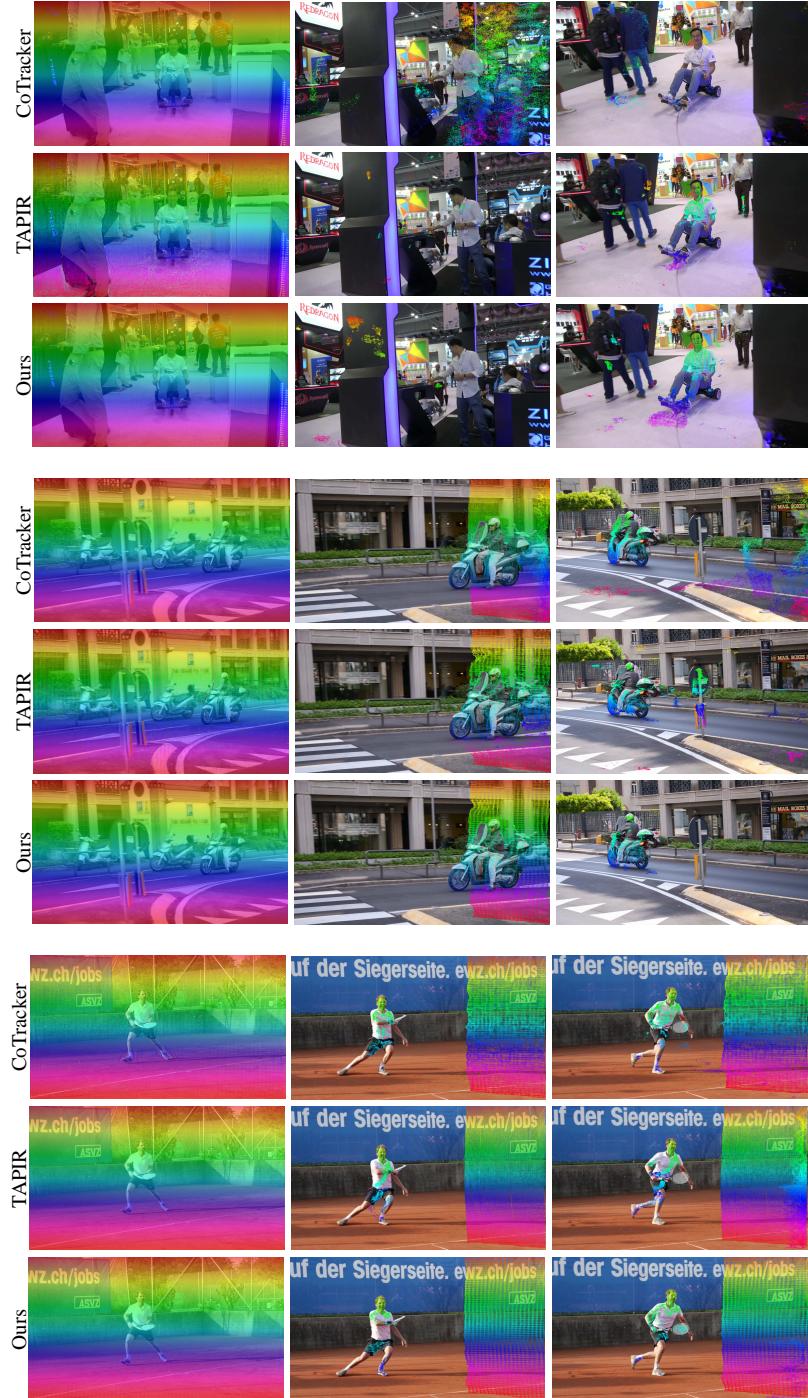


Fig. 9: Additional qualitative comparison with state-of-the-art [12, 24].

## References

1. Bay, H., Tuytelaars, T., Van Gool, L.: Surf: Speeded up robust features. In: Computer Vision–ECCV 2006: 9th European Conference on Computer Vision, Graz, Austria, May 7–13, 2006. Proceedings, Part I 9. pp. 404–417. Springer (2006)
2. Bian, W., Huang, Z., Shi, X., Dong, Y., Li, Y., Li, H.: Context-tap: Tracking any point demands spatial context features. arXiv preprint arXiv:2306.02000 (2023)
3. Bradbury, J., Frostig, R., Hawkins, P., Johnson, M.J., Leary, C., Maclaurin, D., Necula, G., Paszke, A., VanderPlas, J., Wanderman-Milne, S., Zhang, Q.: JAX: composable transformations of Python+NumPy programs (2018), <http://github.com/google/jax>
4. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation. arXiv preprint arXiv:1706.05587 (2017)
5. Cheng, H.K., Schwing, A.G.: Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model. In: European Conference on Computer Vision. pp. 640–658. Springer (2022)
6. Cho, S., Hong, S., Jeon, S., Lee, Y., Sohn, K., Kim, S.: Cats: Cost aggregation transformers for visual correspondence. Advances in Neural Information Processing Systems **34**, 9011–9023 (2021)
7. Cho, S., Hong, S., Kim, S.: Cats++: Boosting cost aggregation with convolutions and transformers. IEEE Transactions on Pattern Analysis and Machine Intelligence **45**(6), 7174–7194 (2022)
8. Cho, S., Huang, J., Kim, S., Lee, J.Y.: Flowtrack: Revisiting optical flow for long-range dense tracking. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 19268–19277 (2024)
9. Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y.: Deformable convolutional networks. In: Proceedings of the IEEE international conference on computer vision. pp. 764–773 (2017)
10. DeTone, D., Malisiewicz, T., Rabinovich, A.: Superpoint: Self-supervised interest point detection and description. In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops. pp. 224–236 (2018)
11. Doersch, C., Gupta, A., Markeevea, L., Recasens, A., Smaira, L., Aytar, Y., Carreira, J., Zisserman, A., Yang, Y.: Tap-vid: A benchmark for tracking any point in a video. Advances in Neural Information Processing Systems **35**, 13610–13626 (2022)
12. Doersch, C., Yang, Y., Vecerik, M., Gokay, D., Gupta, A., Aytar, Y., Carreira, J., Zisserman, A.: Tapir: Tracking any point with per-frame initialization and temporal refinement. arXiv preprint arXiv:2306.08637 (2023)
13. Dusmanu, M., Rocco, I., Pajdla, T., Pollefeys, M., Sivic, J., Torii, A., Sattler, T.: D2-net: A trainable cnn for joint detection and description of local features. arXiv preprint arXiv:1905.03561 (2019)
14. Greff, K., Belletti, F., Beyer, L., Doersch, C., Du, Y., Duckworth, D., Fleet, D.J., Gnanapragasam, D., Golemo, F., Herrmann, C., et al.: Kubric: A scalable dataset generator. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3749–3761 (2022)
15. Harley, A.W., Fang, Z., Fragkiadaki, K.: Particle video revisited: Tracking through occlusions using point trajectories. In: European Conference on Computer Vision. pp. 59–75. Springer (2022)
16. Hartley, R., Zisserman, A.: Multiple view geometry in computer vision. Cambridge university press (2003)

17. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
18. Hong, S., Cho, S., Kim, S., Lin, S.: Unifying feature and cost aggregation with transformers for semantic and visual correspondence. In: The Twelfth International Conference on Learning Representations (2024)
19. Hong, S., Cho, S., Nam, J., Lin, S., Kim, S.: Cost aggregation with 4d convolutional swin transformer for few-shot segmentation. In: European Conference on Computer Vision. pp. 108–126. Springer (2022)
20. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: International conference on machine learning. pp. 448–456. pmlr (2015)
21. Janai, J., Güney, F., Behl, A., Geiger, A., et al.: Computer vision for autonomous vehicles: Problems, datasets and state of the art. Foundations and Trends® in Computer Graphics and Vision **12**(1–3), 1–308 (2020)
22. Jiang, W., Trulls, E., Hosang, J., Tagliasacchi, A., Yi, K.M.: Cotr: Correspondence transformer for matching across images. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6207–6217 (2021)
23. Kang, D., Kwon, H., Min, J., Cho, M.: Relational embedding for few-shot classification. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8822–8833 (2021)
24. Karaev, N., Rocco, I., Graham, B., Neverova, N., Vedaldi, A., Rupprecht, C.: Cotracker: It is better to track together. arXiv preprint arXiv:2307.07635 (2023)
25. Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al.: The kinetics human action video dataset. arXiv preprint arXiv:1705.06950 (2017)
26. Lee, A.X., Devin, C.M., Zhou, Y., Lampe, T., Bousmalis, K., Springenberg, J.T., Byravan, A., Abdolmaleki, A., Gileadi, N., Khosid, D., et al.: Beyond pick-and-place: Tackling robotic stacking of diverse shapes. In: 5th Annual Conference on Robot Learning (2021)
27. Lee, J., Kim, D., Ponce, J., Ham, B.: Sfnet: Learning object-aware semantic correspondence. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2278–2287 (2019)
28. Liu, C., Yuen, J., Torralba, A.: Sift flow: Dense correspondence across scenes and its applications. IEEE transactions on pattern analysis and machine intelligence **33**(5), 978–994 (2010)
29. Loshchilov, I., Hutter, F.: Sgdr: Stochastic gradient descent with warm restarts. arXiv preprint arXiv:1608.03983 (2016)
30. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017)
31. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. International journal of computer vision **60**, 91–110 (2004)
32. Manuelli, L., Li, Y., Florence, P., Tedrake, R.: Keypoints into the future: Self-supervised correspondence in model-based reinforcement learning. arXiv preprint arXiv:2009.05085 (2020)
33. Melekhov, I., Tiulpin, A., Sattler, T., Pollefeys, M., Rahtu, E., Kannala, J.: Dgcnet: Dense geometric correspondence network. In: 2019 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 1034–1042. IEEE (2019)
34. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM **65**(1), 99–106 (2021)

35. Min, J., Kang, D., Cho, M.: Hypercorrelation squeeze for few-shot segmentation. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 6941–6952 (2021)
36. Moing, G.L., Ponce, J., Schmid, C.: Dense optical tracking: Connecting the dots. arXiv preprint arXiv:2312.00786 (2023)
37. Nam, J., Lee, G., Kim, S., Kim, H., Cho, H., Kim, S., Kim, S.: Diffmatch: Diffusion model for dense matching. arXiv preprint arXiv:2305.19094 (2023)
38. Neoral, M., Šerých, J., Matas, J.: Mft: Long-term tracking of every pixel. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 6837–6847 (2024)
39. Oh, S.W., Lee, J.Y., Xu, N., Kim, S.J.: Video object segmentation using space-time memory networks. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9226–9235 (2019)
40. Pollefeys, M., Nistér, D., Frahm, J.M., Akbarzadeh, A., Mordohai, P., Clipp, B., Engels, C., Gallup, D., Kim, S.J., Merrell, P., et al.: Detailed real-time urban 3d reconstruction from video. International Journal of Computer Vision **78**, 143–167 (2008)
41. Pont-Tuset, J., Perazzi, F., Caelles, S., Arbeláez, P., Sorkine-Hornung, A., Van Gool, L.: The 2017 davis challenge on video object segmentation. arXiv preprint arXiv:1704.00675 (2017)
42. Press, O., Smith, N.A., Lewis, M.: Train short, test long: Attention with linear biases enables input length extrapolation. arXiv preprint arXiv:2108.12409 (2021)
43. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. The Journal of Machine Learning Research **21**(1), 5485–5551 (2020)
44. Rocco, I., Arandjelovic, R., Sivic, J.: Convolutional neural network architecture for geometric matching. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 6148–6157 (2017)
45. Rocco, I., Arandjelović, R., Sivic, J.: Efficient neighbourhood consensus networks via submanifold sparse convolutions. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16. pp. 605–621. Springer (2020)
46. Rocco, I., Cimpoi, M., Arandjelović, R., Torii, A., Pajdla, T., Sivic, J.: Neighbourhood consensus networks. Advances in neural information processing systems **31** (2018)
47. Sarlin, P.E., Cadena, C., Siegwart, R., Dymczyk, M.: From coarse to fine: Robust hierarchical localization at large scale. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12716–12725 (2019)
48. Sarlin, P.E., DeTone, D., Malisiewicz, T., Rabinovich, A.: Superglue: Learning feature matching with graph neural networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4938–4947 (2020)
49. Schonberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4104–4113 (2016)
50. Shaw, P., Uszkoreit, J., Vaswani, A.: Self-attention with relative position representations. arXiv preprint arXiv:1803.02155 (2018)
51. Sun, D., Herrmann, C., Reda, F., Rubinstein, M., Fleet, D.J., Freeman, W.T.: Disentangling architecture and training for optical flow. In: European Conference on Computer Vision. pp. 165–182. Springer (2022)

52. Sun, J., Shen, Z., Wang, Y., Bao, H., Zhou, X.: Loftr: Detector-free local feature matching with transformers. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8922–8931 (2021)
53. Tancik, M., Srinivasan, P., Mildenhall, B., Fridovich-Keil, S., Raghavan, N., Singh, U., Ramamoorthi, R., Barron, J., Ng, R.: Fourier features let networks learn high frequency functions in low dimensional domains. Advances in Neural Information Processing Systems **33**, 7537–7547 (2020)
54. Teed, Z., Deng, J.: Raft: Recurrent all-pairs field transforms for optical flow. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16. pp. 402–419. Springer (2020)
55. Tolstikhin, I.O., Houlsby, N., Kolesnikov, A., Beyer, L., Zhai, X., Unterthiner, T., Yung, J., Steiner, A., Keysers, D., Uszkoreit, J., et al.: Mlp-mixer: An all-mlp architecture for vision. Advances in neural information processing systems **34**, 24261–24272 (2021)
56. Torr, P.H., Zisserman, A.: Feature based methods for structure and motion estimation. In: International workshop on vision algorithms. pp. 278–294. Springer (1999)
57. Truong, P., Danelljan, M., Gool, L.V., Timofte, R.: Gocor: Bringing globally optimized correspondence volumes into your neural network. Advances in Neural Information Processing Systems **33**, 14278–14290 (2020)
58. Truong, P., Danelljan, M., Timofte, R.: Glu-net: Global-local universal network for dense flow and correspondences. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 6258–6268 (2020)
59. Truong, P., Danelljan, M., Timofte, R., Van Gool, L.: Pdc-net+: Enhanced probabilistic dense correspondence network. IEEE Transactions on Pattern Analysis and Machine Intelligence (2023)
60. Truong, P., Danelljan, M., Van Gool, L., Timofte, R.: Learning accurate dense correspondences and when to trust them. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5714–5724 (2021)
61. Ulyanov, D., Vedaldi, A., Lempitsky, V.: Instance normalization: The missing ingredient for fast stylization. arXiv preprint arXiv:1607.08022 (2016)
62. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)
63. Vecerik, M., Doersch, C., Yang, Y., Davchev, T., Aytar, Y., Zhou, G., Hadsell, R., Agapito, L., Scholz, J.: Robotap: Tracking arbitrary points for few-shot visual imitation. arXiv preprint arXiv:2308.15975 (2023)
64. Wang, Q., Chang, Y.Y., Cai, R., Li, Z., Hariharan, B., Holynski, A., Snavely, N.: Tracking everything everywhere all at once. arXiv preprint arXiv:2306.05422 (2023)
65. Woo, S., Park, J., Lee, J.Y., Kweon, I.S.: Cbam: Convolutional block attention module. In: Proceedings of the European conference on computer vision (ECCV). pp. 3–19 (2018)
66. Wu, Y., He, K.: Group normalization. In: Proceedings of the European conference on computer vision (ECCV). pp. 3–19 (2018)
67. Xiao, J., Chai, J.x., Kanade, T.: A closed-form solution to non-rigid shape and motion recovery. In: Computer Vision-ECCV 2004: 8th European Conference on Computer Vision, Prague, Czech Republic, May 11–14, 2004. Proceedings, Part IV 8. pp. 573–587. Springer (2004)

68. Yi, K.M., Trulls, E., Ono, Y., Lepetit, V., Salzmann, M., Fua, P.: Learning to find good correspondences. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2666–2674 (2018)
69. Zheng, Y., Harley, A.W., Shen, B., Wetzstein, G., Guibas, L.J.: Pointodyssey: A large-scale synthetic dataset for long-term point tracking. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 19855–19865 (2023)