# Tracking People by Predicting 3D Appearance, Location & Pose

Jathushan Rajasegaran, Georgios Pavlakos, Angjoo Kanazawa, Jitendra Malik
UC Berkeley

## Abstract

*In this paper, we present an approach for tracking people in monocular videos, by predicting their future 3D representations. To achieve this, we first lift people to 3D from a single frame in a robust way. This lifting includes information about the 3D pose of the person, his or her location in the 3D space, and the 3D appearance. As we track a person, we collect 3D observations over time in a tracklet representation. Given the 3D nature of our observations, we build temporal models for each one of the previous attributes. We use these models to predict the future state of the tracklet, including 3D location, 3D appearance, and 3D pose. For a future frame, we compute the similarity between the predicted state of a tracklet and the single frame observations in a probabilistic manner. Association is solved with simple Hungarian matching, and the matches are used to update the respective tracklets. We evaluate our approach on various benchmarks and report state-of-the-art results.*

## 1. Introduction

When we watch a video, we can segment out individual people, cars, or other objects and track them over time. The corresponding task in computer vision has been studied for several decades now, with a fundamental choice being whether to do the tracking in 2D in the image plane, or of 3D objects in the world. The former seems simpler because it obviates the need for inferring 3D, but if we do take the step of back-projecting from the image to the world, other aspects such as dealing with occlusion become easier. In the 3D world the tracked object doesn't disappear, and even young infants are aware of its persistence behind the occluder. A recent paper, Rajasegaran *et al.* [32] argues convincingly on the 3D side of this debate for people tracking, and presents experimental evidence that indeed performance is better with 3D representations. In this paper, we will take this as granted, and proceed to develop a system in the 3D setting of the problem. While our approach broadly applies to any object category where parameterized 3D models are available and can be inferred from images, we will limit ourselves in this paper to studying people, the



Figure 1. **Tracking people by predicting and matching in 3D.** The top row shows our tracking results at three different frames. The results are visualized by a colored head-mask for unique identities. The second and third rows show renderings of the 3D states of the two people in their associated tracklets. The bottom row shows the bottom-up detections in each image frame which, after being lifted to 3D, will be matched with the 3D predictions of each tracklet in the corresponding frame. Note how in the middle frame of second row, the 3D representation of the person persists even though he is occluded in the image. More videos at project site.

most important case in practice.

Once we have accepted the philosophy that we are tracking 3D objects in a 3D world, but from 2D images as raw data, it is natural to adopt the vocabulary from control theory and estimation theory going back to the 1960s. We are interested in the "state" of objects in 3D, but all we have access to are "observations" which are RGB pixels in 2D. In an online setting, we observe a person across multiple time frames, and keep recursively updating our estimate of the person's state — his or her appearance, location in the world, and pose (configuration of joint angles). Since we have a dynamic model (a "tracklet"), we can also predict states at future times. When the next image frame comes in, we detect the people in it, lift them to 3D, and in that setting solve the association problem between these bottom-up detections and the top-down predictions of the different tracklets for this frame. Once the observations have been associated with the tracklets, the state of each person is re-estimated and the process continues. Fig. 1 shows this pro-
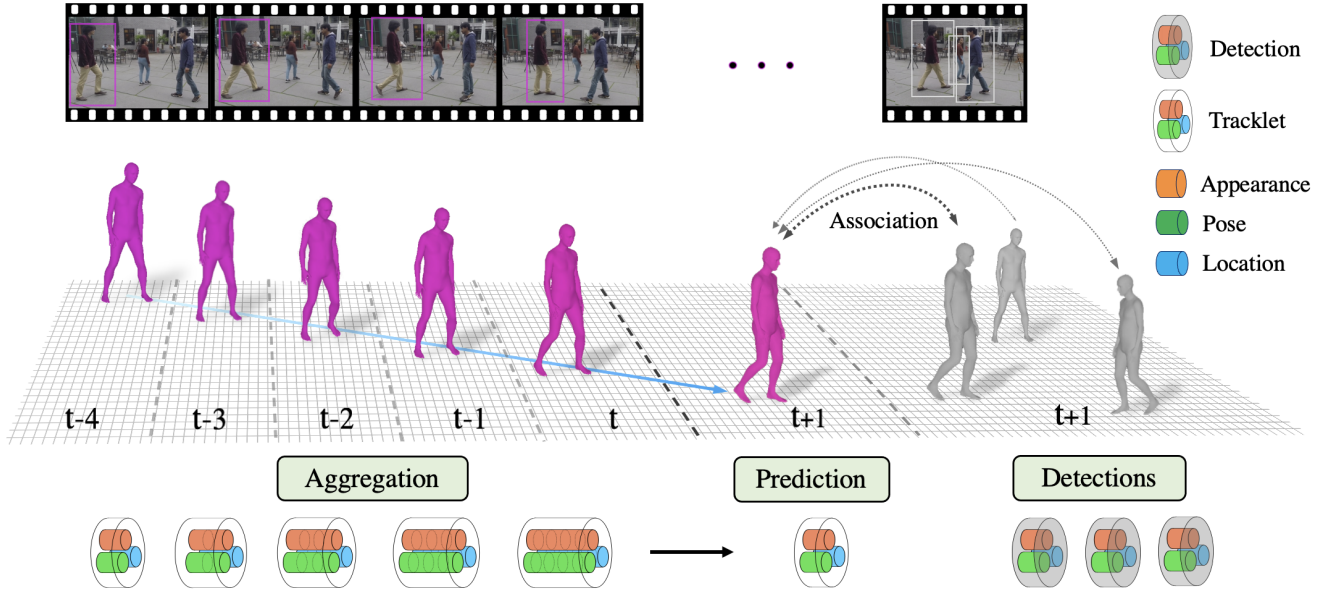
Figure 2. **PHALP: Predicting Human Appearance, Location and Pose for Tracking**: We perform tracking of humans in 3D from monocular video. For every input bounding box, we estimate a 3D representation based on the 3D appearance, 3D pose and 3D location of the person. During tracking, these are integrated to form corresponding tracklet-based representations. We perform tracking by predicting the future representation of each person and using it to solve for association given the detected bounding boxes of a future frame.

cess at work on a real video. Note that during a period of occlusion of a tracklet, while no new observations are coming in, the state of the person keeps evolving following his or her dynamics. It is not the case that "Out of sight, out of mind"!

In an abstract form, the procedure sketched in the previous paragraph is basically the same as that followed in multiple computer vision papers from the 1980s and 1990s. The difference is that in 2021 we can actually make it work thanks to the advances brought about by deep learning and big data, that enable consistent and reliable lifting of people to 3D. For this initial lifting, we rely on the HMAR model [32]. This is applied on every detected bounding box of the input video and provides us with their initial, single frame, observations for 3D pose, appearance as well as location of the person in the 3D space.

As we link individual detections into tracklets, these representations are aggregated across each tracklet, allowing us to form temporal models, *i.e.*, functions for the aggregation and prediction of each representation separately (see left side of Fig. 2). More specifically, for appearance, we use the canonical UV map of the SMPL model to aggregate appearance, and employ its most recent version as a prediction of a person's appearance. For pose, we aggregate information using a modification of the HMMR model [15], where through its "movie strip" representation, we can produce 3D pose predictions. Finally, for 3D location, we use linear regression to predict the future location of the person.

This modeling enables us to develop our tracking system, PHALP (Predicting Human Appearance, Location and Pose for tracking), which aggregates information over time, uses it to predict future states, and then associates the predictions with the detections. First, we predict the 3D location, 3D pose and 3D appearance for each tracklet for a short period of time (right side of Fig. 2). For a future frame, these predictions need to be associated with the detected people of the frame. To measure similarity, we adopt a probabilistic interpretation and compute the posterior probabilities of every detection belonging to each one of the tracklets, based on the three basic attributes. With the appropriate similarity metric, association is then easily resolved by means of the Hungarian algorithm. The newly linked detections can now update the temporal model of the corresponding tracklets for 3D pose, 3D appearance and 3D location in an online manner and we continue the procedure by rolling-out further prediction steps. The final output is an identity label for each detected bounding box in the video. Notably, this approach can also be applied on videos with shot changes, e.g., movies [10], with minor modifications. Effectively, we only modify our similarity to include only appearance and 3D pose information for these transitions, since they (unlike location) are not affected by the shot boundary.

## 2. Related work

**Tracking.** Object tracking is studied in various settings such as single object tracking, multi-object tracking for hu-

mans, and multi-object tracking for vehicles etc. The tracking literature is vast and we refer readers to [6, 7, 42] for a comprehensive summary. In general tracking can be designed for any generic category, however, in this section we discuss the methods that focus on tracking humans. These approaches mostly work in a tracking by detection setting, where 2D location, key-point features [9, 34, 38] and 2D appearance [4, 29, 40, 41] is used to associate detections over time. Quality of the detection plays a key role in tracking-by-detection setting and many works jointly learn or fine-tune their own detection models [4, 29]. In this work, we are interested in the effectiveness of 3D representations for tracking and thus assume that detection bounding boxes are provided, which we associate through our representations. On the other hand, tracking by regression predicts future locations using the knowledge of the past detections. While this alleviate the requirement for good quality detections, most of the works regress in the image plane. The projection from 3D world to the image plane makes it hard to make this prediction, therefore these methods need to learn non-linear motion models [2, 4, 45]. Compared to these methods, PHALP predicts short-term location in 3D coordinates, by simple linear regression. Additionally, we also predict appearance and pose features for better association.

Finally, there are methods that incorporate 3D information in tracking, however these approaches assume multiple input cameras [24, 44] or 3D point cloud observation from lidar data [37]. In this paper we focus on the setting where the input is a monocular video. Some recent works tracks occluded people based on the object permanence [17, 35]. These methods learn complex functions to predict the locations of occluded people. However, by placing humans in 3D space and predicting their location, pose and appearance, object permanence is already built into our system.

**Monocular 3D human reconstruction.** Although there is a long history of methods for 3D human reconstruction from monocular images, *e.g.*, [5, 11], here we focus on more recent works. Many of the relevant approaches rely on the SMPL model [25], which offers a low dimensional parameterization of the human body. HMR [14] has been one of the most notable ones, using a neural network to regress the parameters of a SMPL body from a single image. Follow-up works have improved the robustness of the original model [19, 22], or added additional features like estimation of camera parameters [20], or probabilistic estimation of pose [23]. Recently, Rajasegaran *et al*. [32] introduced HMAR, by extending the model with an appearance head. Other works have focused on extending HMR to the temporal dimension, *e.g.*, HMMR [15], VIBE [18], MEVA [27] and more [31]. In this work, we make use of a modification of the HMAR model [32] as the main feature backbone, while also employing a model that follows the HMMR prin-

ciples [15] for temporal pose prediction, but instead, using a transformer [36] to aggregate pose information over time. Regarding human motion prediction, Kanazawa *et al*. [15], regress future poses from the temporal pose representation of HMMR, the "movie-strip". Zhang *et al*. [43] extend this to PHD, employing autoregressive prediction of human motion. Aksan *et al*. [1] also regress future human motion in an autoregressive manner, using a transformer.

## 3. Method

Tracking humans using 3D representations has significant advantages, including that appearance is independent of pose variations and the ability to have amodal completion for humans during partial occlusion. Our tracking algorithm accumulates these 3D representations over time, to achieve better association with the detections. PHALP has three main stages: 1) lifting humans into 3D representations in each frame, 2) aggregating single frame representations over time and predicting future representations, 3) associating tracks with detections using predicted representations in a probabilistic framework. We explain each stage in the next sections.

### 3.1. Single-frame processing

The input to our system is a set of person detections along with their estimated segmentation masks, provided by conventional detection networks, like Mask-RCNN [12]. Each detection is processed by our feature extraction backbone that computes the basic representations for pose, appearance and location on a single-frame basis. For this feature extraction we use a modification of the HMAR model [32]. HMAR returns a feature representation for the 3D pose $p$, for appearance $a$, while it can recover an estimate for the 3D location $l$ for the person.

The standard HMAR model takes as input the pixels in the bounding box corresponding to a detected person. This means that in a crowded, multi-person scenario, the input will contain pixels corresponding to more than one person in the bounding box, potentially confusing the network. To deal with this problem, we modify HMAR to take as additional input, the pixel level mask of the person of interest (this is readily available as part of the output of Mask R-CNN) and re-train HMAR. Obviously, we cannot expect this step to be perfect, since there can be inaccuracies in the bounding box detections or mask segmentations. However, we observed that the model gives more robust results in the case of close person-person interactions, which are common in natural videos.

### 3.2. 3D tracklet prediction

The 3D estimates for each detection provide a rich and expressive representation for each bounding box. However,

they are only the result of single-frame processing. During tracking, as we expand each tracklet, we have access to more information that is representative of the state of the tracklet along the whole trajectory. To properly leverage this information, our tracking algorithm builds a tracklet representation during every step of its online processing, which allows us to also predict the future states for each tracklet. In this section we describe how we build this tracklet representation, and more importantly, how we use it to *predict* the future state of each tracklet.

**Appearance:** The appearance pathway is used to integrate appearance information for each person over multiple frames. The single frame appearance representation for the person $i$ at time step $t$, $\mathbf{A}_t^i$, is taken from the HMAR model by combining the UV image of that person $\mathbf{T}_t^i \in \mathcal{R}^{3 \times 256 \times 256}$ and the corresponding visibility map $\mathbf{V}_t^i \in \mathcal{R}^{1 \times 256 \times 256}$ at time step $t$:

$$\mathbf{A}_t^i = [\mathbf{T}_t^i, \mathbf{V}_t^i] \in \mathcal{R}^{4 \times 256 \times 256}$$

Note that the visibility mask $\mathbf{V}_t^i \in [0, 1]$ indicates whether a pixel in the UV image is visible or not, based on the estimated mask from Mask-RCNN. Now, if we assume that we have established the identity of this person in neighboring frames, we can integrate the partial appearance information coming from the independent frames to an overall tracklet appearance for the person. Using the set of single frame appearance representations $\mathcal{A}^i = \{\mathbf{A}_t^i, \mathbf{A}_{t-1}^i, \mathbf{A}_{t-2}^i, ...\}$, after every new detection we create a singe per-tracklet appearance representation:

$$\widehat{\mathbf{A}}_t^i = \Phi_A(\widehat{\mathbf{A}}_{t-1}^i, \mathbf{A}_t^i) = (1 - \alpha) * \widehat{\mathbf{A}}_{t-1}^i + \alpha \mathbf{A}_t^i$$

$$\text{where,} \quad \alpha = \begin{cases} \alpha_0, & \text{if } \widehat{\mathbf{V}}_{t-1}^i = 1 \text{ and } \mathbf{V}_t^i = 1 \\ 1, & \text{if } \widehat{\mathbf{V}}_{t-1}^i = 0 \text{ and } \mathbf{V}_t^i = 1 \\ 0, & \text{if } \widehat{\mathbf{V}}_{t-1}^i = 1 \text{ and } \mathbf{V}_t^i = 0. \end{cases}$$

Here, $\Phi_A$ is the appearance aggregation function, which takes a weighted sum of the previous tracklet appearance representation and the new detection appearance representation. Note that, at the start of the tracklet we simply assign the initial single-frame representation to the tracklet representation ($\widehat{\mathbf{A}}_0^i = \mathbf{A}_0^i$). With this definition of $\Phi_A$, we can aggregate appearance information over time, while allowing the representation to change slowly to account for slight appearance changes of the person during a video. Moreover, the UV image provides appearance of each point on the body surface independently of body pose and shape which enables the simple summation operation on the pixel space, without any learnable components. Figure 3 shows how the UV image of the person is aggregated over time and used for association of new detections.

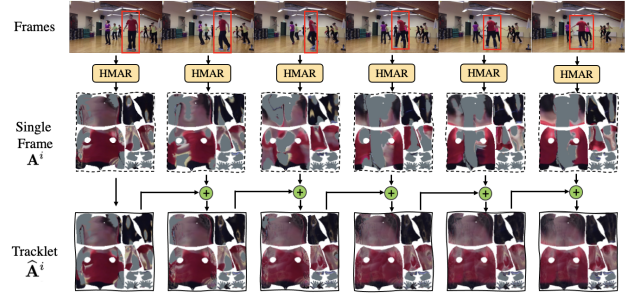For appearance prediction, we make the realistic assumption that human appearance will not change rapidly



Figure 3. **Prediction of appearance:** We show how the single frame appearance $\mathbf{A}^i$ is aggregated over time for the prediction of the tracklet appearance $\widehat{\mathbf{A}}^i$. At the start, we only see the front side of the *person indicated* in the frame, however as he moves his visibility changes, and we only see his back side. With the single frame appearance, we can see that the visibility changes corresponding to the visibility of the person in the frame. However, in the tracklet, the appearance is accumulated over time, and even if the front side is not visible in the last frame, we can see that the tracklet has predicted these regions using the past frames.

over time. Then, the appearance of the tracklet $\widehat{\mathbf{A}}_t^i$ can function as a reasonable prediction for the future appearance of the person. Therefore, we use $\widehat{\mathbf{A}}_t^i$ as the prediction for appearance and use it to measure similarity against a detection in the future frames.

**Location:** Lifting humans from pixels into the 3D space allows us to place them in the global 3D location. Let us assume that a person $i$ at time $t$ has an estimated 3D location $\mathbf{l}_t^i$. Although, we can get an estimate for the location of the person in the global camera frame, this tends to be noisy, particularly along the $z$-axis. To avoid any instabilities when it comes to predicting future location, instead of performing our prediction on the Euclidean $(X, Y, Z)^T$ space, we express our locations in an equivalent $\mathbf{l}_t^i = (x, y, n)^T$ space where $(x, y)$ is the location of the root of the person in the pixel space and $n$ is *nearness*, defined as log inverse depth $n = \log(1/z)$. Nearness is a natural parameterization of depth in multiple vision settings, *e.g.*, [21], because of the $1/z$ scaling of perspective projection. In our case it corresponds to the scale of the human figures that we estimate directly from images. We independently linearly regress the location predictions for $x, y$ and $n$. This is somewhat like the Constant Velocity Assumption (CVA) used in past tracking literature, but there is a subtlety here because constant velocity in 3D need not give rise to constant velocity in 2D (a person would appear to speed up as she approaches the camera). But local linearization is always a reasonable approximation to make, which is what we do.

Let us assume that a tracklet has a set of past locations $\mathcal{L}^i = \{\mathbf{l}_t^i, \mathbf{l}_{t-1}^i, \mathbf{l}_{t-2}^i, ...\}$. Then, the prediction of the loca-

tion for time step $t + 1$ is given by:

$$\widehat{\mathbf{l}}^i_{t+1} = (\widehat{x}^i_{t+1}, \widehat{y}^i_{t+1}, \widehat{n}^i_{t+1})^T$$
$$\text{where, } \widehat{x}^i_{t+1} = \Phi_L(\{x^i_t, x^i_{t-1}, x^i_{t-2}, ..., x^i_{t-w}\}, t+1).$$

Here, $\Phi_L$ is the location aggregation function and we use a simple linear regression for prediction in our tracking algorithm. $\widehat{y}^i_{t+1}$ and $\widehat{n}^i_{t+1}$ are also predicted in a similar fashion. $\Phi_L$ takes the last $w$ observations to fit a line by least squares and regress the future location for $x, y$ and $n$ independently. From the standard theory of linear regression, the prediction interval for $x$ at a time step $t'$ is given by the equation below:

$$\delta_x(t') = t_{(1-\alpha/2)} \times \sqrt{MSE \times \left(1 + \frac{1}{w} + \frac{(t' - \bar{t})^2}{\sum(t - \bar{t})^2}\right)}.$$

Here, $t_{(1-\alpha/2)}$ is the Student's $t$ distribution with confidence $\alpha$ and degree of freedom $w - 2$. $MSE$ is the mean squared error on the predicted locations and $\bar{t}$ is the mean of the time stamps for the previous observations. In a similar manner, we can compute prediction intervals $\Delta_y, \Delta_n$ for $y$ and $n$ respectively.

**Pose:** For the pose pathway, we need to integrate pose information across the tracklet and be able to predict future poses for the near future. To do this, we borrow ideas from the HMMR architecture [15]. Effectively, we learn a function $\Phi_P$ that takes as input a series of pose embeddings of a person $\mathcal{P}^i = \{\mathbf{p}^i_t, \mathbf{p}^i_{t-1}, \mathbf{p}^i_{t-2}, ...\}$ and computes a temporal pose embedding $\widehat{\mathbf{p}}_t$. We train this temporal pose aggregation function $\Phi_P$ to smooth the pose $\widehat{\mathbf{p}}^i_t$ at frame $t$, and regress the future pose representations $\{\widehat{\mathbf{p}}^i_{t+1}, \widehat{\mathbf{p}}^i_{t+2}, ..., \widehat{\mathbf{p}}^i_{t+c}\}$ (typically for up to $c = 12$ frames in the future). We use a transformer [36] to compute $\Phi_P$. This choice allows for some additional flexibility, since sometimes we are not able to detect an identity in some frames (*e.g.*, due to occlusions), which can be handled gracefully by the transformer, by masking out the attention for the representation of the corresponding frame.

### 3.3. Tracking with predicted 3D representations

Given the bounding boxes and their single-frame 3D representations, our tracking algorithm associates identities across frames in an online manner. At every frame, we make future predictions for each tracklet and we measure the similarity with the detected single-frame representation. More specifically, let us assume that we have a tracklet $T_i$, which has been tracked for a sequence of frames and has information for appearance, pose and location. The tracklet predicts its appearance $\widehat{\mathbf{A}}$, location $\widehat{\mathbf{l}}$ and pose $\widehat{\mathbf{p}}$ for the next frame, and we need to measure a similarity score between these predictions of the track $T_i$ and a detection $D_j$ to make an association. Our tracklet representation has three different attributes (appearance, location and pose), so, directly
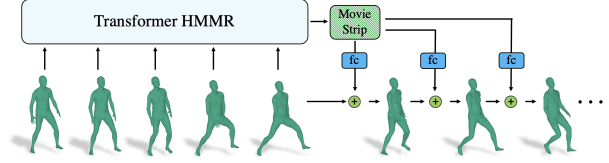


Figure 4. **Prediction of Pose:** We use a modified version of HMMR [15] with transformer backbone. Having transformer as the backbone gives us the flexibility to have missing people in the tracklet (by masking the attention maps), while still allowing us to predictions of future poses. Finally, the transformer give us a movie-strip representation and that is used to regress future poses.
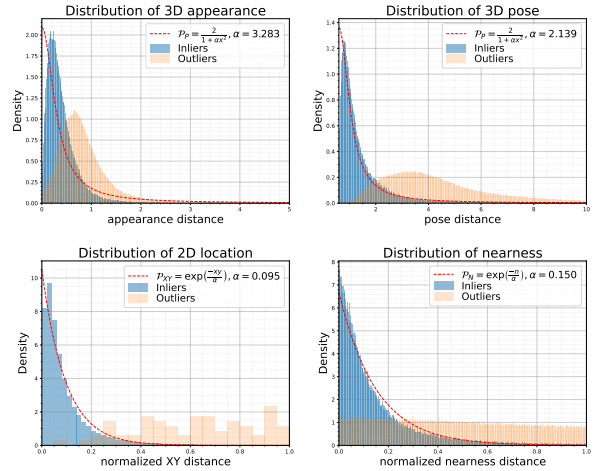


Figure 5. **Conditional distributions of the attribute distances:** We plot the data for the distances between the tracklet prediction and the single frame detection using the ground truth data from PoseTrack [3]. The curves show how the correct matches (inliers) and the incorrect matches (outliers) are distributed. Note that, for 2D location and nearness we plot the distances normalized by the prediction interval.

combining their similarities/distances would not be ideal, since, each attribute has different characteristics. Instead, we investigate the conditional distributions of inliers and outliers of the attributes. Figure 5 presents the corresponding probability distributions for the PoseTrack dataset [3]. The characteristics of these distributions motivate our design decisions for our further modeling.

Assume that tracklet $T_i$ has an appearance representation $\widehat{\mathbf{A}}^i_t$. On the detection side, the detection $D_j$ has a single-frame appearance representation $\mathbf{A}^j_t$. Both of these representations are in the pixel space, therefore we first encode them into an embedding space using the HMAR appearance-encoder network. This gives us an appearance embedding $\widehat{\mathbf{a}}^i_t$ and $\mathbf{a}^j_t$ for the prediction of the tracklet $T_i$ and detection $D_j$, respectively. We are interested in estimating the posterior probability of the event where the detection $D_j$ belongs to the track $T_i$, given some distance measure of

the appearance feature ($\Delta_a$). Assuming that the appearance distance is $\Delta_a = ||\widehat{\mathbf{a}}_t^i - \mathbf{a}_t^j||_2^2$, then the posterior probability is proportional to the conditional probability of the appearance distances, given correct assignments based on Bayes rule. We model this conditional probability as a Cauchy distribution, based on the observations from the inlier distribution of appearance distances (see also Fig 5):

$$\mathcal{P}_A(D_j \in T_i | d_a = \Delta_a) \propto \frac{1}{1 + \beta_a \Delta_a}$$

The distribution has one scaling hyper-parameter $\beta_a$.

Similarly, for pose, we use Cauchy distribution to model the conditional probability of inlier distances. We measure pose distance $\Delta_p = ||\widehat{\mathbf{p}}_t^i - \mathbf{p}_t^j||_2^2$ between the predicted pose representation $\mathbf{p}_t^i$ from the track $T_i$ and the pose representation $\mathbf{p}_t^j$ of detection $D_j$. The posterior probability that the detection belongs to the track, given the pose distance is:

$$\mathcal{P}_P(D_j \in T_i | d_p = \Delta_p) \propto \frac{1}{1 + \beta_p \Delta_p}$$

Here, $\Delta_p = ||\widehat{\mathbf{p}}_t^i - \mathbf{p}_t^j||_2^2$ and $\beta_p$ is the scaling factor.

For location, let us assume the track $T_i$ has predicted a location $\widehat{\mathbf{l}}_t^i = (\widehat{x}_t^i, \widehat{y}_t^i, \widehat{n}_t^i)^T$ with a set of prediction intervals $\{\delta_x, \delta_y, \delta_n\}$, and the detection $D_j$ is at a 3D location $\mathbf{l}_t^j = (x_t^j, y_t^j, n_t^j)^T$. We treat the 3D coordinates $x, y$ and the nearness term $n$ coordinates independently, and compute the posterior probabilities of the detection belongs to the tracklet given the location distance. We model the conditional probability distribution as an exponential distribution, based on the findings from the empirical data. The Fig 5 shows the distribution of 2D distance and nearness distance, scaled by the confidence interval, of inliers approximately follow the exponential distribution.

$$\mathcal{P}_{XY}(D_j \in T_i | d_{xy} = \Delta_{xy}) \propto \frac{1}{\beta_{xy}} \exp\left(\frac{-\Delta_{xy}}{\beta_{xy}\delta_{xy}}\right)$$

Here, $\beta_{xy}$ is a scaling parameter for the exponential distribution, $\Delta_{xy}$ is the 2D pixel distance between the predicted track and the detection and $\delta_{xy} = \sqrt{\delta_x^2 + \delta_y^2}$ is the prediction interval for the 2D location prediction. We use a similar form of exponential distribution for the posterior probability for nearness $\mathcal{P}_N$:

$$\mathcal{P}_N(D_j \in T_i | d_n = \Delta_n) \propto \frac{1}{\beta_n} \exp\left(\frac{-\Delta_n}{\beta_n\delta_n}\right)$$

Here, $\beta_n$ is the scaling parameter for the exponential distribution, $\delta_n$ is the confidence interval for the nearness prediction, and $\Delta_n$ is the $L_1$ distance between the nearness of the tracklet prediction and the detection.

Now that we have computed the conditional probabilities of the detection belonging to a track conditioned on the

---

**Algorithm 1** Tracking Algorithm

1: **procedure** PHALP TRACKING
2: **Require:** All active tracklets $\mathcal{T}$, all detections and their single frame 3D representations at time $t$, $\mathcal{D}$ and maximum age of a track $t_{max}$.
3:     **for** $T_j \in \mathcal{T}$ **do**
4:         # predict all attributes for the next frame.
5:         $\widehat{\mathbf{A}}_t^j \leftarrow \Phi_A(\{\mathbf{A}_{t-1}^j, \mathbf{A}_{t-2}^j, ...\})$
6:         $\widehat{\mathbf{p}}_t^j \leftarrow \Phi_P(\mathbf{p}_{t-1}^j, \mathbf{p}_{t-1}^j, ...\})$
7:         $\widehat{\mathbf{l}}_t^j \leftarrow \Phi_L(\{\mathbf{l}_{t-1}^j, \mathbf{l}_{t-2}^j, ...\})$
8:     # Compute the cost between tracks and detections.
9:     $\mathbf{C}_{i,j} \leftarrow \Phi_C(D_i, T_j)$ for all $D_i \in \mathcal{D}$ and $T_j \in \mathcal{T}$
10:     # Hungarian to assign detections to tracklets.
11:     $\mathcal{M}, \mathcal{T}_u, \mathcal{D}_u \leftarrow \texttt{Assignment}(\mathbf{C})$
12:     # Update the matched tracks.
13:     $\mathcal{T} \leftarrow \{T_j(D_i), \ \forall(i,j) \in \mathcal{M}\}$
14:     # Increase the age of unmatched tracks.
15:     $\mathcal{T} \leftarrow \{T_j(age)+ = 1, \ \forall(j) \in \mathcal{T}_u\}$
16:     # Make new tracks with unmatched detections.
17:     $\mathcal{T} \leftarrow \{T_j(D_i), \ \forall(i) \in \mathcal{D}_u, \ j = |\mathcal{T} + 1|\}$
18:     Kill the tracks with age $\geq t_{max}$.
19: **return** Tracklets $\mathcal{T}$

---

individual cues of appearance, location and pose, we can compute the overall conditional probability of the detection $D_j$ belonging to the track $T_i$, given all the cues together (assumed to be independent):

$$\mathcal{P}(D_j \in T_i | \Delta_a, \Delta_p, \Delta_{xy}, \Delta_n) \propto \mathcal{P}_A \mathcal{P}_P \mathcal{P}_{XY} \mathcal{P}_N$$

This allow us to estimate how probable an association is based on various attribute distances. Finally, we map the similarity measures (probability values up to a scale), to cost values, for solving association. The cost function between the detection representations and a predicted representations of the tracklet is defined as:

$$\begin{aligned}\Phi_C(D_j, T_i) &= -\log(\mathcal{P}(D_j \in T_i)) \\ &= -\log(\mathcal{P}_A) - \log(\mathcal{P}_P) - \log(\mathcal{P}_{XY}) - \log(\mathcal{P}_N),\end{aligned}$$

where the second equality is up to an additive constant. Once the cost between all the tracks and the detection is computed, we simply pass it to the Hungarian algorithm for solving the association.

**Estimating the parameters of the cost function:** The cost function $\Phi_C$ has 4 parameters ($\beta_a, \beta_p, \beta_{xy}$ and $\beta_n$). Additionally, the Hungarian algorithm has one parameter $\beta_{th}$ to decide whether the track is not a match to the detection. Therefore, overall we have five parameters for the whole association part of our tracking system. Now, we treat this as an empirical risk minimization problem and optimize the $\beta$ values based on a loss function. We initialize

$\beta_a, \beta_p, \beta_{xy}$ and $\beta_n$ with the values from the estimated density functions and use frame level association error as a loss function for the optimization. We use the Nelder–Mead [30] algorithm for this optimization. Finally, the optimized $\beta$ values are used for the cost function across all the datasets, and a simple tracking algorithm is used to associate detections with tracklet predictions. The sketch of the tracking algorithm is shown in Algorithm 1.

### 3.4. Extension to shot changes

Our framework can easily be extended to also handle shot changes, which are common in edited media, like movies, TV shows, but also sports. Since shot changes can be detected relatively reliably, we use an external shot detector [13] to identify frames that indicate shot changes. Informed by the detection of this boundary, during tracking, we update the distance metric accordingly. More specifically, since appearance and 3D pose are invariant to the viewpoint, we keep these factors in the distance computation, while we drop the location distance from the distance metric, because of the change in the camera location. Then, the association is computed based on this updated metric. We use the AVA dataset [10] to demonstrate this utility of our tracking system and present results in Section 4.

### 4. Experiments

In this section, we present the experimental evaluation of our approach. We report results on three datasets: PoseTrack [3], MuPoTS [28] and AVA [10], which capture a diverse set of sequences, including sports, casual interactions and movies. Our method operates on detections and masks coming from an off-the-shelf Mask-RCNN network [12], and returns the identity label for each one of them. Therefore, the metrics we use to report results also focus on identity tracking at the level of the bounding box. More specifically, we report results using Identity switches (IDs), Multi-Object Tracking Accuracy (MOTA) [16], ID F1 score (IDF1) [33] and HOTA [26]. In all cases, we adopt the protocols of Rajasegaran *et al*. [32] for evaluation.

First, we ablate the main components of our approach. Specifically, we investigate the effect of each one of the tracking cues we employ, *i.e*., appearance, 3D location and 3D pose, and how they affect the overall tracking pipeline. For this comparison, we report results on the Posetrack dataset [3]. The full results are presented in Table 1. As we can see, removing each one of the main cues leads to degradation in the performance of the system, where 3D location seems to have the largest effect on the performance, followed by appearance and 3D pose. Moreover, this ablation also highlights the importance of having the nearness term in the cost function, a feature that is not available to purely 2D tracking methods.

| Method | PoseTrack | | |
|---|---|---|---|
| | IDs↓ | MOTA↑ | IDF1↑ |
| w/o 3D appearance | 632 | 58.4 | 74.9 |
| w/o 3D pose | 558 | 58.9 | 76.2 |
| w/o location | 948 | 57.3 | 71.6 |
| w/o nearness | 622 | 58.5 | 74.8 |
| Full system | **541** | **58.9** | **76.4** |

Table 1. **Ablation of the main components of PHALP on Pose-Track [3].** Removing each tracking cue (3D appearance, 3D pose or 3D location) leads to degradation in the performance.

Next, we evaluate our approach in comparison with the state-of-the-art methods. The results are presented in Table 2. We report results on PoseTrack [3], MuPoTS [28] and AVA [10]. Our method outperforms the previous baselines, as well as the state-of-the-art approach of Rajasegaran [32]. The gains are significant across all metrics. Our method also outperforms the other approaches in the HOTA metric.

Finally, we also show qualitative results of our method on multiple datasets in Fig 6. These results show that our method performs reliably even in very hard occlusion cases, while it is able to recover the correct identity over multiple successive occlusions. Fig 6 also shows the robustness of our method in complex motion sequences, shot changes and long trajectories.

### 5. Discussion

We presented PHALP, an approach for monocular people tracking, by predicting appearance, location and pose in 3D. Our method relies on a powerful backbone for 3D human mesh recovery, modeling on the tracklet level for collecting information across the tracklet's detections, and eventually predicting the future states of the tracklet. One of the main benefits of PHALP is that the association aspect requires tuning of only five parameters, which makes it very friendly for training on multi-object tracking datasets, where annotating the identity of every person in a video can be expensive. We should note that our approach can be naturally extended to make use of more attributes, *e.g*., a face embedding, which could be useful for cases with close-ups, like movies. The main assumptions for PHALP are that we have access to a good object detector for the initial bounding box/mask detection, and a strong HMAR network for single-frame lifting of people to 3D. If the performance of these components is not satisfactory, it can also affect PHALP. Regarding societal impact, tracking systems have often been used for human surveillance. We do not condone such use. Instead, we believe that a tracking system will be valuable for studying social-human interactions.

| Method | Posetrack | | | | MuPoTS | | | | AVA | |
|---|---|---|---|---|---|---|---|---|---|---|
| | IDs↓ | MOTA↑ | IDF1↑ | HOTA↑ | IDs↓ | MOTA↑ | IDF1↑ | HOTA↑ | IDs↓ | IDF1↑ |
| Trackformer [29] | 1263 | 33.7 | 64.0 | 46.7 | 43 | 24.9 | 62.7 | 53.2 | 716 | 40.9 |
| Tracktor [4] | 702 | 42.4 | 65.2 | 38.5 | 53 | 51.5 | 70.9 | 50.3 | 289 | 46.8 |
| AlphaPose [8] | 2220 | 36.9 | 66.9 | 37.6 | 117 | 37.8 | 67.6 | 41.8 | 939 | 41.9 |
| FlowPose [39] | 1047 | 15.4 | 64.2 | 38.0 | 49 | 21.4 | 67.1 | 43.0 | 452 | 52.9 |
| T3DP [32] | 655 | 55.8 | 73.4 | 50.6 | 38 | 62.1 | 79.1 | 59.2 | 240 | 61.3 |
| PHALP | **541** | **58.9** | **76.4** | **52.9** | **22** | **66.2** | **81.4** | **59.4** | **227** | **62.7** |

Table 2. **Comparison with state-of-the-art tracking methods.** We compare our method PHALP with various tracking methods in three different datasets. Our approach outperforms the other baselines across all datasets and metrics.
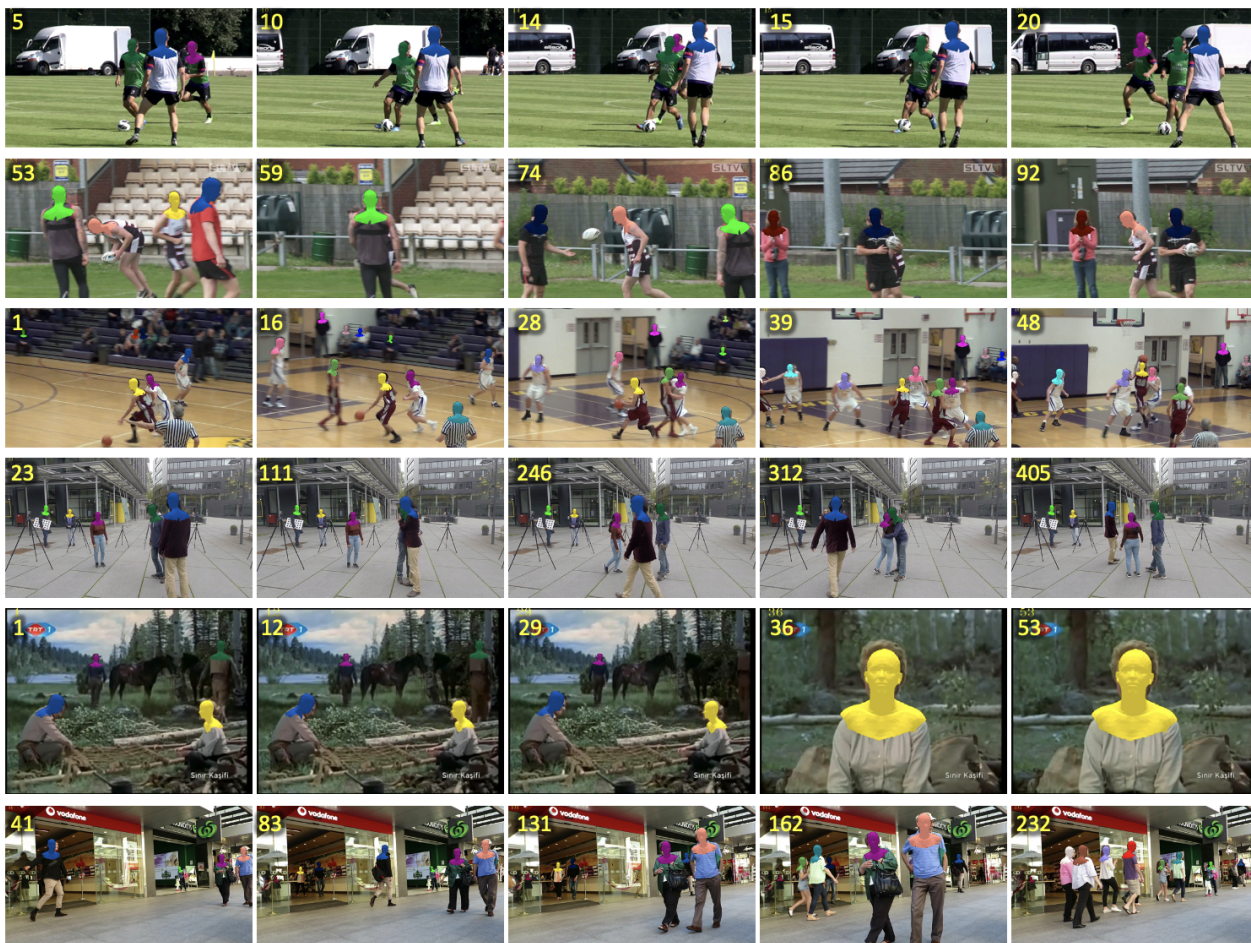


Figure 6. **Qualitative Results:** We show the tracking performance of PHALP in various datasets (frame number is shown at the top left corner). The first three rows are from the PoseTrack dataset [3]. These results show that even during successive occlusions our method is able to track the identity of the correct person. Note that, in the first row, although the green head-mask person and the purple head-mask person have similar appearance, our method can track each one of them successfully. In the second row, the player is going through multiple occlusions, yet recovered correctly. The third row shows the robustness of our linearization approximation for 3D location prediction, even when the motions of the players are very complex. In the MuPoTS dataset [28] (4th row), our method can handle very close interactions between people. This is due to the fact that, our modification of HMAR recovers meshes conditioned on the detected mask. We also show results (5th row) on the AVA dataset [10]. After the 3rd frame, there is a shot change in the video, and the woman is tracked successfully across the shots. Finally, we show qualitative results on a MOT17 sequence. The blue person is tracked for the whole sequence while he is going through multiple occlusions for a long time. More results at the PHALP website.

# References

[1] Emre Aksan, Peng Cao, Manuel Kaufmann, and Otmar Hilliges. A spatio-temporal transformer for 3D human motion prediction. In *3DV*, 2020. 3

[2] Alexandre Alahi, Kratarth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. Social LSTM: Human trajectory prediction in crowded spaces. In *CVPR*, 2016. 3

[3] Mykhaylo Andriluka, Umar Iqbal, Eldar Insafutdinov, Leonid Pishchulin, Anton Milan, Juergen Gall, and Bernt Schiele. PoseTrack: A benchmark for human pose estimation and tracking. In *CVPR*, 2018. 5, 7, 8

[4] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe. Tracking without bells and whistles. In *ICCV*, 2019. 3, 8

[5] Christoph Bregler and Jitendra Malik. Tracking people with twists and exponential maps. In *CVPR*, 1998. 3

[6] Gioele Ciaparrone, Francisco Luque Sánchez, Siham Tabik, Luigi Troiano, Roberto Tagliaferri, and Francisco Herrera. Deep learning in video multi-object tracking: A survey. *Neurocomputing*, 381:61–88, 2020. 3

[7] Patrick Dendorfer, Aljosa Osep, Anton Milan, Konrad Schindler, Daniel Cremers, Ian Reid, Stefan Roth, and Laura Leal-Taixé. MOTChallenge: A benchmark for single-camera multiple target tracking. *IJCV*, 129(4):845–881, 2021. 3

[8] Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. RMPE: Regional multi-person pose estimation. In *ICCV*, 2017. 8

[9] Rohit Girdhar, Georgia Gkioxari, Lorenzo Torresani, Manohar Paluri, and Du Tran. Detect-and-track: Efficient pose estimation in videos. In *CVPR*, 2018. 3

[10] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, Cordelia Schmid, and Jitendra Malik. AVA: A video dataset of spatio-temporally localized atomic visual actions. In *CVPR*, 2018. 2, 7, 8

[11] Peng Guan, Alexander Weiss, Alexandru O Balan, and Michael J Black. Estimating human shape and pose from a single image. In *ICCV*, 2009. 3

[12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017. 3, 7

[13] Qingqiu Huang, Yu Xiong, Anyi Rao, Jiaze Wang, and Dahua Lin. MovieNet: A holistic dataset for movie understanding. In *ECCV*, 2020. 7

[14] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018. 3

[15] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3D human dynamics from video. In *CVPR*, 2019. 2, 3, 5

[16] Rangachar Kasturi, Dmitry Goldgof, Padmanabhan Soundararajan, Vasant Manohar, John Garofolo, Rachel Bowers, Matthew Boonstra, Valentina Korzhova, and Jing Zhang. Framework for performance evaluation of face, text, and vehicle detection and tracking in video: Data, metrics, and protocol. *PAMI*, 2008. 7

[17] Tarasha Khurana, Achal Dave, and Deva Ramanan. Detecting invisible people. In *ICCV*, 2021. 3

[18] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. VIBE: Video inference for human body pose and shape estimation. In *CVPR*, 2020. 3

[19] Muhammed Kocabas, Chun-Hao P Huang, Otmar Hilliges, and Michael J Black. PARE: Part attention regressor for 3D human body estimation. In *ICCV*, 2021. 3

[20] Muhammed Kocabas, Chun-Hao P Huang, Joachim Tesch, Lea Muller, Otmar Hilliges, and Michael J Black. SPEC: Seeing people in the wild with an estimated camera. In *ICCV*, 2021. 3

[21] Jan J Koenderink. Optic flow. *Vision research*, 26(1):161–179, 1986. 4

[22] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3D human pose and shape via model-fitting in the loop. In *ICCV*, 2019. 3

[23] Nikos Kolotouros, Georgios Pavlakos, Dinesh Jayaraman, and Kostas Daniilidis. Probabilistic modeling for human mesh recovery. In *ICCV*, 2021. 3

[24] Oh-Hun Kwon, Julian Tanke, and Juergen Gall. Recursive bayesian filtering for multiple human pose tracking from multiple cameras. In *ACCV*, 2020. 3

[25] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *ACM Transactions on Graphics (TOG)*, 34(6):1–16, 2015. 3

[26] Jonathon Luiten, Aljosa Osep, Patrick Dendorfer, Philip Torr, Andreas Geiger, Laura Leal-Taixé, and Bastian Leibe. HOTA: A higher order metric for evaluating multi-object tracking. *IJCV*, 2021. 7

[27] Zhengyi Luo, S Alireza Golestaneh, and Kris M Kitani. 3D human motion estimation via motion compression and refinement. In *ACCV*, 2020. 3

[28] Dushyant Mehta, Oleksandr Sotnychenko, Franziska Mueller, Weipeng Xu, Srinath Sridhar, Gerard Pons-Moll, and Christian Theobalt. Single-shot multi-person 3D pose estimation from monocular RGB. In *3DV*, 2018. 7, 8

[29] Tim Meinhardt, Alexander Kirillov, Laura Leal-Taixe, and Christoph Feichtenhofer. Trackformer: Multi-object tracking with transformers. *arXiv preprint arXiv:2101.02702*, 2021. 3, 8

[30] John A Nelder and Roger Mead. A simplex method for function minimization. *The computer journal*, 7(4):308–313, 1965. 7

[31] Georgios Pavlakos, Jitendra Malik, and Angjoo Kanazawa. Human mesh recovery from multiple shots. *arXiv preprint arXiv:2012.09843*, 2020. 3

[32] Jathushan Rajasegaran, Georgios Pavlakos, Angjoo Kanazawa, and Jitendra Malik. Tracking people with 3D representations. In *NeurIPS*, 2021. 1, 2, 3, 7, 8

[33] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In *ECCV*, 2016. 7

[34] Michael Snower, Asim Kadav, Farley Lai, and Hans Peter Graf. 15 keypoints is all you need. In *CVPR*, 2020. 3

[35] Pavel Tokmakov, Jie Li, Wolfram Burgard, and Adrien Gaidon. Learning to track with object permanence. In *ICCV*, 2021. 3

[36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017. 3, 5

[37] Xinshuo Weng, Yongxin Wang, Yunze Man, and Kris M Kitani. GNN3DMOT: Graph neural network for 3D multi-object tracking with 2D-3D multi-feature learning. In *CVPR*, 2020. 3

[38] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *ECCV*, 2018. 3

[39] Yuliang Xiu, Jiefeng Li, Haoyu Wang, Yinghong Fang, and Cewu Lu. Pose Flow: Efficient online pose tracking. *arXiv preprint arXiv:1802.00977*, 2018. 8

[40] Jiarui Xu, Yue Cao, Zheng Zhang, and Han Hu. Spatial-temporal relation networks for multi-object tracking. In *ICCV*, 2019. 3

[41] Yihong Xu, Aljosa Osep, Yutong Ban, Radu Horaud, Laura Leal-Taixé, and Xavier Alameda-Pineda. How to train your deep multi-object tracker. In *CVPR*, 2020. 3

[42] Alper Yilmaz, Omar Javed, and Mubarak Shah. Object tracking: A survey. *Acm computing surveys (CSUR)*, 38(4):13–es, 2006. 3

[43] Jason Y Zhang, Panna Felsen, Angjoo Kanazawa, and Jitendra Malik. Predicting 3D human dynamics from video. In *ICCV*, 2019. 3

[44] Yuxiang Zhang, Liang An, Tao Yu, Xiu Li, Kun Li, and Yebin Liu. 4D association graph for realtime multi-person motion capture using multiple video cameras. In *CVPR*, 2020. 3

[45] Yang Zhang, Hao Sheng, Yubin Wu, Shuai Wang, Weifeng Lyu, Wei Ke, and Zhang Xiong. Long-term tracking with deep tracklet association. *IEEE Transactions on Image Processing*, 29:6694–6706, 2020. 3