# CoTracker: It is Better to Track Together

**Nikita Karaev**[1,2]     **Ignacio Rocco**[1]     **Benjamin Graham**[1]     **Natalia Neverova**[1]
**Andrea Vedaldi**[1]     **Christian Rupprecht**[2]

[1] Meta AI          [2] Visual Geometry Group, University of Oxford
nikita@robots.ox.ac.uk

## Abstract

Methods for video motion prediction either estimate jointly the instantaneous motion of all points in a given video frame using optical flow or independently track the motion of individual points throughout the video. The latter is true even for powerful deep-learning methods that can track points through occlusions. Tracking points individually ignores the strong correlation that can exist between the points, for instance, because they belong to the same physical object, potentially harming performance. In this paper, we thus propose CoTracker, an architecture that jointly tracks multiple points throughout an entire video. This architecture combines several ideas from the optical flow and tracking literature in a new, flexible and powerful design. It is based on a transformer network that models the correlation of different points in time via specialised attention layers. The transformer iteratively updates an estimate of several trajectories. It can be applied in a sliding-window manner to very long videos, for which we engineer an unrolled training loop. It can track from one to several points jointly and supports adding new points to track at any time. The result is a flexible and powerful tracking algorithm that outperforms state-of-the-art methods in almost all benchmarks. Project page: https://co-tracker.github.io

## 1 Introduction

Establishing point correspondences is a fundamental problem in computer vision, necessary for many downstream tasks. Here, we are interested in estimating point correspondences in a video containing dynamic objects and a moving camera. This task is often referred to as motion estimation and is interpreted as follows. Given one or more 2D points, which are the projection of certain physical points of the underlying 3D scene, the goal is to find the location of the same physical points in all other frames of the video.

There are two variants of the video motion estimation problem. In *optical flow*, the objective is to estimate the velocity of all points within a video frame. This estimation is performed jointly for all points, but the motion is only predicted at an infinitesimal distance. In *tracking*, the goal is to estimate the motion of points over an extended period. For efficiency and modelling simplicity, tracking methods typically focus on a sparse selection of points and treat them as statistically independent. Even recent techniques such as TAP-Vid [10] and Particle Video Revisited [15], which employ modern deep networks and can track points even in the presence of occlusions, model tracks independently. This approximation is crude because points are often strongly correlated (e.g., because they belong to the same physical object).

In this paper, we hypothesise that accounting for the correlation between tracked points can significantly improve tracking accuracy. Thus, we propose a new neural tracker that supports joint tracking of several points in long video sequences. Our design is inspired by prior work on deep networks for optical flow and point tracking, but with several changes to support joint estimation of multiple tracks and windowed application for processing longer videos.

Figure 1: **Qualitative Examples**. CoTracker predicts the tracks for an arbitrary selection of points in long videos. In these examples from DAVIS [29], the points are selected on an object of interest, but any point can be tracked. Here, we compensate for camera motion by using predicted background tracks and visualize the object tracks with colored lines.

Our neural network takes as input a video and a variable number of starting track locations and outputs the full tracks (Fig. 1). The design is flexible because it allows to track arbitrary points, selected at any spatial location and time in the video.

The network works by taking an initial, approximate version of the tracks and then *incrementally* refining them to better match the content of the video. Tracks are initialized trivially: given a track point or a track fragment, the rest of the track is initialized by assuming that the point remains stationary. In this manner, tracks can be initialized starting from any point, even in the middle of a video, or from the output of the tracker itself when operated in a sliding-window fashion. All these cases are supported seamlessly by the same architecture.

The network itself is a transformer operating on a 2D grid of tokens: the first grid dimension represents time, and the second is the set of tracked points. Via suitable self-attention operators, the transformer can consider each track as a whole for the duration of a window and can exchange information between tracks, exploiting their correlation. We train our network on synthetic TAP-Vid-Kubric [10], which we show to be particularly effective for this task.

We evaluate our tracker on three different benchmarks. We pay particular attention to the design of the evaluation protocol to ensure its fairness. Specifically, we note that existing benchmarks contain ground truth tracks which are concentrated on a few foreground objects, potentially revealing the presence of such objects to a joint tracker like ours. In order to ensure that no ground-truth information is leaked to our tracker, we test different distributions of tracked points. Specifically, we track one benchmark track at a time but add additional points on a grid to allow the model to perform joint tracking. In this way, our comparison is fair to single-point trackers.

Our architecture works well for tracking single points and excels for groups of points, obtaining state-of-the-art tracking performance in several benchmarks. In particular, we improving performance on the TAP-Vid [10] and FastCapture [31] datasets by a large margin, even for long tracks.

## 2   Related work

**Optical flow.** Optical flow estimates dense instantaneous motion; it is one of the longest studied problems in computer vision and was initially approached by solving simple differential equations arising from color constancy [16, 24, 4, 5]. Nowadays, optical flow is approached by using deep learning. Dosovitskiy *et al.* [12] introduced FlowNet, the first end-to-end convolutional network for optical flow, which was later improved in FlowNet2 [18] by applying a stacked architecture with warping. DCFlow [46] proposed to construct a 4D cost volume, which captures the relationship between the features of the source and target images. This approach has since become a critical building block in subsequent flow estimation algorithms [39, 43, 41]. PWC-Net [39] adopted this cost volume but introduced pyramidal processing and warping to reduce the computational cost. Teed and Dang proposed RAFT [41] that preserves high-resolution flow and incrementally updates the flow field. The success of RAFT inspired a wave of subsequent works [21, 45, 21, 47] that attempted to improve model efficiency or flow accuracy.

Transformers [42] have also been applied to the optical flow problem [17, 34, 48]. Flowformer [17] drew inspiration from RAFT and proposed a transformer-based approach that tokenizes the 4D cost volume. GMFlow [48] replaced the update network with a softmax with self-attention for refinement. Perceiver IO [19] proposed a unified transformer for several tasks including optical flow.

Whereas traditional optical flow algorithms predict motion between two frames, our focus is long-term tracking, which optical flow estimators can only approximate by temporal integration.

**Multi-frame optical flow.** There have been several attempts [30, 33, 20, 30] to extend optical flow to multiple frames. Traditionally, Kalman filtering [9, 13] was a staple mechanism to ensure temporal consistency. Modern multi-frame methods produce dense flow. RAFT [41] can be applied in a warm-start manner [37, 40] for multi-frame optical flow estimation. VideoFlow [33] extends optical flow to three and five consecutive frames by integrating forward and backward motion features during flow refinement. However, these methods are not designed for long-term tracking and do not consider points occluded for a long time.

**Tracking Groups.** Before the emergence of deep learning, some authors proposed handcrafted joint trackers [3, 32], but there is less work on doing so with modern deep network architectures like ours. Our work is weakly related to multiple *object* tracking [8] where tracking through occlusions [49], with changes in appearance [25], and with temporal priors [35] have been extensively studied. However, our focus is the tracking of points, not objects.

**Tracking any (physical) point.** TAP-Vid [10] introduced the problem of tracking any physical point in a video and proposed a benchmark and a simple baseline method for it. This method computes a cost volume for each pair of frames independently, feeding it to occlusion and coordinate regression branches. The regression branch is unable to locate occluded points. Particle Video Revisited [15] revisits the classic Particle Video [32] problem by introducing a model for point tracking through occlusions. The method tracks selected points with a fixed sliding window and restarts from the last frame where the point is visible. However, the model loses the target if it stays occluded beyond the size of the window. Furthermore, while the original Particle Video does track points jointly, Particle Video Revisited and TAP-Vid track points independently.

Several concurrent works have been developed roughly at the same time as our work [44, 27, 11]. OmniMotion [44] optimizes a volumetric representation for each video during test-time, refining estimated correspondences in a canonical space. MFT [27] conducts optical flow estimation between distant frames and chooses the most reliable chain of optical flows. TAPIR [11] is a feed-forward point tracker with a matching stage inspired by TAP-Vid [10] and a refinement stage inspired by PIPs [15].

Modern trackers and optical flow estimators are commonly trained using synthetic datasets [22, 26, 12, 38, 6], as annotating real data can be challenging. Synthetic datasets provide accurate annotations, and training on them has demonstrated the ability to generalize to real-world data [26, 38].

# 3 CoTracker

Our goal is to track 2D points throughout the duration of a video. We formalise the problem as follows. A video $V = (I_t)_{t=1}^T$ is a sequence of $T$ RGB frames $I_t \in \mathbb{R}^{3 \times H \times W}$. Tracking amounts to producing a track $P_t^i = (x_t^i, y_t^i) \in \mathbb{R}^2$, $t = t^i, \ldots, T$, $i = 1, \ldots, N$, for each of the $N$ points, where $t^i \in \{1, \ldots, T\}$ denotes the time when the track starts. Furthermore, the *visibility flag* $v_t^i \in \{0, 1\}$ tells whether a point is visible or occluded in a given frame. We assume that a track's starting point identifies the corresponding physical point unambiguously [10], which requires that the point is visible at the start (*i.e.*, $v_{t_i}^i = 1$). Note that tracks can start at any time $1 \leq t^i \leq T$ during the video.

A *tracker* is an algorithm that takes as input the video $V$ as well as the starting locations and times $(P_{t^i}^i, t^i)_{i=1}^N$ of $N$ tracks, and outputs an estimate $\hat{P}_t^i = (\hat{x}_t^i, \hat{y}_t^i)$ of the tracks for all valid ($t \geq t^i$) times. We also task the tracker with predicting an estimate $\hat{v}_i^t$ of the visibility flags. Of these values, only the initial ones $\hat{P}_{t^i}^i = P_{t^i}^i$ and $\hat{v}_{t^i}^i = v_{t^i}^i = 1$ are known to the tracker, and the others must be predicted.

## 3.1 Transformer formulation

We approach this prediction problem using a transformer network $\Psi : G \mapsto O$, which we call CoTracker. The goal of the transformer is to improve a given estimate of the tracks. Tracks are encoded as a grid of input tokens $G_t^i$, one for each track $i = 1, \ldots, N$, and time $t = 1, \ldots, T$. The updated tracks are expressed by a corresponding grid of output tokens $O_t^i$. These tokens are built as
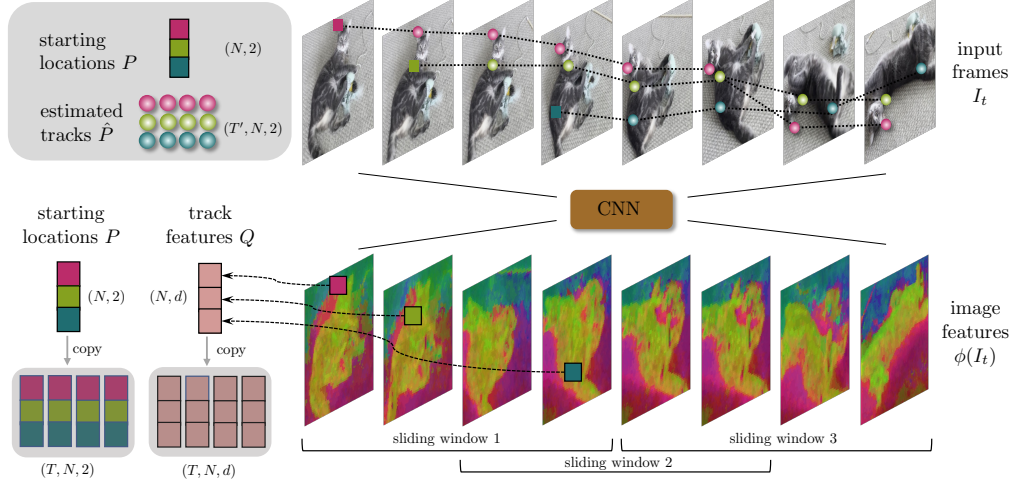
Figure 2: **CoTracker architecture.** We compute convolutional features $\phi(I_t)$ for every frame and process them with sliding windows. In order to initialize track features $Q$, we bilinearly sample from $\phi(I_t)$ with starting point locations $P$. Locations $P$ also serve to initialize estimated tracks $\hat{P}$. See Fig. 3 for a detailed visualization of one sliding window.

follows, using a design inspired by RAFT [41] for optical flow and PIPs [15] for tracking. Please see Fig. 2 for a general overview and Fig. 3 the model components.

**Image features.** We extract dense $d$-dimensional appearance features $\phi(I_t) \in \mathbb{R}^{d \times \frac{H}{k} \times \frac{W}{k}}$ from each video frame $I_t$ by using a convolutional neural network that we train with the transformer end-to-end. Here, $k = 8$ is a downsampling factor, utilised for efficiency. We consider several scaled versions $\phi(I_t; s) \in \mathbb{R}^{d \times \frac{H}{sk} \times \frac{W}{sk}}$, of these features with strides $s = 1, \ldots, S$. These downscaled features are obtained by applying average pooling to the base features in $s \times s$ neighbourhoods. We use $S = 4$ scales.

**Track features.** The appearance of the tracks is captured by feature vectors $Q_t^i \in \mathbb{R}^d$ (these are time-dependent to accommodate changes in the track appearance). They are initialized at first by sampling image features, and then they are updated by the neural network, as explained below.

**Correlation features.** In order to facilitate matching tracks to images, we introduce the correlation features $C_t^i \in \mathbb{R}^S$. These features are obtained by comparing the track features $Q$ and the image features $\phi(I_t; s)$ around the current track locations $\hat{P}_i^t$. Specifically, the vector $C_i^t$ is obtained by stacking the inner products

$$\left\langle Q_t^i, \ \phi(I_t; s)[\hat{P}_t^i/ks + \delta] \right\rangle, \quad s = 1, \ldots, S, \quad \delta \in \mathbb{Z}^2, \quad \|\delta\|_1 \le \Delta, \quad \Delta \in \mathbb{N},$$

where the offsets $\delta$ define a neighborhood of point $\hat{P}_i^t$. The image features $\phi(I_t; s)$ are sampled at non-integer locations by using bilinear interpolation and zero padding. The dimension of $C_t^i$ is $\left(2(\Delta^2 + \Delta) + 1\right)S = 164$ for our choice $S = \Delta = 4$.

**Tokens.** The input tokens $G(\hat{P}, \hat{v}, Q)$ code for position, visibility, appearance, and correlation of the tracks. This information is represented by stacking corresponding features:

$$G_t^i = (\hat{P}_t^i, \ \text{logit}(\hat{v}_t^i), \ Q_t^i, \ C_t^i, \ \eta(\hat{P}_t^i - \hat{P}_1^i)).$$

All the components except the last one have been introduced above. The last component is derived from the estimated position: it is the sinusoidal positional encoding $\eta$ of the track location with respect to the initial location at time $t = 1$. The latter information could be inferred by the transformer by observing $\hat{P}_t^i$ alone, but we found it beneficial to pass it directly as well.

The output tokens $O(\hat{P}', Q')$ contain the updated locations and appearance features only, *i.e.* $O_t^i = (\hat{P}_t^{i\,\prime}, Q_t^{i\,\prime})$.
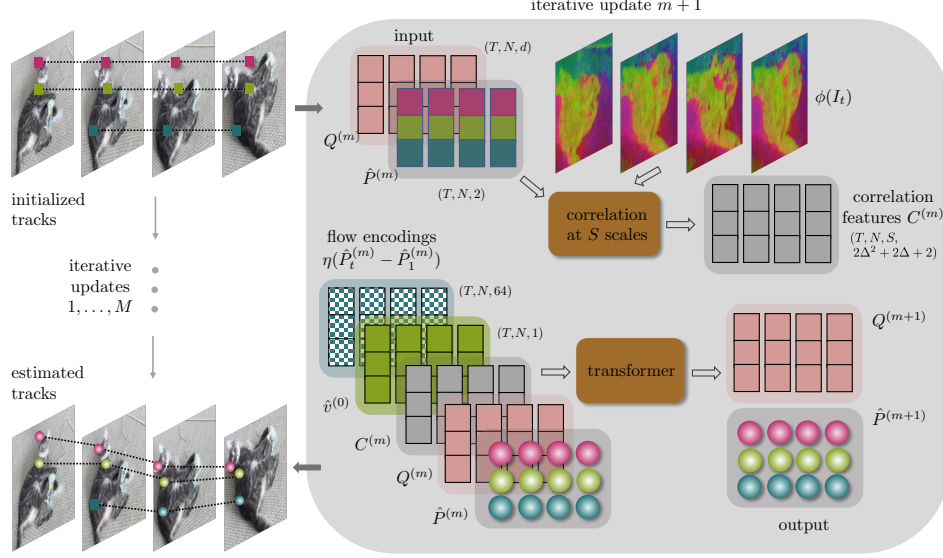
4

Figure 3: **CoTracker architecture.** Visualization of one sliding window with $M$ iterative updates. During one iteration, we update point tracks $\hat{P}^{(m)}$ and track features $Q^{(m)}$. $Q^{(0)}$ is initialized with $Q$ for all sliding windows, $\hat{P}^{(0)}$ with start tracks for the unseen frames and with previous predictions for frames processed in the previous sliding window. Visibility $\hat{v}$ is computed after the last update $M$.

**Iterated transformer applications.** We apply the transformer $M$ times in order to progressively improve the track estimates. Let $m = 0, 1, \ldots, M$ index the estimate, with $m = 0$ denoting initialization. Then:

$$O(\hat{P}^{(m+1)}, Q^{(m+1)}) = \Psi(G(\hat{P}^{(m)}, \hat{v}^{(0)}, Q^{(m)})).$$

Note that the visibility mask $\hat{v}$ is not updated by the transformer; instead, it is updated once at the end of the $M$ applications of the transformer as $\hat{v}^{(M)} = \sigma(WQ^{(M)})$, where $\sigma$ is the sigmoid activation function and $W$ is a learned matrix of weights. We found iterative updates for the visibility did not further improve the performance, likely due to the fact that visibility highly depends on predicting an accurate location first.

The quantities $\hat{P}^{(0)}$, $v^{(0)}$ and $Q^{(0)}$ are initialized from the starting location and time of the tracks. For all tracks $i = 1, \ldots, N$ and times $t = 1, \ldots, T$, we simply set

$$\hat{P}_t^{i,(0)} \leftarrow P_{t^i}^i, \quad \hat{v}_t^{i,(0)} \leftarrow 1, \quad [Q_t^{i,(0)}]_s \leftarrow \phi(I_{t^i}; s)[P_{t^i}^i/ks], \tag{1}$$

effectively broadcasting $t^i$ to all other times $t = 1, \ldots, T$ (both before and after $t^i$).

## 3.2 Windowed inference

An advantage of our transformer design is that it can easily support windowed applications in order to process very long videos. Consider in particular a video $V$ of length $T' > T$ longer than the maximum window length $T$ supported by the architecture. To track points throughout the entire video $V$, we split the video in $J = \lceil 2T'/T - 1 \rceil$ windows of length $T$, with an overlap of $T/2$ frames.[1]

In order to process a video, we apply the transformer $MJ$ times: the output of the first window is used as the input to the second window and so on. Let the superscript $(m, j)$ denote the $m$-th update of the transformer applied to the $j$-th window. We thus have a $M \times J$ grid of quantities $(\hat{P}^{(m,j)}, \hat{v}^{(m,j)}, Q^{(m,j)})$, spanning transformer iterations and windows. Starting from $m = 0$ and $j = 1$, we initialize these quantities using Eq. (1) as before. We then apply the transformer $M$ times to obtain the state $(\hat{P}^{(M,1)}, \hat{v}^{(M,1)}, Q^{(M,1)})$. From this, we initialize $(\hat{P}^{(0,2)}, \hat{v}^{(0,2)}, Q^{(0,2)})$ for the second window, in a manner similar to Eq. (1). Specifically, the first $T/2$ components of $\hat{P}^{(0,2)}$ are

---

[1]We assume that $T$ is even. The last window is shorter if $T/2$ does not divide $T'$.

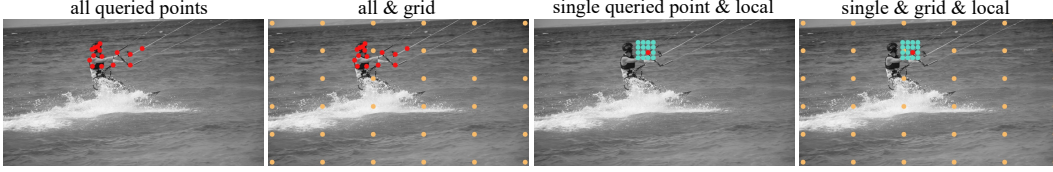| all queried points | all & grid | single queried point & local | single & grid & local |

Figure 5: Sampling strategies for target points and additional points. The "global" and "local" strategies provide the model with global and local context to the target points, respectively.

copies of the last $T/2$ components of $\hat{P}^{(M,1)}$; the last $T/2$ components of $\hat{P}^{(0,2)}$ are instead copies of the last time $t = T/2 - 1$ from $\hat{P}^{(M,1)}$. The same update rule is used for $\hat{v}^{(0,2)}$, while $Q^{(0,j)}$ is always initialized with initial track features $Q$. After initializing $(\hat{P}^{(0,2)}, \hat{v}^{(0,2)}, Q^{(0,2)})$, the transformer is applied $M$ more times to the second window, and the process is repeated for the next window and so on. Finally, any new track is added by extending the token grids using initialization (1).

## 3.3 Unrolled learning

We found it important to learn the windowed transformer in an unrolled fashion in order to properly handle semi-overlapping windows. The primary loss is for track regression, summed over iterated transformer applications and windows:

$$\mathcal{L}_1(\hat{P}, P) = \sum_{j=1}^{J} \sum_{m=1}^{M} \gamma^{M-m} \|\hat{P}^{(m,j)} - P^{(j)}\|, \tag{2}$$

where $\gamma = 0.8$ discounts early transformer updates. Here $P^{(j)}$ contains the ground-truth trajectories restricted to window $j$ (trajectories which start in the middle of the window are padded backwards). The second loss is the cross entropy of the visibility flags $\mathcal{L}_2(\hat{v}, v) = \sum_{j=1}^{J} \mathrm{CE}(\hat{v}^{(M,j)}, v^{(j)})$.

While only a moderate number of windows are used in the loss during training due to the computational cost, at test time we can unroll the windowed transformer applications arbitrarily, thus in principle handling any video length.

Unrolled inference allows tracking points that appear later in the video. We start tracking a point only from the sliding window where it appears first. We also make sure that such points are present in the training data by sampling visible points from the middle frame of a sequence.

## 3.4 Transformer

Our transformer consists of interleaving time and group attention blocks with two linear layers applied to the input and to the output. Factorising the attention [1] across time and point tracks makes the model computationally tractable: the complexity is reduced from $O(N^2 T^2)$ to $O(N^2 + T^2)$. In addition to position encodings $\eta$ for estimated trajectories that the transformer takes as input, we add standard sinusoidal position encodings: 1-dimensional for time and 2-dimensional for space.



Figure 4: **Number of points**. We track queried points from TAP-Vid-DAVIS together with $N$ additional points sampled on the first frame. Adding 16-25 points is best, while more points overwhelm tracking the interest points.

## 3.5 Point Selection

One of the key insights of our method is that tracking multiple points simultaneously allows the model to better reason about the motion in the video and the correlation between tracks.

However, care needs to be applied when the method is evaluated on benchmark datasets. In these datasets, points annotated by humans often lie
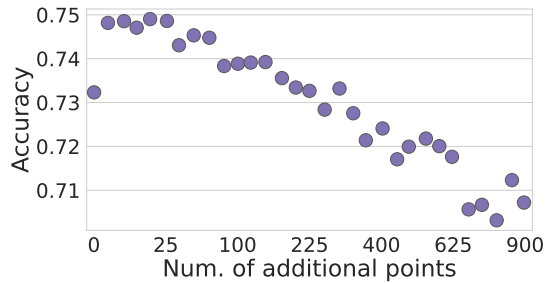
6

Table 1: **Joint tracking**. For time-only attention, we train a variant model with the 6 group attention layers replaced by time attention so that the model has 12 time attention layers in total, maintaining the same model size overall. For "single target points" we track target points separately by combining one target point with additional points sampled on a grid, locally, or both (see Section 3.5 and Fig. 5).

| Attention | | Setting | | DAVIS First | | | BADJA | |
|---|---|---|---|---|---|---|---|---|
| time | group | tgt. points | add. points | AJ | $< \delta^x_{avg}$ | OA | seg-based | $< \delta^{3px}$ |
| ✓ | ✗ | single | - | 56.3 | 72.6 | 84.8 | 59.5 | 15.0 |
| ✓ | ✓ | single | - | 39.8 | 61.6 | 78.8 | 55.7 | 13.8 |
| ✓ | ✓ | single | grid | 54.4 | 70.4 | 86.8 | 60.5 | 16.6 |
| ✓ | ✓ | single | local | 57.8 | 75.1 | 86.8 | 62.9 | 17.2 |
| ✓ | ✓ | single | grid & local | **60.6** | **75.4** | **89.3** | **63.6** | **18.0** |
| ✓ | ✓ | all | - | 54.3 | 73.0 | 84.9 | 64.0 | 18.2 |
| ✓ | ✓ | all | grid | **58.8** | **74.6** | **88.5** | **64.2** | **18.6** |

Table 2: **Evaluation on TAP-Vid**. We report numbers for the "queried first" evaluation protocol and for the easier "queried strided" protocol on DAVIS. $< \delta^x_{avg}$ measures position accuracy of visible points, averaged across different thresholds. OA stands for Occlusion Accuracy, and AJ represents Average Jaccard, which evaluates both occlusion and coordinate accuracy.

| Method | DAVIS First | | | DAVIS Strided | | | Kinetics First | | |
|---|---|---|---|---|---|---|---|---|---|
| | AJ | $< \delta^x_{avg}$ | OA | AJ | $< \delta^x_{avg}$ | OA | AJ | $< \delta^x_{avg}$ | OA |
| TAP-Net [10] | 33.0 | 48.6 | 78.8 | 38.4 | 53.1 | 82.3 | 38.5 | 54.4 | 80.6 |
| PIPs [15] | 42.2 | 64.8 | 77.7 | 52.4 | 70.0 | 83.6 | 31.7 | 53.7 | 72.9 |
| MFT [27] | 47.3 | 66.8 | 77.8 | 56.1 | 70.8 | 86.9 | - | - | - |
| OmniMotion [44] | - | - | - | 51.7 | 67.5 | 85.3 | - | - | - |
| TAPIR [11] | <u>56.2</u> | <u>70.0</u> | <u>86.5</u> | <u>61.3</u> | <u>73.6</u> | **88.8** | **49.6** | <u>64.2</u> | <u>85.0</u> |
| CoTracker (Ours) | **60.6** | **75.4** | **89.3** | **64.8** | **79.1** | <u>88.7</u> | <u>48.7</u> | **64.3** | **86.5** |

in salient locations of moving objects. In order to ensure that our performance figures are robust with respect to point selection, *and* to ensure a rigorously fair comparison with existing methods, we evaluate the model using only a single target ground truth point at a time. This decouples the performance of the model from the point distribution in the dataset.

We experiment with two point selection strategies that are visualized in Fig. 5. With the "global" strategy, we simply select additional points on a regular grid across the whole image. With the "local" strategy, we select additional points close to the target point, using a regular grid around the target point, thus allowing the model to focus on a neighbourhood of it. Point selection is only used at inference time.

## 4 Experiments

Our model is trained with ground truth trajectories on synthetic data. Evaluation is then performed on four benchmark datasets containing manually annotated trajectories in real videos: TAP-Vid-DAVIS [10], TAP-Vid-Kinetics [10], BADJA [2], and FastCapture [31].

### 4.1 Datasets and evaluation protocol

We consider several datasets for training and evaluation. **FlyingThings++** is a version of FlyingThings3D [26] that Particle Video Revisited [32] has re-engineered to learn their tracker PIPs. It contains videos 8 frames long with flying objects and point tracks randomly selected on these objects. Occlusion annotations are created by overlaying additional objects on top of the rendered video, so they are only approximate. **TAP-Vid** contains four benchmark datasets for evaluation and one synthetic dataset for training. The latter, TAP-Vid-Kubric, consists of 24-frame sequences rendered from 3D scenes with objects dropped into them, using the Kubric engine [14]. Point tracks are selected randomly, primarily on objects, and some on the background. Occlusions appear naturally when a point is occluded by another object or moves out of the scene. We found that training on

Table 3: **Point tracking benchmarks**.

| Method | BADJA | | FastCapture |
|---|---|---|---|
| | seg-based | $< \delta^{3px}$ | $< \delta^{3px}$ |
| DINOv2 [28] | <u>64.7</u> | 3.4 | 12.1 |
| RAFT [41] | 49.7 | 7.6 | 32.8 |
| TAP-Net [10] | 54.4 | 6.3 | 43.3 |
| PIPs [15] | 61.9 | 13.5 | 39.9 |
| TAPIR [11] | **66.9** | <u>15.2</u> | <u>63.2</u> |
| CoTracker (Ours) | 63.6 | **18.0** | **68.0** |

Table 4: **Training data**. For a fair comparison, we train PIPs on Kubric which improves performance.

| Method | Training data | DAVIS First | | | Kinetics First | | | BADJA | |
|---|---|---|---|---|---|---|---|---|---|
| | | AJ | $< \delta^x_{avg}$ | OA | AJ | $< \delta^x_{avg}$ | OA | seg-based | $< \delta^{3px}$ |
| PIPs [15] | Flying things++ | 42.2 | 64.8 | 77.7 | 31.7 | 53.7 | 72.9 | 61.9 | 13.5 |
| PIPs [15] | TAP-Vid-Kubric | 43.8 | 66.5 | 77.4 | 32.7 | 55.1 | 72.9 | **64.0** | 14.7 |
| CoTracker (Ours) | TAP-Vid-Kubric | **60.6** | **75.4** | **89.3** | **48.7** | **64.3** | **86.5** | 63.6 | **18.0** |

TAP-Vid-Kubric improves PIPs on all the benchmarks (Table 4), and thus chose this dataset to train our model as well.

We evaluate our models on the TAP-Vid benchmark datasets TAP-Vid-DAVIS (30 videos of ∼100 frames) and TAP-Vid-Kinetics (1144 videos of ∼250 frames from Kinetics [7]). In these benchmarks, points are queried on objects at random frames and the goal is to predict positions and occlusion labels of queried points. In the TAP-Vid "queried first" evaluation protocol, each point is queried only once in the video, at the first frame where it becomes visible. Hence, the model should predict positions only for future frames. However, in the "queried strided" protocol, points are queried every five frames and tracking should be done in both directions. Given that ours (and PIPs) are online methods, which means that they track points only forward, we run them forward and backwards starting from each queried point. As "queried first" requires estimating the longest tracks, it is a more difficult setting than "strided". Moreover, "strided" demands estimating the same track from multiple starting locations and is thus much more computationally expensive. For these two reasons, we focus on the "strided" setting for the evaluation. We evaluate the trackers using the standard metrics for these datasets: Occlusion Accuracy (OA), $< \delta^x_{avg}$ (position accuracy for visible points), and Average Jaccard (AJ), averaged across different thresholds.

**BADJA** [2] is a benchmark dataset of animals in motion with keypoint annotations. These are annotated in the reference frame and the goal is to propagate them to future frames. A keypoint coordinate is predicted accurately if it is within the distance of $0.2\sqrt{A}$ from the ground truth coordinate, where $A$ is the area of the ground-truth segmentation mask on the frame. Accuracy is measured only for visible points. Both BADJA and TAP-Vid-DAVIS are based on the DAVIS dataset [29]. **FastCapture** is a benchmark of dynamic human reconstructions rendered on a black background. We sample twenty points on humans in the first frame and evaluate the 3px accuracy.

## 4.2 Implementation details

In this section, we describe the most important implementation details and refer to the appendix for a complete description.

**Training.** We render 11,000 pre-generated 24-frame sequences for TAP-Vid-Kubric with annotations for 2,000 tracked points per sequence. Points are sampled preferentially on objects. During training, for each sequence, we randomly sample 256 points that are visible either in the first or in the middle frames of the sequence.

We train CoTracker on TAP-Vid-Kubric sequences of size $T' = 24$ with sliding windows of size $T = 8$ for 50,000 iterations using 32 NVIDIA TESLA Volta V100 32GB GPUs.
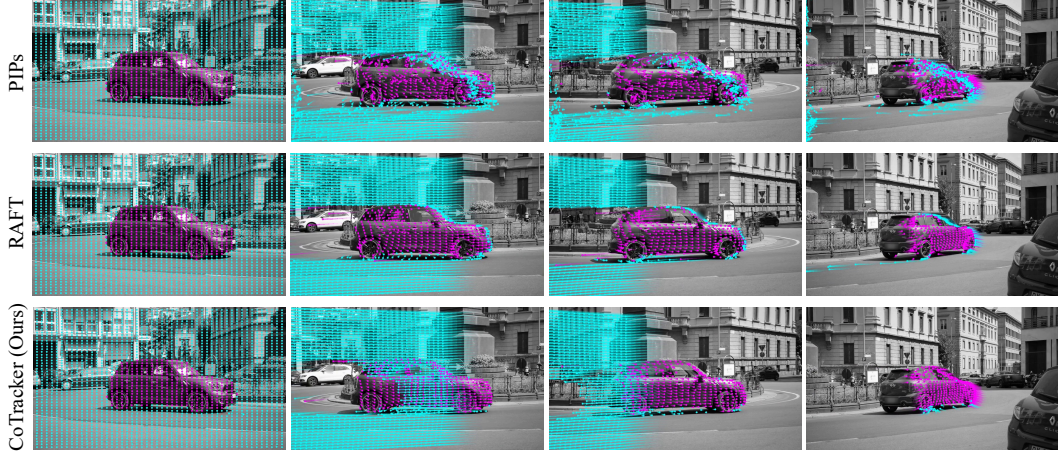
Figure 6: **Qualitative Results**. For PIPs (top), many points are incorrectly tracked and end up being 'stuck' on the front of the car or the side of the image when they become occluded. RAFT (middle) predictions have less noise, but it fails to handle occlusions, leading to points getting stuck on the car as well. Our CoTracker (bottom) produces cleaner tracks. The tracks are also more 'linear' than PIPs tracks, which is accurate as the primary motion is a homography (the observer does not translate).

Table 5: **Unrolling.** Comparing unrolling the sliding window through time during training shows that it considerably improves the performance, as the model learns to track points over long periods of time.

| Unrolled | DAVIS First | | | BADJA | |
|---|---|---|---|---|---|
| | AJ | $< \delta_{avg}^{x}$ | OA | seg-based | $< \delta^{3px}$ |
| ✗ | 46.9 | 69.1 | 80.1 | 55.8 | 13.6 |
| ✓ | **60.6** | **75.4** | **89.3** | **63.6** | **18.0** |

It is possible to train the model on GPUs with less memory by decreasing the number of sampled points per batch.

**Evaluation.** We evaluate CoTracker with 6 global and 6 local grid points (Fig. 5). All the models are evaluated on videos with a maximum resolution $384 \times 512$. For the DINOv2 baseline, we compute features using DINOv2 at a resolution of $378 \times 504$. We then bi-linearly sample from the queried feature map and compute the dot product between the queried feature and all the other feature maps. After that, we apply softmax to predict the coordinates of the queried point. PIPs is trained at a resolution of $384 \times 512$ but evaluated with $320 \times 512$. Using the released code, we found that $384 \times 512$ slightly improves PIPs performance on all the benchmarks except BADJA. For the TAP-Vid benchmarks, we follow the standard protocol and downsample videos to $256 \times 256$ before passing them to the model. As PIPs and CoTracker are trained on $384 \times 512$ videos, we resize downsampled frames back to $384 \times 512$ for them, and to $378 \times 504$ for DINOv2. All the TAP-Vid metrics are then computed in $256 \times 256$.

### 4.3 Results

**CoTracker: joint tracking and support grids.** First, in Table 1 we compare different support grids (Section 3.5). Using the uncorrelated 'single target point' protocol, reasoning about tracks jointly (group attn.) in addition to along tracks (time attn.) improves results provided that the correct contextual points are considered (combining local and global grids performs best). Performance increases even further for the 'all target points' protocol, likely because the different target points are indeed highly correlated. We do not use this protocol when comparing to prior work for fairness. It is, however, a realistic scenario because a segmentation model can be used to pick such correlated points automatically.

9

In Tables (2, 3) we compare our model (single target point evaluation) to the prior state of the art on the four benchmark datasets. Our approach is able to track points and their visibility more accurately than previous work, supporting the intuition that long-term tracking of points in groups is beneficial to single point models (such as TAP-Net and PIPs) and to dense but short-term optical flow methods (such as RAFT) that tend to accumulate drift.

**Training Data.** We evaluate the influence of the training data on the model in Table 4. Kubric is better for training than FlyingThings++ for the prior state of the art. FlyingThings++ sequences are too short to train our model which relies on training with sliding windows. Additionally, Kubric is overall more realistic, as it is created by rendering 3D scenes with naturally occluded objects and physics simulation.

**Sliding Window.** In Table 5, we evaluate the importance of unrolling the sliding window scheme during training. As the benchmark sequences during evaluation are often much longer ($> 10\times$), our model greatly benefits from learning to propagate information between windows.

### 4.4 Limitations

While CoTracker improves over the prior state of the art, several limitations remain. As a sliding-window-based method, it cannot track points through occlusions that are longer than the size of a single window. Additionally, the transformer complexity is quadratic in the number of tracked points, which prevents a trivial application of this technique to dense prediction.

## 5 Conclusions

We have presented CoTracker, a new paradigm for long-term tracking of points in videos. Our method simultaneously tracks groups of points and thus learns to account for their correlation. Additionally, we show that training through an unrolled temporally sliding window mechanism allows the model to generalise to long video sequences. Our method outperforms the current state of the art in point tracking on various benchmarks and produces tracks which are qualitatively cleaner.

## 6 Acknowledgments

## References

[1] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *Proc. ICML*, July 2021. 6

[2] Benjamin Biggs, Thomas Roddick, Andrew Fitzgibbon, and Roberto Cipolla. Creatures great and SMAL: Recovering the shape and motion of animals from video. In *Proc. ACCV*, 2018. 7, 8

[3] Stanley T Birchfield and Shrinivas J Pundlik. Joint tracking of features and edges. In *Proc. CVPR*, 2008. 3

[4] M. J. Black and P. Anandan. A framework for the robust estimation of optical flow. In *Proc. ICCV*, 1993. 2

[5] Andrés Bruhn, Joachim Weickert, and Christoph Schnörr. Lucas/kanade meets horn/schunck: Combining local and global optic flow methods. *IJCV*, 61, 2005. 2

[6] Daniel J Butler, Jonas Wulff, Garrett B Stanley, and Michael J Black. A naturalistic open source movie for optical flow evaluation. In *Proc. ECCV*, 2012. 3

[7] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proc. CVPR*, 2017. 8

[8] Fei Chen, Xiaodong Wang, Yunxiang Zhao, Shaohe Lv, and Xin Niu. Visual object tracking: A survey. *CVIU*, 222, 2022. 3

[9] Toshio M Chin, William Clement Karl, and Alan S Willsky. Probabilistic and sequential computation of optical flow using temporal coherence. *IEEE Trans. on Image Processing*, 3(6), 1994. 3

[10] Carl Doersch, Ankush Gupta, Larisa Markeeva, Adrià Recasens, Lucas Smaira, Yusuf Aytar, João Carreira, Andrew Zisserman, and Yi Yang. Tap-vid: A benchmark for tracking any point in a video. *arXiv*, 2022. 1, 2, 3, 7, 8, 14

[11] Carl Doersch, Yi Yang, Mel Vecerik, Dilara Gokay, Ankush Gupta, Yusuf Aytar, Joao Carreira, and Andrew Zisserman. Tapir: Tracking any point with per-frame initialization and temporal refinement, 2023. 3, 7, 8

[12] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. In *Proc. ICCV*, 2015. 2, 3

[13] Michael Elad and Arie Feuer. Recursive optical flow estimation—adaptive filtering approach. *J. of Visual Communication and Image Representation*, 9(2), 1998. 3

[14] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J Fleet, Dan Gnanapragasam, Florian Golemo, Charles Herrmann, Thomas Kipf, Abhijit Kundu, Dmitry Lagun, Issam Laradji, Hsueh-Ti (Derek) Liu, Henning Meyer, Yishu Miao, Derek Nowrouzezahrai, Cengiz Oztireli, Etienne Pot, Noha Radwan, Daniel Rebain, Sara Sabour, Mehdi S. M. Sajjadi, Matan Sela, Vincent Sitzmann, Austin Stone, Deqing Sun, Suhani Vora, Ziyu Wang, Tianhao Wu, Kwang Moo Yi, Fangcheng Zhong, and Andrea Tagliasacchi. Kubric: a scalable dataset generator. In *Proc. CVPR*, 2022. 7

[15] Adam W Harley, Zhaoyuan Fang, and Katerina Fragkiadaki. Particle video revisited: Tracking through occlusions using point trajectories. In *Proc. ECCV*, 2022. 1, 3, 4, 7, 8, 13, 14

[16] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3), 1981. 2

[17] Zhaoyang Huang, Xiaoyu Shi, Chao Zhang, Qiang Wang, Ka Chun Cheung, Hongwei Qin, Jifeng Dai, and Hongsheng Li. Flowformer: A transformer architecture for optical flow. In *Proc. ECCV*, 2022. 2

[18] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *Proc. CVPR*, 2017. 2

[19] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, Olivier J. Hénaff, Matthew M. Botvinick, Andrew Zisserman, Oriol Vinyals, and João Carreira. Perceiver IO: A general architecture for structured inputs & outputs. In *Proc. ICLR*, 2022. 2

[20] Joel Janai, Fatma Guney, Anurag Ranjan, Michael Black, and Andreas Geiger. Unsupervised learning of multi-frame optical flow with occlusions. In *Proc. ECCV*, 2018. 3

[21] Shihao Jiang, Yao Lu, Hongdong Li, and Richard Hartley. Learning optical flow from a few matches. In *Proc. CVPR*, 2021. 2

[22] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. Dynamicstereo: Consistent dynamic depth from stereo videos. In *Proc. CVPR*, 2023. 3

[23] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 13

[24] Bruce D Lucas and Takeo Kanade. An iterative image registration technique with an application to stereo vision. In *Proc. IJCAI*, volume 2, 1981. 2

[25] Lain Matthews, Takahiro Ishikawa, and Simon Baker. The template update problem. *PAMI*, 26(6), 2004. 3

[26] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *Proc. CVPR*, 2016. 3, 7

[27] Michal Neoral, Jonáš Šerých, and Jiří Matas. Mft: Long-term tracking of every pixel, 2023. 3, 7

[28] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision. *arXiv*, 2023. 8, 13

[29] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv*, 2017. 2, 8

[30] Zhile Ren, Orazio Gallo, Deqing Sun, Ming-Hsuan Yang, Erik B Sudderth, and Jan Kautz. A fusion approach for multi-frame optical flow estimation. In *Proc. WACV*, 2019. 3

[31] Ignacio Rocco, Iurii Makarov, Filippos Kokkinos, David Novotny, Benjamin Graham, Natalia Neverova, and Andrea Vedaldi. Real-time volumetric rendering of dynamic humans. *arXiv*, 2023. 2, 7

[32] Peter Sand and Seth Teller. Particle video: Long-range motion estimation using point trajectories. *IJCV*, 80, 2008. 3, 7

[33] Xiaoyu Shi, Zhaoyang Huang, Weikang Bian, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei Qin, Jifeng Dai, and Hongsheng Li. Videoflow: Exploiting temporal cues for multi-frame optical flow estimation. *arXiv*, 2023. 3

[34] Xiaoyu Shi, Zhaoyang Huang, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei Qin, Jifeng Dai, and Hongsheng Li. Flowformer++: Masked cost volume autoencoding for pretraining optical flow estimation. *arXiv*, 2023. 2

[35] Hedvig Sidenbladh, Michael J Black, and David J Fleet. Stochastic tracking of 3d human figures using 2d image motion. In *Proc. ECCV*, 2000. 3

[36] Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial intelligence and machine learning for multi-domain operations applications*, volume 11006, pages 369–386. SPIE, 2019. 13

[37] Xiuchao Sui, Shaohua Li, Xue Geng, Yan Wu, Xinxing Xu, Yong Liu, Rick Goh, and Hongyuan Zhu. Craft: Cross-attentional flow transformer for robust optical flow. In *Proc. CVPR*, 2022. 3

[38] Deqing Sun, Daniel Vlasic, Charles Herrmann, Varun Jampani, Michael Krainin, Huiwen Chang, Ramin Zabih, William T Freeman, and Ce Liu. Autoflow: Learning a better training set for optical flow. In *Proc. CVPR*, 2021. 3

[39] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *Proc. CVPR*, 2018. 2

[40] Shangkun Sun, Yuanqi Chen, Yu Zhu, Guodong Guo, and Ge Li. Skflow: Learning optical flow with super kernels. *arXiv*, 2022. 3

[41] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Proc. ECCV*, 2020. 2, 3, 4, 8, 14

[42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proc. NeurIPS*, 2017. 2

[43] Jianyuan Wang, Yiran Zhong, Yuchao Dai, Kaihao Zhang, Pan Ji, and Hongdong Li. Displacement-invariant matching cost learning for accurate optical flow estimation. *Proc. NeurIPS*, 33, 2020. 2

[44] Qianqian Wang, Yen-Yu Chang, Ruojin Cai, Zhengqi Li, Bharath Hariharan, Aleksander Holynski, and Noah Snavely. Tracking everything everywhere all at once, 2023. 3, 7

[45] Haofei Xu, Jiaolong Yang, Jianfei Cai, Juyong Zhang, and Xin Tong. High-resolution optical flow from 1d attention and correlation. In *Proc. CVPR*, 2021. 2

[46] Jia Xu, Rene Ranftl, and Vladlen Koltun. Accurate optical flow via direct cost volume processing. In *Proc. CVPR*, July 2017. 2

[47] Feihu Zhang, Oliver J Woodford, Victor Adrian Prisacariu, and Philip HS Torr. Separable flow: Learning motion cost volumes for optical flow estimation. In *Proc. CVPR*, 2021. 2

[48] Shiyu Zhao, Long Zhao, Zhixing Zhang, Enyu Zhou, and Dimitris Metaxas. Global matching with overlapping attention for optical flow estimation. In *Proc. CVPR*, 2022. 2

[49] Tao Zhao and Ramakant Nevatia. Tracking multiple humans in crowded environment. In *Proc. CVPR*, volume 2, 2004. 3

# A  Implementation details

In this section, we complete the description of implementation details from the main paper. Code and models are available on the project webpage: https://co-tracker.github.io.

**Feature CNN**  Given a sequence of $T'$ frames with a resolution 384×512, we compute features for each frame with a 2-dimensional CNN that downsamples the input image by a factor of 8 and outputs features with 128 channels. Our CNN is the same as in PIPs [15]: it consists of one $7 \times 7$ convolution with stride 2, eight residual blocks with $3 \times 3$ kernels and instance normalization, and two final convolutions with $3 \times 3$ and $1 \times 1$ kernels.

**Sliding windows**  When passing information from one sliding window to the next, we concatenate binary masks of shape $(N, T)$ with visibility logits that indicate where the model needs to make predictions. For example, masks in the first sliding window would be equal to 1 from the frame where we start tracking the point, and 0 before that. Masks for all the subsequent sliding windows will be equal to 1 for the first overlapping $T/2$ frames, and 0 for the remaining $T/2$. During training, tracking starts either from the first or from a random frame where the point is visible. If a point is not visible in the first sliding window, it will be added when it becomes visible first.

**Iterative updates**  We train the model with $M = 4$ iterative updates, and evaluate with $M = 6$. This setting is a reasonable trade-off between speed and accuracy as evaluated in Fig. 7. Interestingly, the performance is stable across a factor of 2 (2-8) around the training choice $M = 4$.
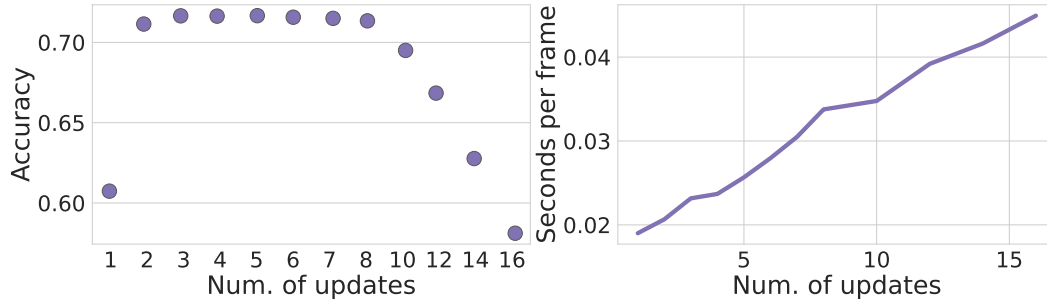


Figure 7: **Inference iterative updates.** We ablate the number of iterative updates $M$ during inference. The network is trained with $M = 4$.

**Training**  CoTracker is trained with a batch size of 32, distributed across 32 GPUs. After applying data augmentations, we randomly sample 256 trajectories for each batch, with points visible either in the first or in the middle frame. We train the model for 50,000 iterations with a learning rate of $5e^{-4}$ and a linear 1-cycle [36] learning rate schedule, using the AdamW [23] optimizer.

**Augmentations**  During training, we employ Color Jitter and Gaussian Blur to introduce color and blur variations. We augment occlusions by either coloring a randomly chosen rectangular patch with its mean color or replacing it with another patch from the same image. Random scaling across height and width is applied to each frame to add diversity.

# B  Additional ablations

Here we provide additional ablation experiments that supplement the model choices from the main paper.

**Feature network**  To understand the impact of learning image features, we replace our simple feature CNN with a frozen DINOv2 [28] ViT-s encoder. We take the output of an intermediate layer with 384 channels and linearly project it to 128 channels, as in our Feature CNN. The model is then trained with a resolution of 378x504, ensuring that both height and width are divisible by 14, which corresponds to the DINO patch size.

Table 6 shows that generic features are not precise enough for long-term tracking of points. Additionally, DINO features have been shown to contain strong semantic components which might hinder tracking performance.

Table 6: **Feature network**. Replacing the image encoder with a frozen pre-trained Vision Transformer decreases the performance.

| Feature Network | DAVIS First | | | BADJA | |
|---|---|---|---|---|---|
| | AJ | $< \delta_{avg}^x$ | OA | seg-based | $< \delta^{3px}$ |
| **CNN** | **51.2** | **68.5** | **85.2** | 65.5 | **15.4** |
| DINOv2 | 38.8 | 56.2 | 83.6 | **67.3** | 9.1 |

Table 7: **Model stride**. Higher resolution features help to make much more accurate predictions.

| Stride | DAVIS First | | | Kinetics First | | | BADJA | |
|---|---|---|---|---|---|---|---|---|
| | AJ | $< \delta_{avg}^x$ | OA | AJ | $< \delta_{avg}^x$ | OA | seg-based | $< \delta^{3px}$ |
| 8 | 51.2 | 68.5 | 85.2 | 43.3 | 61.4 | 83.3 | **65.5** | 15.4 |
| **4** | **60.6** | **75.4** | **89.3** | **48.7** | **64.3** | **86.5** | 63.6 | **18.0** |

**Model stride ablation.** In Table 7, we train CoTracker with two different feature downsampling factors (strides). The model with stride 4 has features with twice the resolution as the model with stride 8, resulting in better performance.

**Training sliding window size.** We evaluate the impact of the size of the sliding window (Table 8) during training. A longer sliding window length can improve the model's ability to handle occluded points. However, we are limited by the training data which only contains sequences of 24 frames, so learning is less effective for larger window sizes. Overall, a sliding window length of 8 yields the best results in our setting.

**Inference sliding window size.** In Table 9 we inspect how the sliding window size of a model trained with sliding window $T = 16$ affects model performance. We find that the model prefers training and evaluation with the same size across two different benchmarks. If a varying window size is needed, the model could potentially benefit from also training with a variable size.

## C   Efficiency

In Fig. 8 we evaluate the efficiency of PIPs [15], RAFT [41] and CoTracker with stride 8 on TAP-Vid-DAVIS [10] by running them on a V100 GPU. Points in this dataset can be queried on any frame. Since PIPs estimates the trajectories of every target independently and can only share computation between points when they are queried on the same frame, we run it separately for every point. RAFT estimates dense optical flow between every pair of consecutive frames. We obtain final tracks by composing estimated flows for each queried point starting from the corresponding frame. For CoTracker, we add points on a regular grid and track them all jointly. We show that inference time increases quadratically with the number of tracked points due to the group attention layers. Nevertheless, in this setting CoTracker is faster than PIPs and RAFT when tracking $< 800$ points.
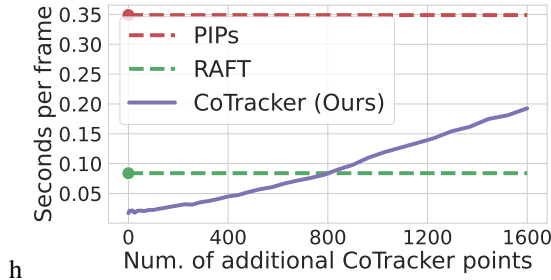


Figure 8: **Efficiency**. We track queried points from the TAP-Vid-DAVIS [10] dataset together with $N^2$ points sampled on the first frame, where $N$ is the grid size. When tracking $< 800$ points queried at different frames, CoTracker is more efficient than PIPs [15] and RAFT [41].

14

Table 8: **Training sliding window size**. We train and evaluate CoTracker with the varying window lengths. As the training data only contains sequences of length 24, larger window sizes are difficult to train.

| Window size | DAVIS First | | | BADJA | |
|---|---|---|---|---|---|
| | AJ | $< \delta_{avg}^x$ | OA | seg-based | $< \delta^{3px}$ |
| 4 | 57.9 | 73.2 | 87.0 | 64.6 | 18.2 |
| **8** | **60.6** | **75.4** | **89.3** | 63.6 | 18.0 |
| 12 | 58.8 | **75.4** | 87.6 | **65.1** | 17.7 |
| 16 | 57.8 | 75.2 | 86.7 | 64.1 | **18.4** |
| 20 | 57.3 | 74.0 | 85.9 | 64.9 | 18.3 |

Table 9: **Inference sliding window size**. The model is trained with stride 8 and sliding window size $T = 16$. The same window size gives the best results.

| Window size | DAVIS First | | | BADJA | |
|---|---|---|---|---|---|
| | AJ | $< \delta_{avg}^x$ | OA | seg-based | $< \delta^{3px}$ |
| 8 | 45.9 | 64.5 | 83.2 | 60.4 | 13.9 |
| 12 | 50.7 | 67.9 | 84.8 | 64.4 | 15.2 |
| **16** | **51.2** | **68.5** | **85.2** | **65.5** | 15.4 |
| 20 | 50.9 | **68.5** | 84.7 | 63.0 | **16.4** |
| 24 | 50.2 | 67.8 | 83.9 | 60.2 | 16.1 |

CoTracker has fewer parameters than PIPs (24M vs 29M) but more than RAFT (5M).

# D   Broader societal impact

Motion estimation, whether in the guise of point tracking or optical flow, is a fundamental, low-level computer vision task. Many other tasks in computer vision build on it, from 3D reconstruction to video object segmentation. Ultimately, motion estimation algorithms are an important component of a very large number of applications of computer vision in many different areas. Point tracking has no *direct* societal impact; however, positive or negative effects on society can materialise through its use in other algorithms, depending on the final application.