
TAPVid-3D: A Benchmark for Tracking Any Point in 3D

Skanda Koppula^{1,2*}, Ignacio Rocco^{1*}, Yi Yang¹, Joe Heyward¹,
João Carreira¹, Andrew Zisserman^{1,3}, Gabriel Brostow², Carl Doersch¹
¹Google DeepMind ²University College London ³University of Oxford

Abstract

We introduce a new benchmark, *TAPVid-3D*, for evaluating the task of long-range Tracking Any Point in 3D (TAP-3D). While point tracking in two dimensions (TAP) has many benchmarks measuring performance on real-world videos, such as TAPVid-DAVIS, three-dimensional point tracking has none. To this end, leveraging existing footage, we build a new benchmark for 3D point tracking featuring 4,000+ real-world videos, composed of three different data sources spanning a variety of object types, motion patterns, and indoor and outdoor environments. To measure performance on the TAP-3D task, we formulate a collection of metrics that extend the Jaccard-based metric used in TAP to handle the complexities of ambiguous depth scales across models, occlusions, and multi-track spatio-temporal smoothness. We manually verify a large sample of trajectories to ensure correct video annotations, and assess the current state of the TAP-3D task by constructing competitive baselines using existing tracking models. We anticipate this benchmark will serve as a guidepost to improve our ability to understand precise 3D motion and surface deformation from monocular video. Code for dataset download, generation, and model evaluation is available at <https://tapvid3d.github.io/>.

1 Introduction

For robots, humans, and other agents to effectively interact with the physical 3D world, it is necessary to understand a scene’s structure and dynamics. This is a key ingredient to any general embodied intelligence: the ability to learn a world model to understand and predict the structure and motion of arbitrary scenes. It is attractive to leverage the vast amounts of monocular video data that is available cheaply, and use such signals to understand the geometry and 3D motion in real-world footage. But how well can current perception algorithms actually do this?

The field has seen many efforts to measure 3D and motion understanding from videos, each contributing a part of the overall goal. For example, monocular depth estimation is a widely recognized task [7, 36, 63]. However, success in depth estimation alone does not reveal whether the model understands how surfaces move from one frame to the next, and may not respect temporal continuity. For instance, a basketball spinning on its axis is often not visible in a sequence of depth maps. On the other end of the generality spectrum, 3D pose tracking, e.g. for rigid objects [14, 34] and people [12, 49], evaluates precise motion, but requires a known 3D model of each object with articulations. Building parametric models and pose estimation methods for all animal classes is infeasible, much less for all the objects that a robot might encounter in, for example, an unseen, busy construction site.

An alternative and more general approach for dynamic scene understanding observes that the world consists of particles, each of which individually follows a 3D trajectory through space and time [58, 62]. Measuring the motion of these points provides a way to measure 3D motion *without requiring any*

*Equal contribution. Corresponding author: skandak@google.com

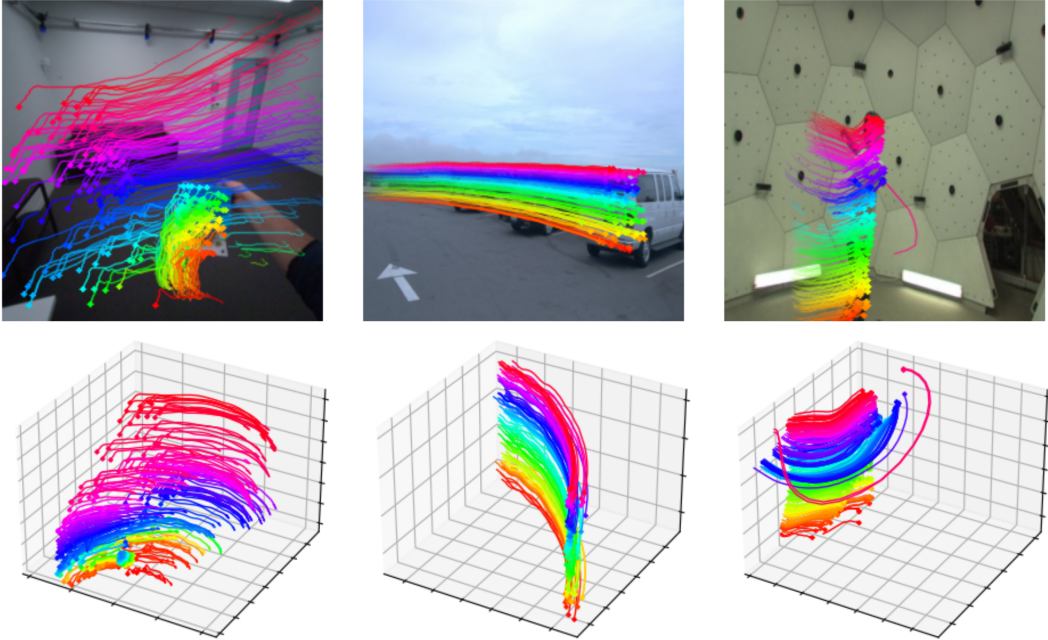


Figure 1: Random samples from TAPVid-3D: on the top row, we visualize the point trajectories projected into the 2D video frame; on the bottom row, we visualize the metric 3D point trajectories. From left to right, we show one example from each constituent data source: ADT, DriveTrack and Panoptic Studio.

3D shapes to be known a priori. To this end, in TAPVid-3D, we focus on providing the community with a benchmark consisting of real world videos and three-dimensional point tracking annotations, spanning a wide variety of objects, scenes, and motion patterns.

Prior work on 2D understanding has shown that this kind of point-wise motion can be extremely valuable: both optical flow and long term two-dimensional point tracking (TAP) tasks have been applied to video editing [66], controllable video generation [59], robotic manipulation [54, 60], and more [8, 17]. TAP is an occlusion-aware, temporal extension of optical flow, which itself has a 3D extension called *scene flow* [37, 39, 40]. While useful, scene flow suffers from the same challenges as optical flow, namely, that it captures instantaneous motion and does not evaluate correct, occlusion-aware association of pixels over long sequences. A new three-dimensional TAP benchmark would provide a way to measure progress on many of these tasks: our TAPVid-3D benchmark targets evaluating general motion understanding, for models performing both dense and sparse particle tracking, in two and three dimensions.

Unfortunately, all currently used evaluations for TAP on real-world videos assess only 2D tracking ability (e.g. the TAP-Vid suite [9], BADJA [4], CroHD [52], and JHMDB [18]), and cannot evaluate the performance of 3D point tracking due to lack of access to the ground-truth metric position trajectories. While evaluations based on synthetic environments, like Kubric [13], RGB-Stacking [9], and Point Odyssey [70], could potentially provide 3D point tracking annotations, these introduce a significant domain gap with real-world scenes and are therefore less representative of model performance on real-world tasks.

Many applications stand to benefit from direct evaluation of three-dimensional point tracking capabilities and improvements to such models. For example, robotic manipulation tasks are likely to be easier with accurate 3D motion estimates, to understand the changing relative world position of the gripper, any objects, and the background. Video generation models would be more useful if creators were able to condition on exact motion tracks describing the 3D movement of both objects and the camera, as a director would do on a stage. Standard scene understanding tasks like velocity estimation, motion prediction, and object parts segmentation are simpler given the 3D motion tracks of individual points. Many visual odometry, mapping, and structure from motion pipelines rely on accurate 3D correspondences; with the ability to track 3D point correspondences from *any* pixel, such

pipelines could be made more robust and accurate, even with many moving objects. Overall, the task of three dimensional point tracking provides a strict superset of information as compared to its 2D counterpart, is likely to be more useful in downstream applications, and provides a greater test of physical world motion understanding.

To this end, we introduce TAPVid-3D: a *real-world* benchmark for evaluating the Tracking Any Point in 3D (TAP-3D) task. The benchmark contributes: (1) a unification of three distinct real-world video data sources, with pipelines to generate, standardize, and validate consistent ground-truth (x, y, z) 3D trajectories and occlusion information, (2) formalization of the TAP-3D task, with new metrics to measure accuracy of 3D track estimation, and (3) an assessment of the current state of TAP-3D, formed from the first real-world video evaluations of the nascent set of early 3D TAP models. Code for dataset download, generation, and model evaluation is available at <https://tapvid3d.github.io/>.

2 Related work

Tracking Any Point. In recent years, long-range tracking of local image points has been formalized as the Tracking Any Point task (TAP) and evaluated using the, now standard, TAPVid benchmark [9], among others. Success on TAP consists of tracking the 2D trajectory of any given (x, y) *query point*, defined at a particular frame t , throughout the rest of a given video. By definition, the query point is considered visible at the query frame, and associated to the material point of the scene which is observed at (x, y, t) . However, this material point may become occluded or go out of the image boundaries. To handle this, TAP models also must predict a binary visibility flag v for each timestamp of the video. The tracked (x, y) positions and visibility estimates are scored jointly using an all-encompassing average Jaccard metric. Most current TAP models [9–11, 16, 21] are limited to tracking in 2D pixel space. Recently, some works have started exploring the extension of the TAP problem to 3D (TAP-3D) [57, 62]. However, most TAP benchmarks containing real-world videos, such as TAPVid-DAVIS [9], Perception Test [43], CroHD [52] and BADJA [4], don’t have 3D annotations, and therefore evaluations are still performed on the 2D tracking task. Concurrent to our work, Wang et al. [57] introduced both a synthetic (LSFOdyssey) and a real benchmark (LSFDiving) for evaluating TAP-3D. However, their real benchmark is limited to the driving domain and only contains 180 test clips with 40 frames each. Our proposed TAPVid-3D benchmark is substantially larger and more diverse, containing 4000+ clips with durations between 25 and 300 frames.

Scene flow estimation. The scene flow estimation problem, introduced by Vedula et al. [55], seeks to obtain a dense, instantaneous, 3D motion field of a 3D scene, analogously to the way optical flow estimates a dense 2D motion field across consecutive frame pairs of a 2D video. The TAP-3D task is related to the scene flow problem in a similar way in which the TAP task is related to the optical flow problem. While scene flow seeks to obtain dense instantaneous motion estimation, TAP-3D seeks to obtain longer-range tracking, spanning tens or hundreds of frames. Furthermore, TAP-3D does not seek to produce dense fields of tracks, but is rather interested in tracking a sparse set of query points, which is more computationally tractable. From work in TAP, we have observed that having motion representations with longer temporal context is useful for downstream tasks such as robotic manipulation, while having sparser spatial coverage is sufficient for many tracking applications.

Pose estimation. Closely related to point tracking is pose estimation and keypoint tracking. Many benchmarks have been proposed for 2D and 3D pose estimation [15, 26, 56, 68]. 3D pose estimation tasks and benchmarks largely focus on specific categories of moving objects, and for objects that are articulated: e.g. humans [19, 56], hands [53], animals [38, 67], and even jointed furniture [31]. For general motion and scene understanding, we aim to learn motion estimation across any object or scene pixel, expanding the generality of the task.

Static scene reconstruction. Static scene reconstruction, a fundamental problem in computer vision, has been advanced through techniques like Structure-from-Motion (SfM) and monocular depth estimation. Significant contributions include COLMAP [45] for state-of-the-art reconstructions and MVSNet [65], which enhances multi-view stereo depth estimation with deep learning. These studies collectively advance robust and precise static scene reconstruction. Evaluation of the local features and depth are crucial for static scene reconstruction methods. [46] provided a comparative evaluation of hand-crafted and learned features, while MegaDepth [27] improved single-view depth prediction using large scale multi-view Internet photo collections. Despite significant progress, static

Benchmark Type	Long Term Continuity	3D	Pixel Level Occlusion	Pixel Level Motion	Non-rigid Surfaces	No Prior 3D Model
Monodepth [6, 27, 50]	✗	✓	✗	✗	✓	✓
2D Point Tracking (TAP) [9]	✓	✗	✓	✓	✓	✓
Scene flow [23, 61]	✗	✓	✗	✓	✓	✓
3D Pose Tracking [33, 56]	✓	✓	✗	✓	✗	✗
3D Object Box Tracking [30, 51]	✓	✓	✗	✗	✗	✓
3D Point Tracking (TAP-3D)	✓	✓	✓	✓	✓	✓

Table 1: The proposed TAPVid-3D benchmark provides a unique set of characteristics, not covered in previous tasks or benchmarks. It extends the temporal continuity, occlusion modeling, and motion estimation capabilities of TAP benchmarks into 3D. Each row describes aspects tested in each task.

scene reconstruction struggles with dynamic environments, highlighting the need for dynamic scene methods.

Dynamic scene reconstruction. 3D reconstruction of dynamic scenes and objects is a widely studied problem in computer vision. Over the years, several methods have proposed solutions to this problem, starting with Non-rigid Structure-from-Motion methods (NRSfM) [1, 5]. While these methods have shown some success modelling simple motions like facial expressions and skeletal motions, they fail to generalize to arbitrary motions. More recently, deep-learning based methods have been used to perform 3D reconstruction of dynamic scenes. One line of works exploit Monocular Depth Estimation models [44] and performs test-time optimization on each given video to obtain smoother reconstructions, under the assumption that the frame rate of the camera is high with respect to the speed of the depicted object (quasi-static scene assumption) [25, 36, 69]. While these methods are more general than the classic NRSfM counterparts, they still require costly test-time optimization and fail to model rapid motions. Other lines of work attempt to fit a neural-scene representation to each video, such as neural-radiance fields [28, 29] or 3D Gaussian Splatting [64]. However, these methods require a costly per-video optimization, and typically make smoothness and local-rigidity assumptions about the motion of the points in the scene. We believe the development of TAP-3D methods should significantly help for the problem of dynamic scene reconstruction, as these models can run in a feed-forward manner, without requiring test-time optimization, and don’t need any explicit motion prior assumptions as they can learn these from data.

Table 1 summarizes the focus areas of measurement for common scene understanding benchmarks. On the bottom row is the proposed TAP-3D task, and corresponding benchmark, TAPVid-3D, which brings to the table a more complete test of dynamic scene understanding in one simple evaluation.

3 TAPVid-3D

We build a real-world benchmark for evaluating Tracking Any Point in 3D (TAP-3D) models. To do this, we leveraged three publicly available datasets: (i) Aria Digital Twin [42], (ii) DriveTrack [2, 51] and (iii) Panoptic Studio [20]. These data sources span different application domains, environments, and video characteristics, deriving ground truth tracking trajectories from different sensor types. For instance, Aria Digital Twin is a dataset of egocentric video, and is more close to bimanual robotic manipulation problems, where the camera is robot mounted and sees the actions from a first person view. DriveTrack features footage captured from a Waymo car navigating outdoor scenes, akin to tasks in robotic navigation and outdoor, rigid-body scene understanding tasks. Finally, Panoptic Studio captures third-person views of people performing diverse actions within an instrumented dome. It presents complex human movement which aligns more closely with NRSfM-adjacent tasks. We believe that, combined, these data sources present a diverse and comprehensive benchmark of TAP-3D capabilities for many potential downstream tasks. We describe our pipeline to extract ground truth metric 3D point trajectories from each source in the next sections, with samples in Figure 1.

Table 2 shows the summary statistics across the entire dataset, and for the dataset subdivisions corresponding to each constituent data source. As there are comparatively a large number of videos (for reference, the commonly used TAPVid-DAVIS has 30 videos, whereas TAPVid-3D has two orders of magnitude more), we release two splits: a `minival` split, with 50 videos from each data

source, and a `full_test` split, containing all videos in the benchmark, without the `minimal` videos. The `minimal` is intentionally lightweight, and likely most useful for online evaluation during training.

3.1 Aria Digital Twin

This split employs real videos captured with the Aria glasses [41] inside different recording studios, which mimic household environments. Accurate digital replicas of these studios, created in 3D modeling software, are used to obtain pseudo-ground truth annotations of the original footage. This includes annotations such as segmentation masks, 3D object bounding boxes and depth maps. We leverage these annotations to extract 3D trajectories from given 3D query points. In particular, given a video $V = \{I_t\}_{t=1, \dots, T}$ with T frames and $W \times H$ spatial resolution, with corresponding object segmentation masks $\{S_t\} \subset \mathbb{Z}^{W \times H}$, and depth maps $\{D_t\} \subset \mathbb{R}^{W \times H}$, extrinsic world-to-camera poses* $\{(P_{cam}^w)_t\} \subset \mathbb{R}^{3 \times 4}$, camera intrinsics $K \in \mathbb{R}^{3 \times 3}$, and a query point $q = (x_q, y_q, t_q)$, we first extract the query point’s 3D position Q_{cam}^\dagger in the camera coordinate frame, with

$$(Q_{cam})_{t_q} = K^{-1}(x_q, y_q, 1)^T \cdot D_{t_q}(x_q, y_q). \quad (1)$$

Additionally, we obtain the object ID of the query point from the segmentation mask $q_{id} = S_{t_q}(x_q, y_q)$, and use it to retrieve the 3D object pose of the query object $(P_{obj}^w)_{t_q}$, which converts points from world coordinate frame to the object coordinate frame. This allows us to compute the query point position in the object’s frame of reference as

$$Q_{obj} = (P_{obj}^w)_{t_q} (P_w^{cam})_{t_q} (Q_{cam})_{t_q}. \quad (2)$$

In this way, we *fix* the query point to the corresponding object, and then obtain its track across the whole video by leveraging the object’s pose annotation. For any timestamp t , the position of the query point can be then obtained by

$$(Q_{cam})_t = (P_{cam}^w)_t (P_w^{obj})_t Q_{obj}. \quad (3)$$

To compute the visibility v of the query point at any time, we first employ a pretrained semantic segmentation model to obtain the semantic masks of the operator’s hands $\{H_t\}$, as these are not modelled in the digital replica. Then, we compute the visibility v by verifying that the query point’s depth is close to the observed depth, and it does not lie on the hands segmentation mask H , so

$$v_t = \mathbf{1}(|Z(Q_{cam}) - D_t(u, v)| < \delta) \cdot (1 - H_t(u, v)), \quad (4)$$

where $(u, v) = \Pi_K((Q_{cam})_t)$ is the projection of the query point $(Q_{cam})_t$ to the image plane according to the given camera intrinsics K , and $Z((x, y, z)) = z$ is the function that extracts the z-component of a 3D point. This approach allows us to compute the 3D trajectory $\{(Q_{cam})_t\}$ and visibility flag $\{v_t\}$ of the query point across the whole video.

3.2 DriveTrack

The DriveTrack split is based on videos from the Waymo Open dataset [51], and the 2D point trajectory pipeline used in DriveTrack [2]. In particular, each frame I_t in a video sequence V has a corresponding, time-synchronized point cloud $\{C_t\}$ from the Waymo car’s LIDAR. The subset of points that correspond to a randomly selected, single, tracked vehicle $(Q_{cam})_{t_s}$ are subselected from the entire point cloud $(Q_{cam})_{t_s} \subset C_{t_s}$ at a certain sampling time t_s using a manually-annotated 3D bounding box around the chosen object. These object-specific points are then tracked across the whole video using: (i) a vehicle rigidity assumption, and, (ii) the pose and position of the object’s annotated 3D bounding box through the entire video sequence. Specifically, the tracked points in object coordinate frame Q_{obj} are first computed using (2), and then the trajectories in camera coordinate frame $\{(Q_{cam})_t\}$ are obtained using (3).

*We use the notation P_b^a to represent the SE(3) transform from coordinate frame a to frame b .

†We use the notation Q_{cam} and Q_{obj} to denote the position of the 3D point Q in the camera and object coordinate frames, respectively.

Visibility flag is estimated by first computing the dense depth map D_t of the corresponding camera video frames through interpolation of sparse LIDAR values as in [2]. This is compared to the point’s depth computed from the 3D point trajectory given by $(Q_{cam})_t$, as in the first term of (4). If the point depth (distance from the camera center to the query point) is greater than the depth provided by the depth map (with a 5% relative threshold margin), it is marked as not visible.

Finally, to determine the 2D query points $q = (x_q, y_q, t_q)$ we sample t_q uniformly among the visible timestamps v_t , and then obtain $(x_q, y_q) = \Pi_K((Q_{cam})_{t_q})$.

3.3 Panoptic Studio

The original Panoptic Studio dataset [20] consists in video sequences captured inside a recording dome using stationary cameras, and depicting different actors performing various actions such as passing a ball or swinging a tennis racket, featuring complex non-rigid motions. To obtain 3D trajectories, we leverage the pretrained dynamic 3D reconstructions from Luiten et al. [35]. These have been obtained by first performing a rigid 3D reconstruction through 3D Gaussian Splatting [22], fitting a set of 3D Gaussians $\{(\mu_i, \Sigma_i)_{t_0}\}_{i=1, \dots, N}$ to the first timestamp t_0 of each sequence using the multiple cameras available in the dome. Then, these Gaussians are displaced and rotated in a as-rigid-as-possible manner to model the motion occurring in the subsequent frames of the video. For more details, please refer to [35]. Using these pretrained splatting models, we render pseudo-ground-truth depth maps $\{D_t\}$ for each sequence. Then, given a query point $q = (x_q, y_q, t_q)$, we unproject the point to 3D following (1), obtaining $(Q_{cam})_{t_q}$, and retrieve the index i^* of the closest Gaussian center at that time using the poses $\{(P_{cam}^w)_t\}$, so that $i^* = \operatorname{argmin}_{i=1, \dots, N} \|(Q_{cam})_{t_q} - (P_{cam}^w)_{t_q}(\mu_i)_{t_q}\|_2$.

Note that due to the distance between $(Q_{cam})_{t_q}$ and $(P_{cam}^w)_{t_q}(\mu_{i^*})_{t_q}$, their projections onto the image plane will not match exactly. To account for this difference, we adjust the query point’s 2D position as $(x_q, y_q) = \Pi_K(P_{cam}^w)_{t_q}(\mu_{i^*})_{t_q}$. Then, the query point’s 3D track come by following the motion of the i^* Gaussian center, so $(Q_{cam})_t \equiv (P_{cam}^w)_t(\mu_{i^*})_t$.

Visibility $\{v_t\}$ is estimated as in (4) (omitting the second term), by comparing the depth of our 3D query point with the observed depth from D_t . We only track points across the foreground deforming characters, as tracking background points is trivial given that the cameras in this dataset are stationary.

3.4 Data Cleanup and Validation

While the aforementioned procedures generally produce high quality trajectories, small amounts of noise from the underlying dataset sources can cause issues in a small fraction of sequences. These minor inaccuracies, for example, can be caused by small misalignment between Aria Digital Twin synthetic annotations and the real world video, LIDAR sensor noise, insufficiently constrained Gaussian splats, or numerical error. We minimized these errors through automated methods, and then manually checked a sampling of the videos to ensure accuracy.

Firstly, since trajectories are descriptors of surface motion, their motion should be localized to their associated object. We use instance segmentation models to generate object masks for each frame [24], filtering out errant trajectories that exceed these boundaries when not occluded. In DriveTrack specifically, the tracked bounding box is an approximation to the true object mesh, but such errors are fixed with tight segmentation masks (trimming 2-3% of initial trajectories).

Secondly, we observed that some trajectories have a ‘flickering’ visibility flag. This is not unique to TAPVid-3D, as we notice this in the widely used Kubric [13] and DriveTrack [2]. However, to mitigate this in our dataset, we oversample trajectories in the annotation generation pipeline, and filter trajectories whose visibility state changes more times than 10% of the total number of video frames. More details can be found in the appendix.

3.5 Metrics

To accompany the TAPVid-3D dataset, we adopt and extend the metrics used in TAP [9] to the 3D tracking domain. These metrics measure the quality of the predicted 3D point trajectories (APD), the ability to predict point visibility (OA), or both simultaneously (AJ).

Dataset split	#clips (minival)	#trajs per clip	#frames per clip	#videos	#scenes	resolution	fps
Aria Digital Twin	1956 (50)	1024	300	215	2	512 × 512	30
DriveTrack	2457 (50)	256	25 – 300	2457	252	1920 × 1280	10
Panoptic Studio	156 (50)	50	150	156	1	640 × 360	30
TAPVid-3D	4569 (150)	50 – 1024	25 – 300	2828	255	Multiple	10/30

Table 2: Overview statistics of the TAPVid-3D benchmark dataset, for all three constituent splits and in total. Clips in the benchmark are temporally sampled from their original video.

The **APD** ($< \delta_{avg}^x$) metric measures the average percent of points within δ error. If \hat{P}_t^i is the i 'th point prediction at time t , P_t^i is the corresponding ground truth 3D point, and v_t^i is the ground-truth point visibility flag, then:

$$\text{APD}_{3D} \equiv \frac{1}{V} \sum_{i,t} v_t^i \cdot \mathbf{1}(\|\hat{P}_t^i - P_t^i\| < \delta_{3D}(P_t^i)), \quad (5)$$

where $\mathbf{1}(\cdot)$ is the indicator function, $\|\cdot\|$ is the Euclidean norm, $V = \sum_{i,t} v_t^i$ is the total number of visibles, and $\delta_{3D}(P_t^i)$ is the threshold value.

The value of this threshold is *relative to ground-truth depth*; and is defined by *unprojecting* a pixel threshold $\delta_{2D} \in \{1, 2, 4, 8, 16\}$ to 3D space using the camera intrinsic parameters: $\delta_{3D}(P_t^i) = Z(P_t^i) \cdot \delta_{2D}/f$, where f is the camera focal length. We argue that points that are closer to the camera are of higher importance than those that are far away, which is enforced by the definition of δ_{3D} [‡]. Note that by using this definition, the APD_{3D} metric defined in (5) is numerically equivalent to the APD_{2D} metric from [9] when the point estimations \hat{P}_t^i all have correct depths.

In addition, it is important to distinguish between occluded and visible points, because downstream algorithms often want to rely exclusively on predictions which are based on visual evidence, rather than on the guesses that might be very wrong. To this end, we adopt the occlusion accuracy metric (**OA**) from [9], which computes the fraction of points on each trajectory where $\hat{v}_t^i = v_t^i$, where \hat{v}_t^i is the model's (binary) visibility prediction.

Finally, we define **3D-AJ**, 3D Average Jaccard, following TAP, which combines OA and APD_{3D} . The AJ metric calculates the number of *true positives* (number of points within the δ_{3D} threshold, predicted correctly to be visible), divided by the sum of *true positives* and *false positives* (predicted visible, but are occluded or farther than the threshold) and *false negatives* (visible points, predicted occluded or predicted to exceed the threshold). Mathematically, it is defined as:

$$\text{AJ}_{3D} \equiv \frac{\sum_{i,t} v_t^i \hat{v}_t^i \alpha_t^i}{\sum_{i,t} v_t^i + \sum_{i,t} ((1 - v_t^i) \hat{v}_t^i) + \sum_{i,t} (v_t^i \hat{v}_t^i (1 - \alpha_t^i))}, \quad (6)$$

where $\alpha_t^i = \mathbf{1}(\|\hat{P}_t^i - P_t^i\| < \delta_{3D}(P_t^i))$ indicates whether the point prediction is below the distance threshold. Note the relationship between this and prior metrics: if the depth estimates for a given point are perfect, then this metric reduces to 2D-AJ, as there will be no depth errors. It is also related to a common Monodepth's metric $\delta < 1.25^k$, which also places a hard threshold on the *relative* depth of the estimated point w.r.t. the ground truth; if the video is a single frame (occlusion-free) and the query points are dense, then there should be no 2D errors, and our metric will behave like a Monodepth metric. For general sequences, however, the algorithm must output correct tracking *and* correct depth—i.e., a full 3D trajectory—in order to be considered correct by our metric.

Finally, one additional complication is *scale ambiguity*. As in monocular depth literature [44], we globally re-scale predictions to match ground truth, before computing the metrics. We do this by multiplying predictions \hat{P}_t^i by the median of the depth ratios $\|P_t^i\|/\|\hat{P}_t^i\|$ over all points and frames, which we call *global median* rescaling. Note, however, that some algorithms may be more adept at estimating the relative depth of *individual* points: algorithms trained in simulation, for instance, may be able to estimate the depth change for a single point just by analyzing the frequencies. With a slight

[‡]A similar depth-adaptive threshold approach is used in Monodepth literature [44].

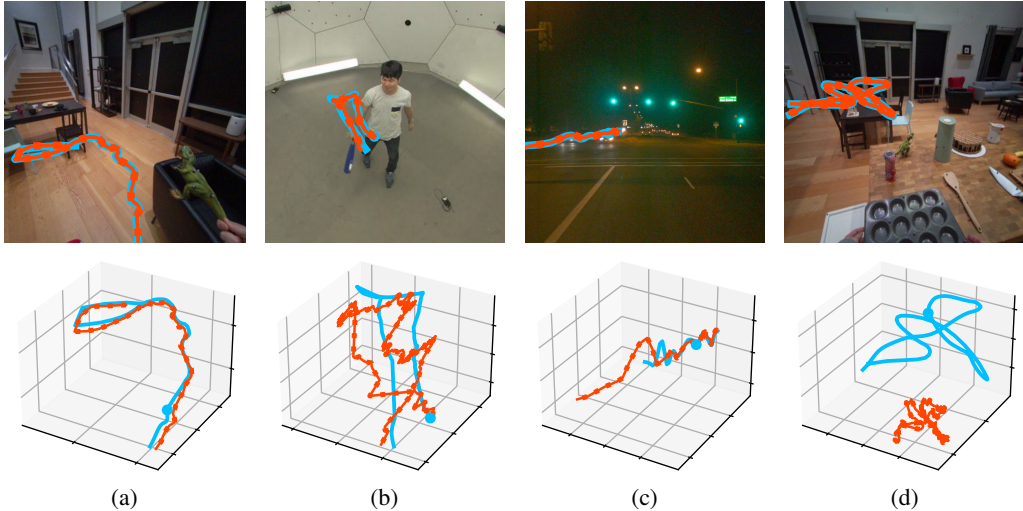


Figure 2: Illustrative TAP-3D results of BootsTAPIR + ZoeDepth. We compare the ground-truth 2D and 3D tracks (blue solid) with the predicted tracks (red dotted). (a) Accurate tracking. (b) Noisy depth estimations result in a noisy 3D track. (c) Inconsistent depth scales across time (scale drift). (d) Inconsistent depth scales across space don’t allow a single global scale factor to properly fit all tracks.

change to the rescaling, we can accommodate such methods even when they don’t produce consistent scale between points. Therefore, we define the *per-trajectory* rescaling, which rescales each track P^i separately multiplying by $\|P_{t_q}^i\|/\|\hat{P}_{t_q}^i\|$, where t_q is the query timestamp. Finally, we explore a hybrid approach (which considers scaling in local neighborhoods to better account for object-object interactions) in the supplementary material.

4 Baselines on TAPVid-3D

We construct one set of baselines by combining state-of-art 2D point trackers and monocular depth estimators. In particular, we use for 2D tracking, several state-of-the-art models such as CoTracker [21], BootsTAPIR model [11], and TAPIR [10]; and for depth regression, both the monocular depth estimation model ZoeDepth [3] and the SfM pipeline COLMAP [47, 48]. To convert the frame pixel space predictions into metric x, y -position coordinates, we unproject using the camera intrinsics and the z -estimate provided by ZoeDepth. In the case of COLMAP, we import the 2D trajectories produced by the TAP methods before running the SfM reconstruction pipeline.

We include the recently released SpatialTracker [62], one of the first models designed specifically for 3D point tracking. We also provide results for a new variant of TAPIR [10] that we built for 3D tracking, that we call TAPIR-3D (more details in the appendix). Uniquely, this is the only baseline that is trained only on synthetic videos. Finally, we include a static baseline, which projects the pixel query point into 3D, and assumes no motion of that 3D point. Results for all these baseline methods are shown in Table 3, for the `full_eval` split, and three major methods on `minival` split. Results in Table 3 uses median depth scaling; results for other scaling options are included in the appendix.

Firstly, comparing Tables 3 and 4, we find that the 3D tracking performance of our baselines are significantly lower compared to their 2D tracking abilities. We show examples in Figure 2, illustrating common failure modes regressing 3D trajectories, noting that while the 2D trajectories look accurate, their understanding of total scene geometry and correct 3D motion is poor.

Secondly, we find that having the three different data sources provides a video diversity to the benchmark that enables better assessment of the strengths and weakness of the video model under test. For example, ZoeDepth-based methods perform comparatively well on the PStudio subset, but significantly underperforms COLMAP on the DriveTrack subset (where monocular depthing might be harder, with the outdoor complex scenes, in which scale-consistent depthing may be difficult). On the other hand, COLMAP struggles with the PStudio subset, where nearly all tracks and majority of the scene is occupied by the main moving objects. This underperformance is likely because COLMAP fails to reconstruct moving content. SpatialTracker uses ZoeDepth under the hood, so its pattern of results across the subsets are similar to 2D Tracking + ZoeDepth. Our TAPIR-3D experimental model,

Baseline	Aria			DriveTrack			PStudio			Average		
	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow
Static Baseline	4.9	10.2	55.4	3.9	6.5	80.8	5.9	11.5	75.8	4.9	9.4	70.7
TAPIR + COLMAP	7.1	11.9	72.6	8.9	14.7	80.4	6.1	10.7	75.2	7.4	12.4	76.1
CoTracker + COLMAP	8.0	12.3	78.6	11.7	19.1	81.7	8.1	13.5	77.2	9.3	15.0	79.1
BootsTAPIR + COLMAP	9.1	14.5	78.6	11.8	18.6	83.8	6.9	11.6	81.8	9.3	14.9	81.4
TAPIR + ZoeDepth	9.0	14.3	79.7	5.2	8.8	81.6	10.7	18.2	78.7	8.3	13.8	80.0
CoTracker + ZoeDepth	10.0	15.9	87.8	5.0	9.1	82.6	11.2	19.4	80.0	8.7	14.8	83.4
BootsTAPIR + ZoeDepth	9.9	16.3	86.5	5.4	9.2	85.3	11.3	19.0	82.7	8.8	14.8	84.8
TAPIR-3D	2.5	4.8	86.0	3.2	5.9	83.3	3.6	7.0	78.9	3.1	5.9	82.8
SpatialTracker [62]	9.9	16.1	89.0	6.2	11.1	83.7	10.9	19.2	78.6	9.0	15.5	83.7
BootsTAPIR + COLMAP*	7.3	11.5	76.3	9.3	15.1	83.5	6.2	10.6	78.7	7.6	12.4	79.5
BootsTAPIR + ZoeDepth*	8.6	14.5	86.9	5.1	8.7	83.5	10.2	17.7	82.0	8.0	13.6	84.1
SpatialTracker [62]*	9.2	15.1	89.9	5.8	10.2	82.0	9.8	17.7	78.4	8.3	14.3	83.4

Table 3: **Using median depth scaling.** We compare the performance of 2D-TAP models [10, 11, 21] combined with ZoeDepth [3] and COLMAP [47]. We also include SpatialTracker [62], and a static point baseline. We report the proposed 3D-AJ as well as the APD and OA metrics. The top nine rows are evaluated on the `full_eval` split, while the bottom three rows (marked with *) indicate evaluation on the `minival` split.

	Aria	DriveTrack	PStudio	Total		
	2D-AJ \uparrow	2D-AJ \uparrow	2D-AJ \uparrow	2D-AJ \uparrow	APD \uparrow	OA \uparrow
TAPIR [10]	48.6	57.2	48.7	53.2	67.4	80.5
CoTracker [21]	54.2	59.8	51.0	57.2	74.2	84.5
BootsTAPIR [11]	54.7	62.9	52.4	59.1	74.7	85.6

Table 4: Evaluating the 2D point tracking performance of our baseline models on TAPVid-3D data, by projecting the ground truth 3D trajectories onto the 2D frame.

like TAPIR, outputs trajectories independently, so scale consistency across the entire scene/across trajectories is lacking; it is thus correspondingly penalized, and exhibits poor 3D-AJ. The varying strengths and weaknesses of these different methods indicate that the best performance could be achieved by integrating elements from all three.

Limitations and Responsible Usage. The three data sources in our benchmark do not cover all possible domains where users may want to infer dynamics, and automatic annotations may be imperfect, e.g. where the underlying sensor readings have noise (despite our filtering). Furthermore, we inherit some limitations from TAP and monocular depth: we only evaluate tracking for solid, opaque objects. From an ethical perspective, Aria and Panoptic were collected in lab settings with consenting participants, while DriveTrack’s Waymo videos come from public roads in 6 US cities. This paper may inherit biases from these datasets, e.g., participants are lab researchers or the populations of those 6 US cities. This dataset is not intended for training, but care should be taken in training data to avoid biases. As a benchmark, the broader impacts are similar to those in prior vision and tracking works: there may be very long term applications to activity recognition and surveillance.

5 Conclusion

We introduce TAPVid-3D, a new benchmark for evaluating the nascent Tracking Any Point in 3D (TAP-3D) task. We contributed (1) a pipeline to annotate three distinct real-world video datasets to produce 3D correspondence annotations per video, (2) the first metrics for the TAP-3D task, which measure multiple axis of accuracy of 3D track estimates, and (3) an assessment of the current state of TAP-3D, by evaluating commonly used tracking models. We believe this benchmark will accelerate research on the TAP-3D task, allowing the development of models with greater dynamic scene understanding from monocular video.

References

- [1] Ijaz Akhter, Yaser Sheikh, Sohaib Khan, and Takeo Kanade. Nonrigid structure from motion in trajectory space. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *NeurIPS*, 2008.

- [2] Arjun Balasingam, Joseph Chandler, Chenning Li, Zhoutong Zhang, and Hari Balakrishnan. Drivetrack: A benchmark for long-range point tracking in real-world videos. *arXiv preprint arXiv:2312.09523*, 2023.
- [3] Shariq Farooq Bhat, Reiner Birkel, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023.
- [4] 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*, 2019.
- [5] Christoph Bregler, Aaron Hertzmann, and Henning Biermann. Recovering non-rigid 3d shape from image streams. In *Proc. CVPR*, 2000.
- [6] Daniel J Butler, Jonas Wulff, Garrett B Stanley, and Michael J Black. A naturalistic open source movie for optical flow evaluation. In *Computer Vision—ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7–13, 2012, Proceedings, Part VI 12*, pages 611–625. Springer, 2012.
- [7] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 8001–8008, 2019.
- [8] Haiyang Chao, Yu Gu, and Marcello Napolitano. A survey of optical flow techniques for robotics navigation applications. *Journal of Intelligent & Robotic Systems*, 73:361–372, 2014.
- [9] 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. *NeurIPS*, 2022.
- [10] 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. *arXiv preprint arXiv:2306.08637*, 2023.
- [11] Carl Doersch, Yi Yang, Dilara Gokay, Pauline Luc, Skanda Koppula, Ankush Gupta, Joseph Heyward, Ross Goroshin, João Carreira, and Andrew Zisserman. Bootstrap: Bootstrapped training for tracking-any-point. *arXiv preprint arXiv:2402.00847*, 2024.
- [12] Shubham Goel, Georgios Pavlakos, Jathushan Rajasegaran, Angjoo Kanazawa, and Jitendra Malik. Humans in 4d: Reconstructing and tracking humans with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14783–14794, 2023.
- [13] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J Fleet, Dan Gnanapragasam, Florian Golemo, Charles Herrmann, et al. Kubric: A scalable dataset generator. In *Proc. CVPR*, 2022.
- [14] Yang Hai, Rui Song, Jiaojiao Li, Mathieu Salzmann, and Yinlin Hu. Rigidity-aware detection for 6d object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8927–8936, 2023.
- [15] Albert Haque, Boya Peng, Zelun Luo, Alexandre Alahi, Serena Yeung, and Li Fei-Fei. Towards viewpoint invariant 3d human pose estimation. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, pages 160–177. Springer, 2016.
- [16] Adam W Harley, Zhaoyuan Fang, and Katerina Fragkiadaki. Particle video revisited: Tracking through occlusions using point trajectories. In *Proc. ECCV*, 2022.
- [17] Chong Huang and Kazuhito Koishida. Improved active speaker detection based on optical flow. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 950–951, 2020.

- [18] H. Jhuang, J. Gall, S. Zuffi, C. Schmid, and M. J. Black. Towards understanding action recognition. In *International Conf. on Computer Vision (ICCV)*, pages 3192–3199, 2013.
- [19] Sam Johnson and Mark Everingham. Clustered pose and nonlinear appearance models for human pose estimation. In *bmvc*, volume 2, page 5. Aberystwyth, UK, 2010.
- [20] Hanbyul Joo, Tomas Simon, Xulong Li, Hao Liu, Lei Tan, Lin Gui, Sean Banerjee, and Timothy Godisart. Panoptic studio: A massively multiview system for social interaction capture. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(1), 2019.
- [21] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. CoTracker: It is better to track together. *arXiv preprint arXiv:2307.07635*, 2023.
- [22] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 2023.
- [23] Ishan Khatri, Kyle Vedder, Neehar Peri, Deva Ramanan, and James Hays. I can’t believe it’s not scene flow! *arXiv preprint arXiv:2403.04739*, 2024.
- [24] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.
- [25] Johannes Kopf, Xuejian Rong, and Jia-Bin Huang. Robust consistent video depth estimation. In *Proc. CVPR*, 2021.
- [26] Jiefeng Li, Can Wang, Hao Zhu, Yihuan Mao, Hao-Shu Fang, and Cewu Lu. Crowdpose: Efficient crowded scenes pose estimation and a new benchmark. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10863–10872, 2019.
- [27] Zhengqi Li and Noah Snavely. Megadepth: Learning singleview depth prediction from internet photos. *ieee*. In *Proc. CVPR*, 2018.
- [28] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proc. CVPR*, 2021.
- [29] Zhengqi Li, Qianqian Wang, Forrester Cole, Richard Tucker, and Noah Snavely. Dynibar: Neural dynamic image-based rendering. In *Proc. CVPR*, 2023.
- [30] Yiyi Liao, Jun Xie, and Andreas Geiger. Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. *IEEE PAMI*, 45(3), 2022.
- [31] Joseph J Lim, Hamed Pirsiavash, and Antonio Torralba. Parsing ikea objects: Fine pose estimation. In *Proceedings of the IEEE international conference on computer vision*, pages 2992–2999, 2013.
- [32] Hao Liu, Yanni Ma, Qingyong Hu, and Yulan Guo. Centertube: Tracking multiple 3d objects with 4d tubelets in dynamic point clouds. *IEEE Transactions on Multimedia*, 2023.
- [33] Yunze Liu, Yun Liu, Che Jiang, Kangbo Lyu, Weikang Wan, Hao Shen, Boqiang Liang, Zhoujie Fu, He Wang, and Li Yi. Hoi4d: A 4d egocentric dataset for category-level human-object interaction. In *Proc. CVPR*, 2022.
- [34] Manolis Lourakis and Xenophon Zabulis. Model-based pose estimation for rigid objects. In *International conference on computer vision systems*, pages 83–92. Springer, 2013.
- [35] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. In *3DV*, 2024.
- [36] Xuan Luo, Jia-Bin Huang, Richard Szeliski, Kevin Matzen, and Johannes Kopf. Consistent video depth estimation. *ACM Transactions on Graphics*, 2020.

- [37] Wei-Chiu Ma, Shenlong Wang, Rui Hu, Yuwen Xiong, and Raquel Urtasun. Deep rigid instance scene flow. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3614–3622, 2019.
- [38] Alexander Mathis, Thomas Biasi, Steffen Schneider, Mert Yuksekgonul, Byron Rogers, Matthias Bethge, and Mackenzie W Mathis. Pretraining boosts out-of-domain robustness for pose estimation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 1859–1868, 2021.
- [39] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4040–4048, 2016.
- [40] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3061–3070, 2015.
- [41] Meta. Introducing Project Aria Glasses, from Meta, 2023. URL <https://www.projectaria.com/glasses/>.
- [42] Xiaqing Pan, Nicholas Charron, Yongqian Yang, Scott Peters, Thomas Whelan, Chen Kong, Omkar Parkhi, Richard Newcombe, and Yuheng Carl Ren. Aria digital twin: A new benchmark dataset for egocentric 3d machine perception. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20133–20143, 2023.
- [43] Viorica Patraucean, Lucas Smaira, Ankush Gupta, Adria Recasens, Larisa Markeeva, Dylan Banarse, Skanda Koppula, Mateusz Malinowski, Yi Yang, Carl Doersch, et al. Perception test: A diagnostic benchmark for multimodal video models. In *Perception Test*, 2024.
- [44] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE PAMI*, 2020.
- [45] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proc. CVPR*, 2016.
- [46] Johannes L Schonberger, Hans Hardmeier, Torsten Sattler, and Marc Pollefeys. Comparative evaluation of hand-crafted and learned local features. In *Proc. CVPR*, pages 1482–1491, 2017.
- [47] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [48] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016.
- [49] Cristian Sminchisescu. 3d human motion analysis in monocular video: techniques and challenges. *Human Motion: Understanding, Modelling, Capture, and Animation*, pages 185–211, 2008.
- [50] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 567–576, 2015.
- [51] Pei Sun, Henrik Kretschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2446–2454, 2020.
- [52] Ramana Sundararaman, Cedric De Almeida Braga, Eric Marchand, and Julien Pettre. Tracking pedestrian heads in dense crowd. In *Proc. CVPR*, 2021.
- [53] Jonathan Tompson, Murphy Stein, Yann Lecun, and Ken Perlin. Real-time continuous pose recovery of human hands using convolutional networks. *ACM Transactions on Graphics*, 33, August 2014.

- [54] Mel Vecerik, Carl Doersch, Yi Yang, Todor Davchev, Yusuf Aytar, Guangyao Zhou, Raia Hadsell, Lourdes Agapito, and Jon Scholz. RoboTAP: Tracking arbitrary points for few-shot visual imitation. In *Proc. Intl. Conf. on Robotics and Automation*, 2024.
- [55] Sundar Vedula, Simon Baker, Peter Rander, Robert Collins, and Takeo Kanade. Three-dimensional scene flow. In *Proc. ICCV*, 1999.
- [56] Timo Von Marcard, Roberto Henschel, Michael J Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *Proceedings of the European conference on computer vision (ECCV)*, pages 601–617, 2018.
- [57] Bo Wang, Jian Li, Yang Yu, Li Liu, Zhenping Sun, and Dewen Hu. SceneTracker: Long-term scene flow estimation network. *arXiv preprint arXiv:2403.19924*, 2024.
- [58] Qianqian Wang, Yen-Yu Chang, Ruojin Cai, Zhengqi Li, Bharath Hariharan, Aleksander Holynski, and Noah Snavely. Tracking everything everywhere all at once. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19795–19806, 2023.
- [59] Zhouxia Wang, Ziyang Yuan, Xintao Wang, Tianshui Chen, Menghan Xia, Ping Luo, and Ying Shan. Motionctrl: A unified and flexible motion controller for video generation. *arXiv preprint arXiv:2312.03641*, 2023.
- [60] Chuan Wen, Xingyu Lin, John So, Kai Chen, Qi Dou, Yang Gao, and Pieter Abbeel. Any-point trajectory modeling for policy learning. *arXiv preprint arXiv:2401.00025*, 2023.
- [61] Benjamin Wilson, William Qi, Tanmay Agarwal, John Lambert, Jagjeet Singh, Siddhesh Khandelwal, Bowen Pan, Ratnesh Kumar, Andrew Hartnett, Jhony Kaesemodel Pontes, et al. Argoverse 2: Next generation datasets for self-driving perception and forecasting. *arXiv preprint arXiv:2301.00493*, 2023.
- [62] Yuxi Xiao, Qianqian Wang, Shangzhan Zhang, Nan Xue, Sida Peng, Yujun Shen, and Xiaowei Zhou. Spatialtracker: Tracking any 2d pixels in 3d space. *arXiv preprint arXiv:2404.04319*, 2024.
- [63] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. *arXiv preprint arXiv:2401.10891*, 2024.
- [64] Zeyu Yang, Hongye Yang, Zijie Pan, Xiatian Zhu, and Li Zhang. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. *arXiv preprint arXiv:2310.10642*, 2023.
- [65] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proc. ECCV*, pages 767–783, 2018.
- [66] Emilie Yu, Kevin Blackburn-Matzen, Cuong Nguyen, Oliver Wang, Rubaiat Habib Kazi, and Adrien Bousseau. VideoDoodles: Hand-drawn animations on videos with scene-aware canvases. *ACM Transactions on Graphics*, 42(4):1–12, 2023.
- [67] Hang Yu, Yufei Xu, Jing Zhang, Wei Zhao, Ziyu Guan, and Dacheng Tao. Ap-10k: A benchmark for animal pose estimation in the wild. *arXiv preprint arXiv:2108.12617*, 2021.
- [68] Weiyu Zhang, Menglong Zhu, and Konstantinos G Derpanis. From actemes to action: A strongly-supervised representation for detailed action understanding. In *Proceedings of the IEEE international conference on computer vision*, pages 2248–2255, 2013.
- [69] Zhoutong Zhang, Forrester Cole, Richard Tucker, William T Freeman, and Tali Dekel. Consistent depth of moving objects in video. *ACM Transactions on Graphics*, 2021.
- [70] Yang Zheng, Adam W Harley, Bokui Shen, Gordon Wetzstein, and Leonidas J Guibas. PointOdyssey: A large-scale synthetic dataset for long-term point tracking. In *Proc. CVPR*, 2023.

Appendix for TAPVid-3D

Table of Contents

1. [More Dataset Samples](#)
2. [Dataset Statistics](#)
3. [Metrics using Median, Per-Trajectory, and Local Neighborhood Scaling](#)
4. [Evaluations using Median, Per-Trajectory, and Local Neighborhood Scaling](#)
5. [Evaluations using Fixed Metric Distance Thresholds](#)
6. [Baselines Details and Compute Resources](#)
7. [Filtering Incorrect Trajectories](#)
8. [Dataset Specifications, Metadata, and other Details](#)
9. [Visualized Samples](#)

6 More Dataset Samples

More dataset samples are provided on our website, <https://tapvid3d.github.io>, including interactive 3D visualizations.

Finally, we include static visualizations of trajectories in the figures included in the [Visualized Samples](#) section at the end of this PDF.

7 Dataset Statistics

Figure 3 showcases various summary statistics about the TAPVid-3D datasets and its 3D point trajectories. In the top left, we have the distribution of the number of frames in each video. The ADT-sourced videos contain the longest videos, and clips of 300 frames were extracted. Similarly Panoptic Studio contains clips of 150 frames, while DriveTrack contains clips of varying duration. In the top right, we have the number of point tracks annotated in each clip. In the bottom right, we count the number of ‘static’ trajectories in each video, marking a trajectory as static if the distance between all pairwise locations within a single point’s trajectory is less than 1 centimeter. The roughly 10 DriveTrack videos consisting of static trajectories are usually cars stopped at stoplights. These ‘static’ videos are a small minority of the 4000+ clips in TAPVid-3D. In the bottom right, we show the average velocity of each trajectory in the dataset, noting that trajectories in DriveTrack are the fastest. These histograms convey that there is a diversity of overall trajectory lengths, video lengths, and point velocities in the TAPVid-3D dataset. Additionally, this dataset is larger than two widely used 2D point tracking real-world-video datasets: TAPVid-Kinetics (1,189 videos) and TAPVid-DAVIS (30 videos).

8 Metrics using *Median, Per-Trajectory, and Local Neighborhood Rescaling*

In the results included in the main paper, we compute the 3D Average Jaccard and APD metrics using a *global median* rescaling procedure (L277). To get a good score, the entire scene must be reconstructed up to scale, and dynamic objects must be placed precisely. This is useful for many applications, such as navigation, but for others it may be overly stringent. If there is little camera motion, or if some of the objects have unclear size, it may also be very difficult for models to infer global scene shape.

However, not all applications require such strong global scene shape capabilities. For example, for imitation learning, we may want an agent to simply approach an object. For such applications, measuring the relative depth of estimated 3D locations along a *single trajectory* may be sufficient, and it may be substantially easier, as the (2D) scaling of local textures may provide enough information to solve the problem. More generally, for robotic imitation of an assembly task, it is the *local* consistency that’s most important: as long as points that are near each other in 3D have the correct depth relative

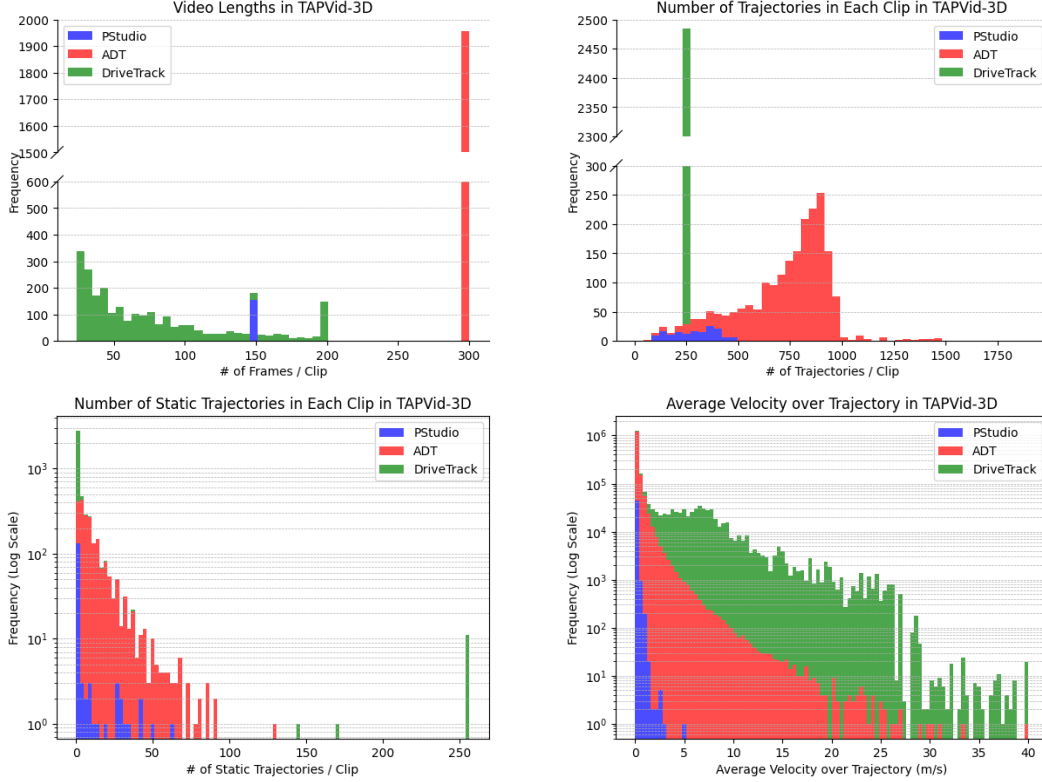


Figure 3: Statistics on TAPVid-3D. Top left: video lengths. Top right: Number of trajectories in each clip. Bottom left: Number of static tracks in each clip. Bottom right: average point velocity.

to one another, then the relative pose of the assembled parts will be clear, especially at the critical stage when the assembled parts are close together.

To enable more rapid progress in such domains, we propose two additional approaches to rescaling estimated trajectories to match the ground-truth 3D point cloud: *Per-Trajectory*, and *Local Neighborhood*. We apply the same Average Jaccard metric regardless of how the points are rescaled, although in the case of Local Neighborhood, the ground truth trajectories are also slightly modified, as explained below. In all cases, users can evaluate the same predictions using any metric without providing any extra information.

Per-Trajectory scaling is computed by rescaling each track P^i separately, multiplying by $\|P_{t_q}^i\|/\|\hat{P}_{t_q}^i\|$, where t_q is the query frame index, and then computing 3D AJ as before. As a result, methods must only compute the *relative* depth for the point at each time, relative to the query frame.

Local Neighborhood is somewhat more involved. Here, the goal is to capture whether *nearby* points are scaled correctly relative to one another, even if distant parts of the scene may not be. For example, if the goal is to understand an action depicted in a video that uses a tool, it is typically important to understand where the hand is relative to the tool, and where the tool is relative to the objects it’s acting on. The distance to the backgrounds—such as the back wall of the room—may not be obvious, especially if there’s relatively little camera motion. However, precisely computing these distances is not relevant to understanding the tool’s motion.

Intuitively, we wish to find an intermediate between two extremes: either rescaling the entire scene with a single scale factor, or rescaling every point with its own scale factor. To this end, we propose to scale each track according to the other track segments that intersect with its *4D tubelet* [32], according to a fixed neighborhood radius.

Specifically, we start by choosing a single neighborhood radius τ specified in meters. For a given trajectory P^i , we first find all points that are within τ meters of the ground truth on any frame, which define the tubelet $\mathcal{T}(P^i)$ associated to the trajectory P^i :

$$\mathcal{T}(P^i) = \{P_t^j \text{ s.t. } \|P_t^j - P_t^i\| < \tau\}.$$

Note that tubelet $\mathcal{T}(P^i)$ includes the trajectory P^i entirely, plus the portions of the trajectories of the other points where they come closer than the tubelet’s radius τ . For each selected ground-truth point in the tubelet, we select the corresponding points from the predictions to construct $\mathcal{T}(P^i; \hat{P}^i)$ as

$$\mathcal{T}(P^i; \hat{P}^i) = \{\hat{P}_t^j \text{ s.t. } \|P_t^j - P_t^i\| < \tau\}.$$

Analogously, we select the predicted and ground-truth visibility $\mathcal{T}(\hat{v}^i)$ and $\mathcal{T}(v^i)$.

Finally, given a predicted tubelet $\mathcal{T}(P^i; \hat{P}^i)$, we rescale all its points together using the ratio of query point distances $\|P_{t_q}^i\|/\|\hat{P}_{t_q}^i\|$, and evaluate the rescaled tubelet set as if it were a single trajectory, by replacing P^i , \hat{P}^i and v_i with their tubelet counterparts in equations (5) or (6) to compute the APD_{3D} and AJ_{3D} respectively. In our experiments, we set the radius τ to 3 centimeters for the PStudio and Aria scenes, and as 10 centimeters for the DriveTrack scenes, across all experiments. This is because tabletop manipulation and human-object motion likely require finer-grained movement than large vehicle movement in the Waymo Open public road scenes.

Baseline	Aria			DriveTrack			PStudio			Average		
	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow
Static Baseline	5.4	11.8	55.4	4.8	8.4	80.8	6.4	12.7	75.8	5.5	11.0	70.7
TAPIR + COLMAP	26.5	37.7	72.6	16.5	24.6	80.4	12.1	19.6	75.2	18.4	27.3	76.1
CoTracker + COLMAP	26.8	38.3	78.6	18.2	28.8	81.7	12.1	19.7	77.2	19.1	28.9	79.1
BootsTAPIR + COLMAP	28.8	41.3	78.6	20.0	29.3	83.8	12.9	20.8	81.8	20.6	30.4	81.4
TAPIR + ZoeDepth	16.2	24.2	79.7	7.4	12.2	81.6	12.0	20.0	78.7	11.9	18.8	80.0
CoTracker + ZoeDepth	17.4	26.3	87.8	6.7	12.3	82.6	12.0	20.8	80.0	12.0	19.8	83.4
BootsTAPIR + ZoeDepth	17.3	27.0	86.5	7.4	12.3	85.3	12.3	20.6	82.7	12.3	20.0	84.8
TAPIR-3D	8.5	14.9	86.0	10.2	17.0	83.3	7.2	13.1	78.9	8.6	15.0	82.8
SpatialTracker [62]	17.4	26.9	89.0	9.0	16.1	83.7	14.2	24.6	78.6	13.6	22.5	83.7

Table 5: **Using per-trajectory depth scaling.** We compare the performance on the `full_eval` split of several 2D-TAP models [10, 11, 21] combined with ZoeDepth [3] and COLMAP [47] on the TAPVid-3D benchmark. We also measure performance on the recently released SpatialTracker [62], and a static point baseline, in which the predicted trajectories are exactly the same as the query point.

Baseline	Aria			DriveTrack			PStudio			Average		
	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow
Static Baseline	5.5	11.8	56.0	4.8	8.4	80.8	6.4	12.6	75.7	5.5	10.9	70.8
TAPIR + COLMAP	23.8	34.4	72.8	12.9	20.3	80.4	9.9	16.5	75.1	15.5	23.7	76.1
CoTracker + COLMAP	24.6	35.3	78.8	15.7	25.2	81.7	10.8	17.7	77.1	17.0	26.1	79.2
BootsTAPIR + COLMAP	26.1	38.0	78.8	16.6	25.1	83.8	10.8	17.8	81.8	17.8	27.0	81.5
TAPIR + ZoeDepth	15.7	23.5	79.8	6.3	10.5	81.6	11.2	18.9	78.7	11.0	17.6	80.1
CoTracker + ZoeDepth	17.0	25.7	88.0	6.0	10.9	82.6	11.4	19.9	80.0	11.4	18.8	83.5
BootsTAPIR + ZoeDepth	16.8	26.3	86.7	6.4	10.9	85.3	11.6	19.6	82.6	11.6	18.9	84.9
TAPIR-3D	7.3	12.9	86.3	5.9	10.5	83.4	5.1	9.6	78.9	6.1	11.0	82.8
SpatialTracker [62]	16.7	25.7	89.3	6.9	12.4	83.7	12.3	21.6	78.5	12.0	19.9	83.8

Table 6: **Using local neighborhood scaling.** We compare the performance on the `full_eval` split of 2D-TAP models [10, 11, 21] combined with ZoeDepth [3] and COLMAP [47]. We also include SpatialTracker [62], and a static point baseline.

9 Evaluations with Median, Per-Trajectory, and Local Neighborhood Scaling

Tables 5, 6, 3 present additional experimental results on the `full_eval` set, on all our baselines. To avoid biasing the results to the TAPVid-3D splits with higher number of videos, these tables present averaged results across the three constituent data sources (weighing each source equally).

As expected, the AJ increases when using the less-strict local rescaling approaches. That is, *per-trajectory* scaling require less scale consistency than the *local neighborhood* metric, which itself is less stringent than the *global median scaling*. However, different methods improve by different amounts. Perhaps most surprisingly, COLMAP gives strong performance with local and per-trajectory rescaling, but underperforms ZoeDepth on Aria and Panoptic Studio when evaluated with global rescaling. This is likely because COLMAP completely fails to reconstruct moving content. For scenes where the majority of tracks are moving, the median rescaling will fail completely; therefore, ZoeDepth giving reasonable estimates for a larger fraction of points gives it an advantage.

TAPIR-3D, unsurprisingly, presents poor performance using global or local scaling, as it does not provide relative depth estimates for different tracks. However, evaluated with per-trajectory scaling, it gives competitive results on DriveTrack, even outperforming SpaTracker. This is somewhat surprising given that it operates on a completely different principle than other methods; it is trained using entirely synthetic data and does not use any geometric constraints, nor relies on monodepth models providing geometric priors. Overall, the different strengths and weaknesses of these highly-diverse methods suggests that the best performance will come from a method that combines ideas from all three. SpaTracker is a step in this direction, using a monodepth initialization while checking the 2D consistency of 3D reconstructions via reprojection, similar to COLMAP, and it often gives competitive performance. We hope that this benchmark can provide a way to quantify how well future methods in this vein accomplish the task.

Finally, we include a static baseline, which predicts the static 3D point $P_q = K^{-1}[x_q, y_q, 1]^T \cdot Z_q$ for all timestamps, in order to quantify the impact of motion. Note that this baseline still requires ground truth depth for the query points, and so isn't trivial to reproduce automatically; however, it still performs very poorly, even for the PStudio dataset where the camera is static (note that the static baseline is static *in the camera coordinate frame*). Thus, we conclude tracking the camera and tracking the objects is important for obtaining strong performance.

10 Evaluations using Fixed Metric Distance Thresholds

In the Average Jaccard formulation in Section 3.5, we describe how we use determine correctly predicted points along a trajectory, using a depth-adaptive radius threshold denoted $\delta_{3D}(P_t^i)$. We also explored using a fixed metric threshold. Specifically, instead of the $\{1, 2, 4, 8, 16\}$ pixel thresholds (projected into 3D), we use a the fixed metric radius thresholds of 1 centimeter, 4 centimeter, 16 centimeters, 64 centimeters, and 2.56 meters. If the predicted point is within this distance to the ground truth point, it is marked as position correct within that threshold. Table 7 describes the model baselines results using this alternative metric.

	Baseline	Aria			DriveTrack			PStudio			Average		
		3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow	3D-AJ \uparrow	APD \uparrow	OA \uparrow
Static Baseline		21.2	40.3	55.4	5.8	9.7	80.8	31.7	46.4	75.8	19.6	32.1	70.7
TAPIR + COLMAP		20.5	32.2	72.6	15.5	21.9	80.4	28.8	41.7	75.2	21.6	31.9	76.1
CoTracker + COLMAP		23.6	33.6	78.6	18.2	25.7	81.7	31.6	44.1	77.2	24.4	34.4	79.1
BootsTAPIR + COLMAP		24.0	35.5	78.6	18.7	25.2	83.8	31.6	43.6	81.8	24.7	34.7	81.4
TAPIR + ZoeDepth		28.2	41.7	79.7	11.0	15.9	81.6	39.0	55.2	78.7	26.1	37.6	80.0
CoTracker + ZoeDepth		32.7	44.0	87.8	10.7	16.2	82.6	40.2	56.1	80.0	27.8	38.8	83.4
BootsTAPIR + ZoeDepth		31.9	45.6	86.5	11.4	16.3	85.3	41.5	55.9	82.7	28.3	39.2	84.8
TAPIR-3D		19.2	29.9	86.0	7.0	10.9	83.3	24.8	36.1	78.9	17.0	25.6	82.8
SpatialTracker [62]		33.1	45.0	89.0	13.1	19.3	83.7	39.5	56.1	78.6	28.5	40.2	83.7

Table 7: **Using fixed metric thresholds for AJ, with median scaling.** We compare the performance on the full_eval split of several 2D-TAP models [10, 11, 21] combined with ZoeDepth [3] and COLMAP [47]. We also include SpatialTracker [62] and the static point baseline.

11 Baselines Details and Compute Resources

CoTracker. We use the pretrained model and PyTorch code from the official CoTracker codebase and run inference enabling the bi-directional tracking mode, with no other modifications to the default parameters. Internally, inference is performed in 512×384 resolution, and the output predictions are rescaled back to the original clip resolutions. Inference is performed using A100 GPUs, and

processing each dataset clip takes about 30s, totaling roughly 38 GPUh for running CoTracker on the whole benchmark.

BootsTAPIR and TAPIR. We use the pretrained models and JAX code from the official codebase and run inference with the default parameters. Internally, inference is performed in 256×256 resolution, and the output predictions are rescaled back to the original clip resolutions. Inference was performed in a CPU cluster using up to 1024 CPUs and totalling about 22800 CPUh.

COLMAP. For running COLMAP, we dumped the 2D tracks estimated by CoTracker, TAPIR and BootsTAPIR as a set of per-frame image feature files and corresponding matches in txt format and imported those in COLMAP using the ‘feature_importer’ and ‘matches_importer’ functionality. We then perform 3D reconstruction through the incremental mapping pipeline (‘mapper’). As each input 2D track can lead to multiple reconstructed 3D points across time (eg. for moving objects), we only keep those with larger “track length” (number of images where that 3D point was reconstructed from). Finally, we transform the resulting reconstructed 3D points positions in world coordinates to camera coordinates using the predicted extrinsic parameters. Inference was performed in a CPU cluster using up to 1024 CPUs and totalling about 14000 CPUh.

ZoeDepth. We used the pretrained models and PyTorch code from the official codebase and run inference with the default parameters. Inference is performed in the native resolution for the ADT and Panoptic Studio clips and in 720×480 for DriveTrack, where the original resolution was too large for running inference in this model on a standard GPU. Inference was performed on 16 V100 GPUs, totalling about 200GPUh.

TAPIR3D

We propose a straightforward extension of TAPIR to 3D by training on 3D ground truth from Kubric [13]. As Kubric is synthetic data, it is straightforward to obtain ground-truth 3D point tracks. However, we don’t expect that Monodepth models trained on Kubric will generalize, as the scenes are very different from real ones. However, we expect that a model trained here might be able to estimate *relative* depth of a point at different times on the trajectory, relative to the query point. Like with point tracking, we expect that low-level texture information may be sufficient to predict the relative depth (specifically, the scaling of the texture elements), and so high-level semantic understanding won’t be necessary, meaning that it can be learned from a semantically meaningless dataset.

Specifically, we train TAPIR to output the log depth of each point on the trajectory, relative to the query point. This is a scalar quantity, and can be predicted using the same network structure as the other scalar quantities, e.g., the occlusion logit. That is, we predict an initial log-space scale factor for every frame by adding an extra head to TAPIR’s occlusion prediction network, which performs convolutions on top of the cost volume for each frame followed by global average pooling. Then we feed this estimate to the iterative refinement steps by concatenating it with the local score maps, the initial occlusion estimate, and so on, passing it through the 1-D convolutional network which produces an update; again, we add an extra head on this network which produces an updated relative depth estimate. We apply an L1 loss on the estimate for both the initialization and each of the four refinement passes, with a weight of 1.0. Otherwise, we train the entire network using the procedure described in [10].

12 Filtering Incorrect Trajectories

We apply different automatic filters for removing problematic tracks. Tracks can present three type of issues: (i) issues with visibility flags, (ii) queries which are outside the moving objects, and (iii) noisy 3D trajectories.

We found visibility issues (i) to be present in all dataset splits, and we remove it simply by oversampling the number of query points and discarding those whose visibility flag changes state more than a 10% of the number of frames in the video.

Issue (ii) was present mostly in the DriveTrack split, where trajectories in a video are localized and describe the motion of exactly one moving object in the scene. In some cases the 3D point-clouds associated with vehicles also contain points that are within the object bounding box, but outside of the object itself, such as in the road. To filter out errant trajectories, we use the Segment Anything

model (SAM) to generate an object mask for each frame [24]. We prompt SAM with a point prompt, computed by taking the geometric median of DriveTrack trajectories at each point in time.

Finally, we found that noisy 3D trajectories (iii) could occur in the Panoptic Studio split, where sometimes the reconstructed 3D Gaussians were not sufficiently constrained due to surfaces having uniform colors. In this case we apply a similar approach as before, and score trajectories based on the percentage of time they are on foreground object masks across all camera viewpoints. We perform a hyperparameter search on the threshold value and select the points that stay on the object masks at least 75% of the time across all masks, which removes most of the problematic points.

13 Dataset Specifications, Metadata, and other Details

For this dataset release, we preserve the licenses for the constituent original data sources, which are non-commercial. For our additions, to the extent that we can, we release under a standard Apache 2.0 license. A full amalgamated license will be available in the open-sourced repository during complete release of the work, after the review period is finished.

We will publicly host the dataset for wide consumption by researchers on Google Cloud Storage indefinitely. Part of the dataset is already hosted in this way (and how the Colab link linked above is able to run). We also intend to open-source code for computing the new 3D-AJ metrics after the camera ready. We anticipate the release will require little maintenance (and the TAPVid-2D dataset release that the team released two years ago is similarly low maintenance), but we are happy to address any emergent issues raised by users.

Specific implementation details on how the dataset can be read are found in the Colab link provided. Each dataset example is provided in a *.npz file, containing the fields: `tracks_xyz` of shape [T, Q, 3] (containing the Q ground truth point tracks for the corresponding video of T frames, with (x, y, z) -coordinates in meters), `query_xyt` of shape [Q, 3] (containing each track's query point position in format (x, y, t) , in (x, y) -pixel space and t as the query frame index), the ground truth `visibility` flags with shape [Q, T], and the `camera_intrinsics` (as $[f_x, f_y, c_x, c_y]$). Each *.npz file is named after its corresponding video in the original data source, which can be loaded by downloading from the original hosting sites [2, 42, 51], respecting their corresponding licenses.

14 Visualized Samples

See Figures 4, 5, 6, 7, 8, 9, 10, 11, and 12 below.

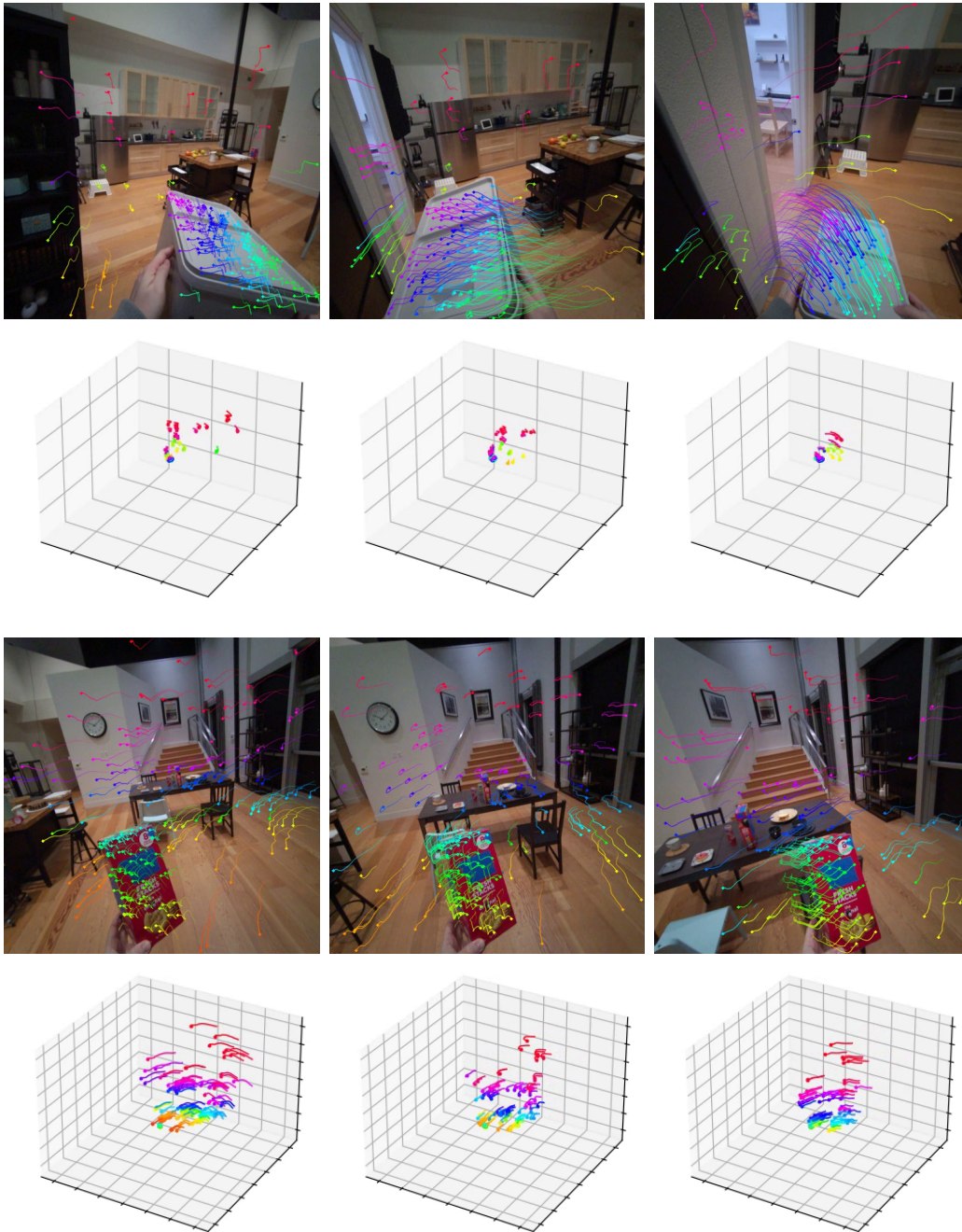


Figure 4: Random samples from ADT subset in TAPVid-3D: on the top row, we visualize the point trajectories projected into the 2D video frame; on the bottom row, we visualize the metric 3D point trajectories. For each video, we show 3 frames sampled at time step 30, 60 and 90.

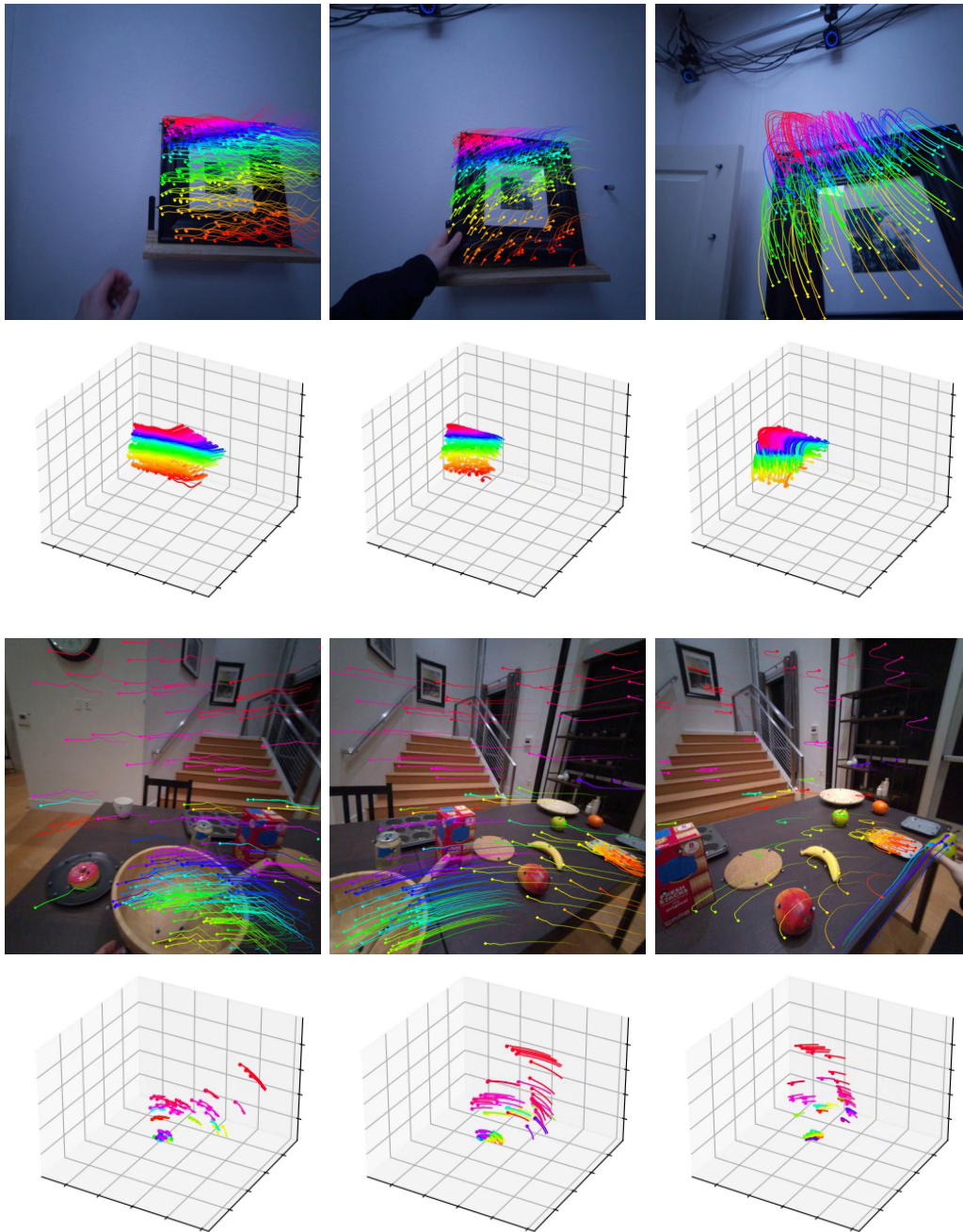


Figure 5: Random samples from ADT subset in TAPVid-3D (cont'd.): on the top row, we visualize the point trajectories projected into the 2D video frame; on the bottom row, we visualize the metric 3D point trajectories. For each video, we show 3 frames sampled at time step 30, 60 and 90.

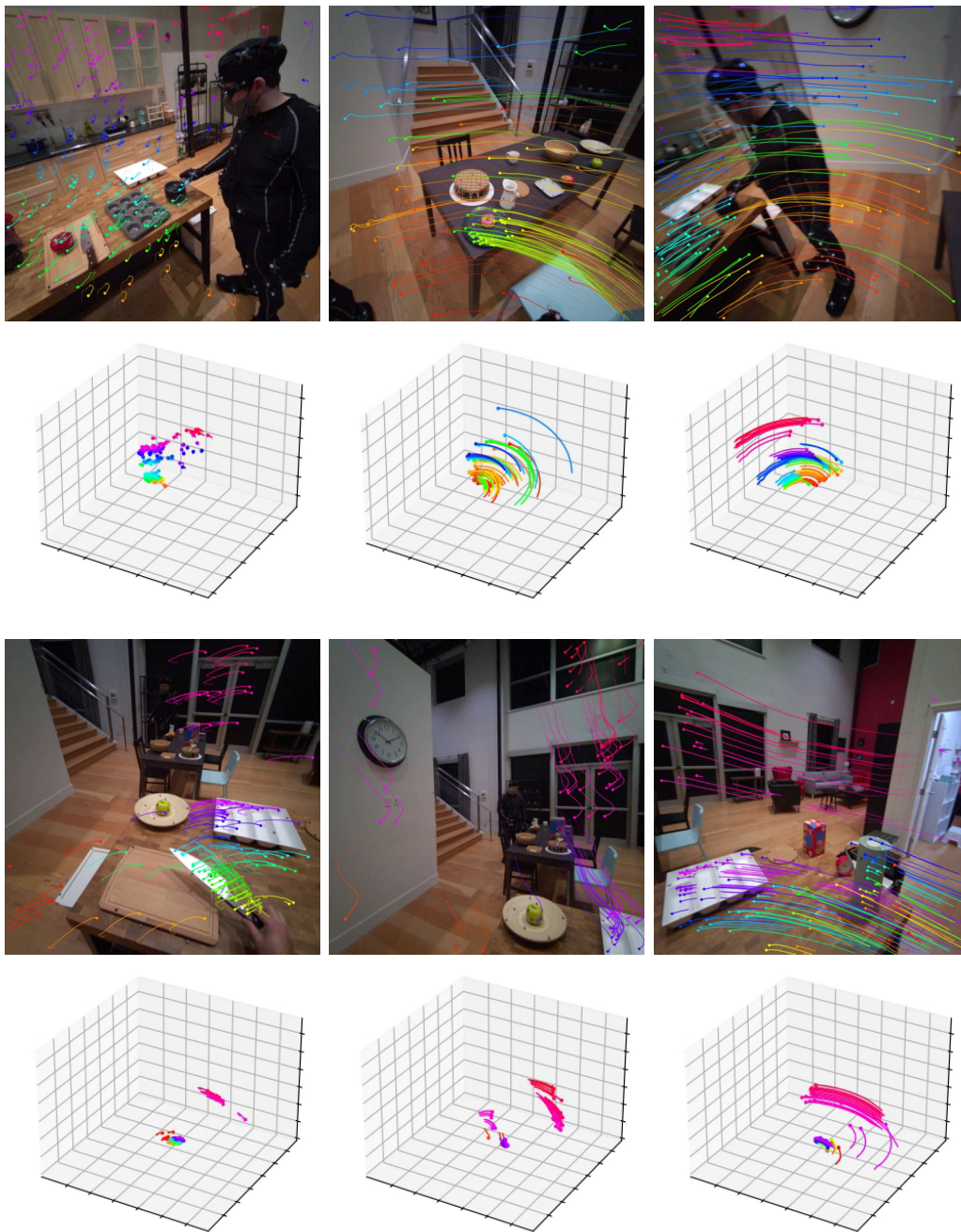


Figure 6: Random samples from ADT subset in TAPVid-3D (cont'd.): on the top row, we visualize the point trajectories projected into the 2D video frame; on the bottom row, we visualize the metric 3D point trajectories. For each video, we show 3 frames sampled at time step 30, 60 and 90.

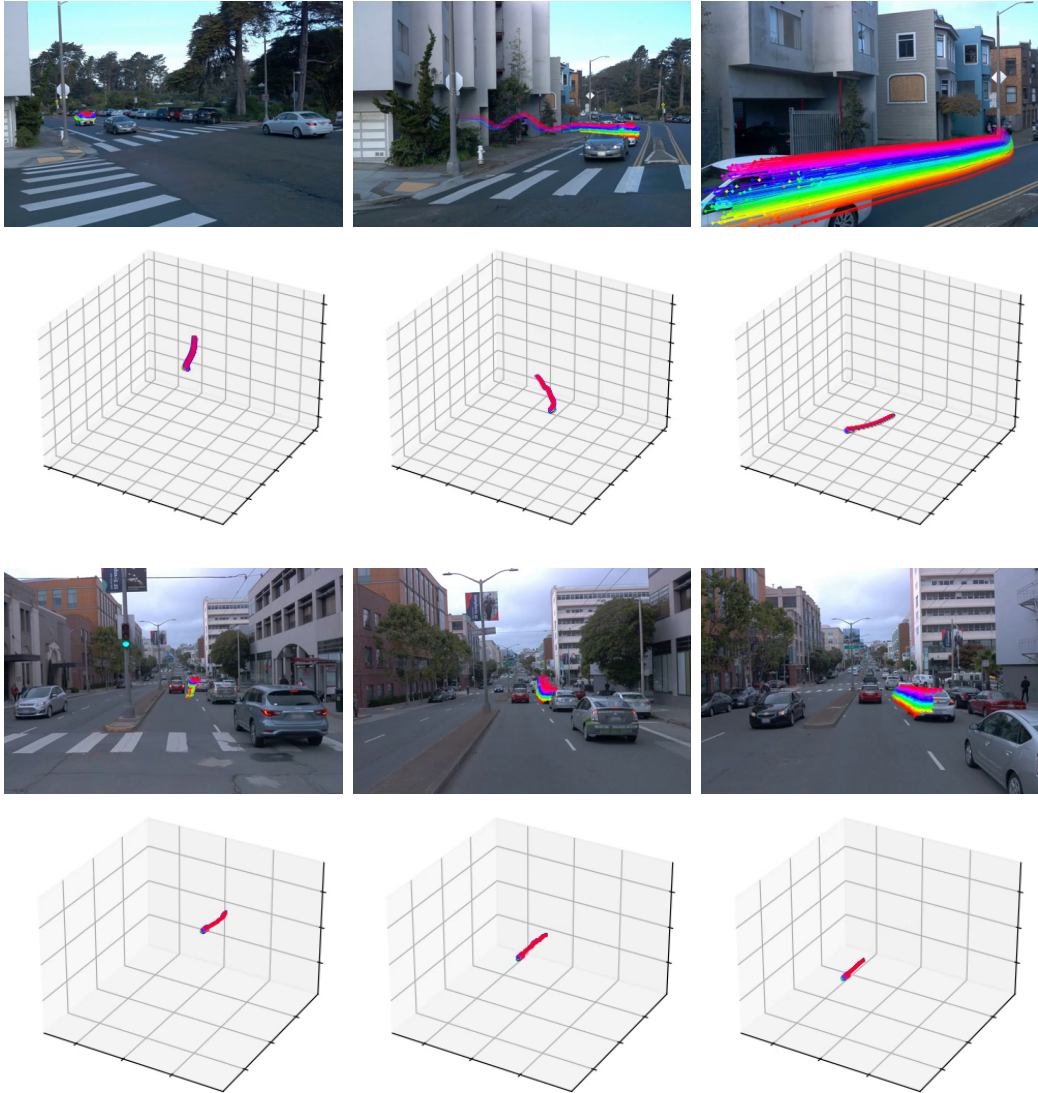


Figure 7: Random samples from DriveTrack subset in TAPVid-3D: on the top row, we visualize the point trajectories projected into the 2D video frame; on the bottom row, we visualize the metric 3D point trajectories. For each video, we show 3 frames sampled at time step 30, 60 and 90.

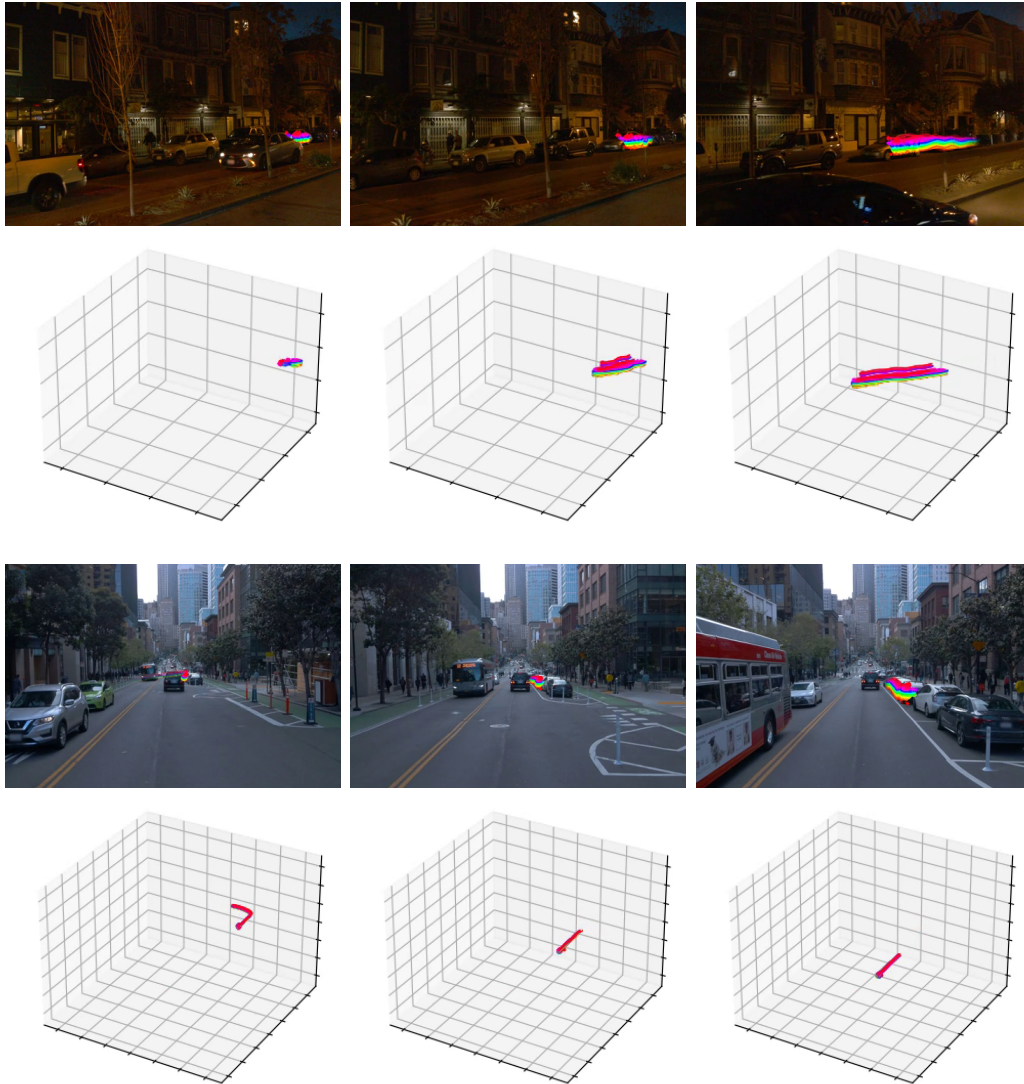


Figure 8: Random samples from DriveTrack subset in TAPVid-3D (cont'd.): on the top row, we visualize the point trajectories projected into the 2D video frame; on the bottom row, we visualize the metric 3D point trajectories. For each video, we show 3 frames sampled at time step 30, 60 and 90.

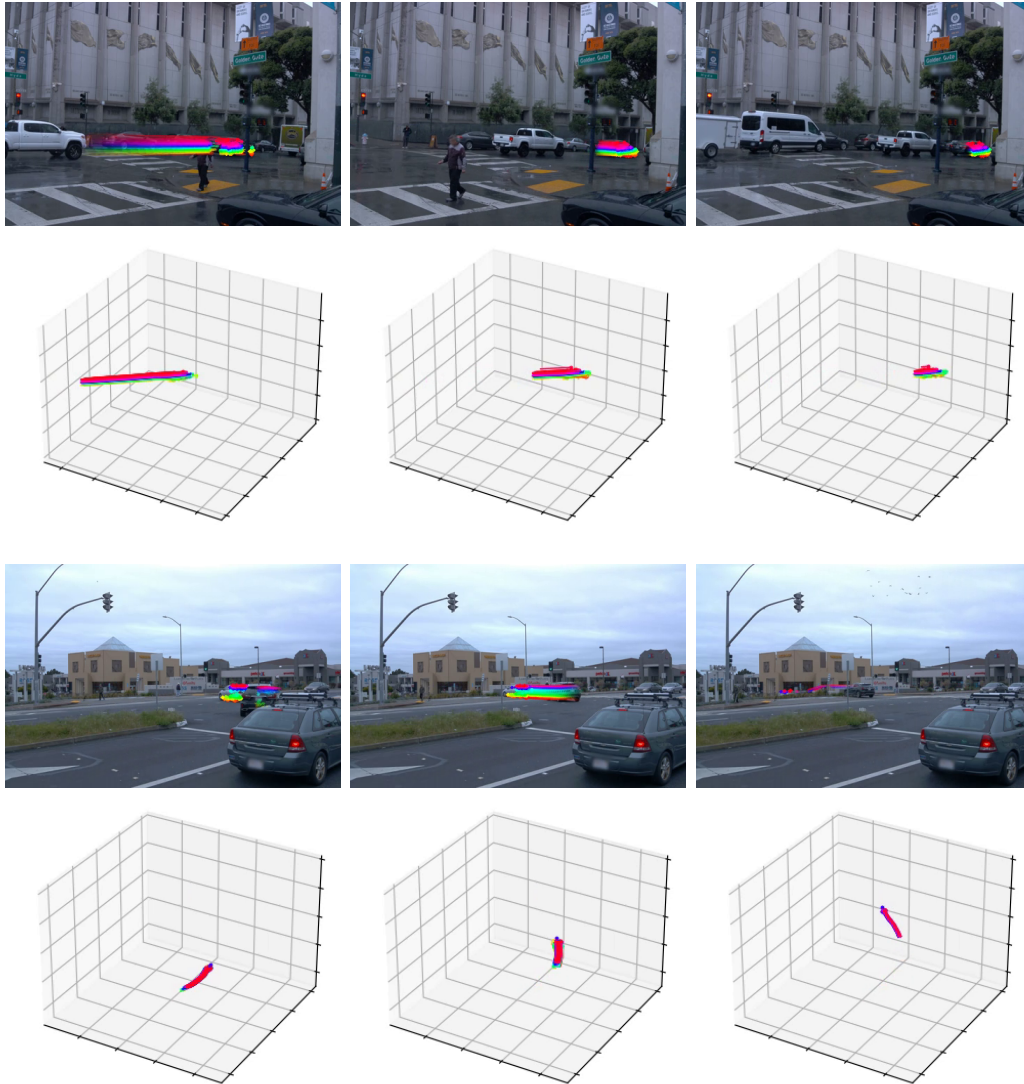


Figure 9: Random samples from DriveTrack subset in TAPVid-3D (cont'd.): on the top row, we visualize the point trajectories projected into the 2D video frame; on the bottom row, we visualize the metric 3D point trajectories. For each video, we show 3 frames sampled at time step 30, 60 and 90.

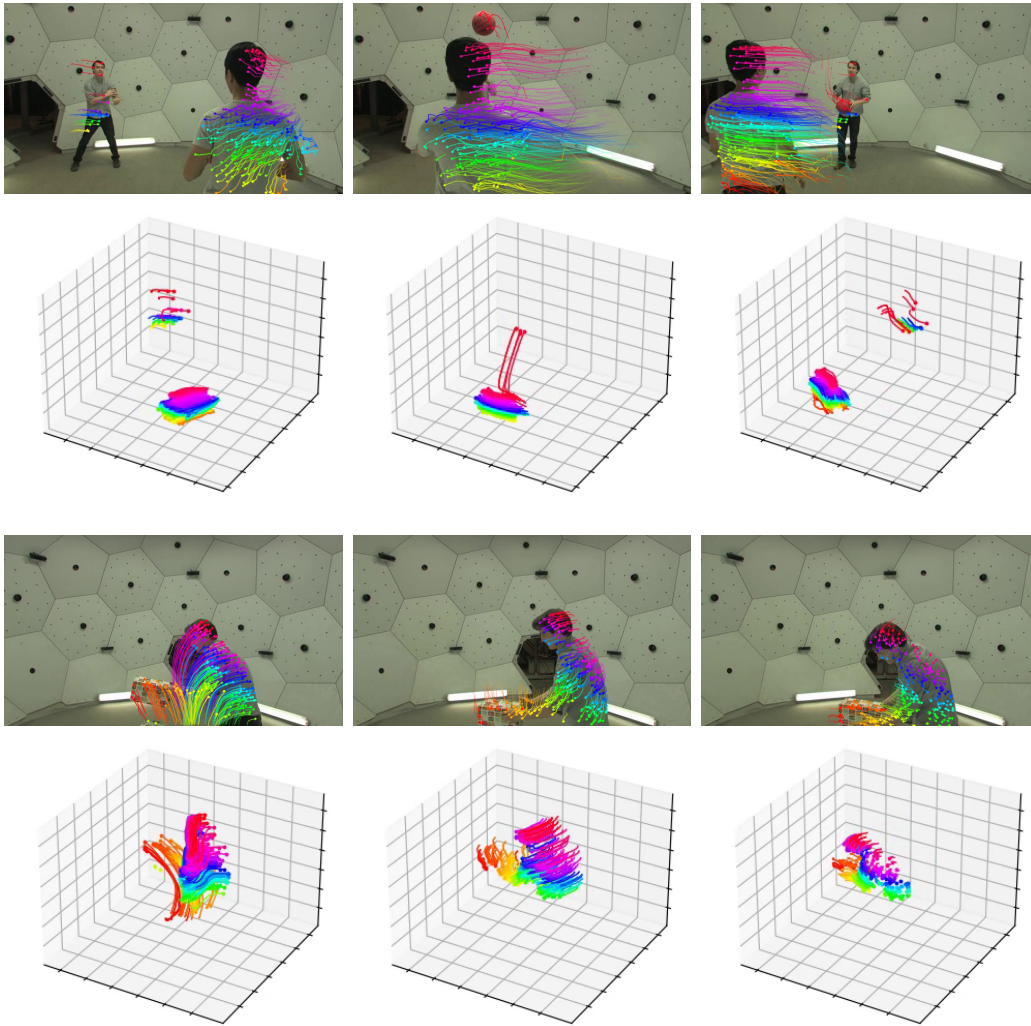


Figure 10: Random samples from Panoptic Studio subset in TAPVid-3D: on the top row, we visualize the point trajectories projected into the 2D video frame; on the bottom row, we visualize the metric 3D point trajectories. For each video, we show 3 frames sampled at time step 30, 60 and 90.

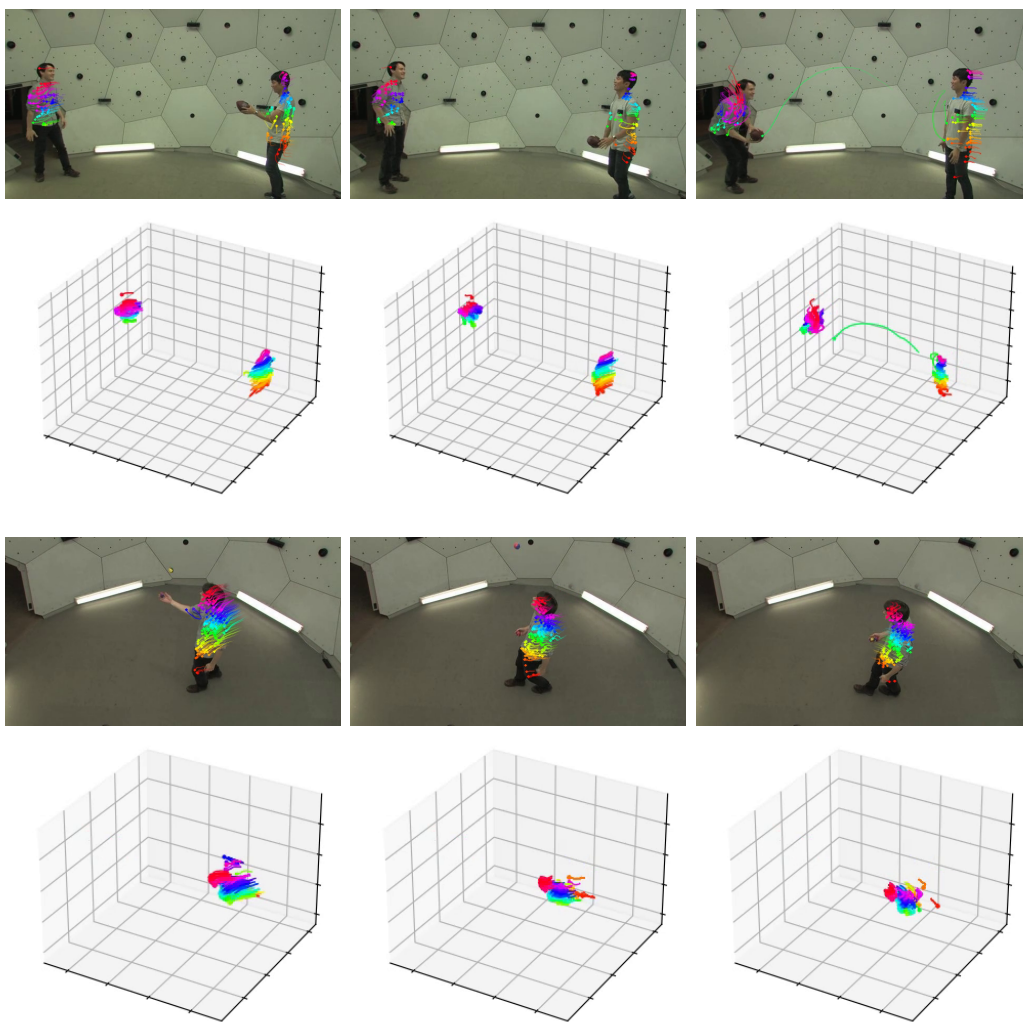


Figure 11: Random samples from Panoptic Studio subset in TAPVid-3D (cont'd.): on the top row, we visualize the point trajectories projected into the 2D video frame; on the bottom row, we visualize the metric 3D point trajectories. For each video, we show 3 frames sampled at time step 30, 60 and 90.

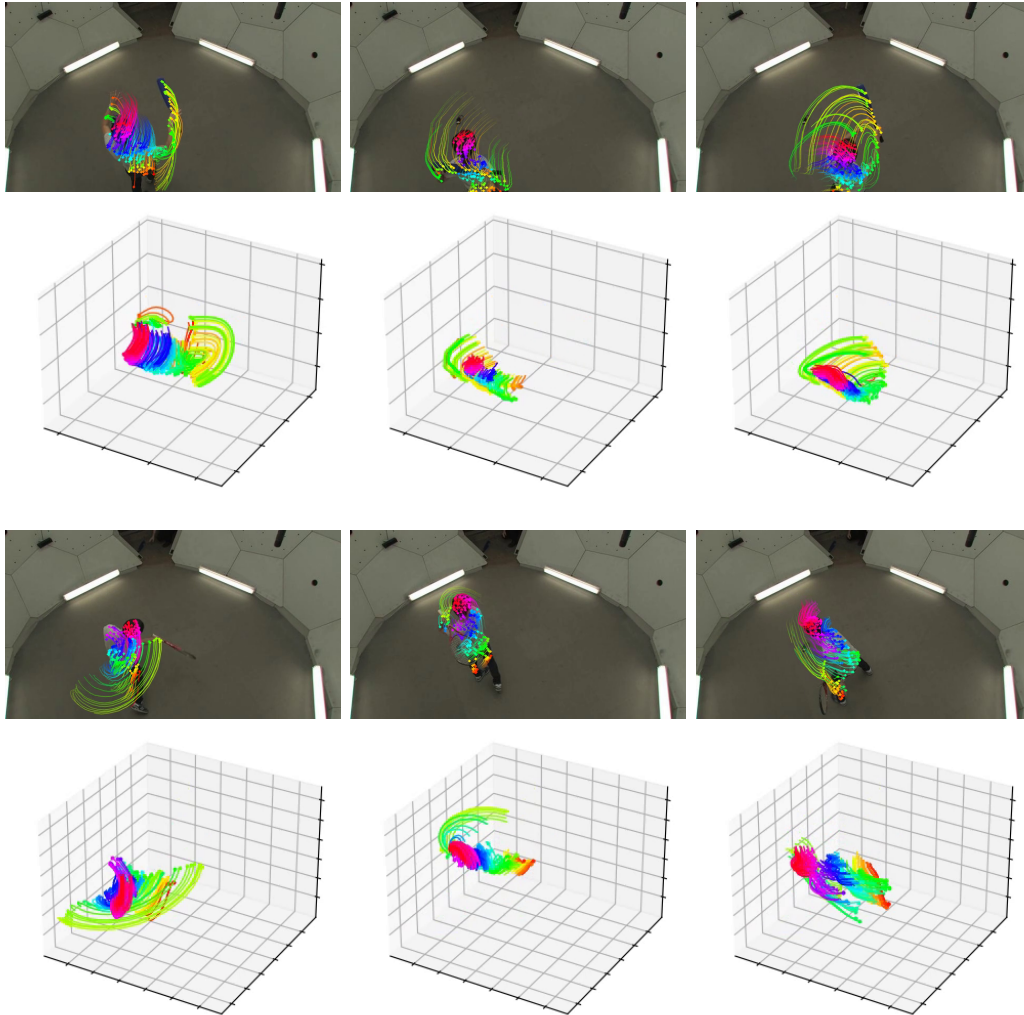


Figure 12: Random samples from Panoptic Studio subset in TAPVid-3D (cont'd.): on the top row, we visualize the point trajectories projected into the 2D video frame; on the bottom row, we visualize the metric 3D point trajectories. For each video, we show 3 frames sampled at time step 30, 60 and 90.