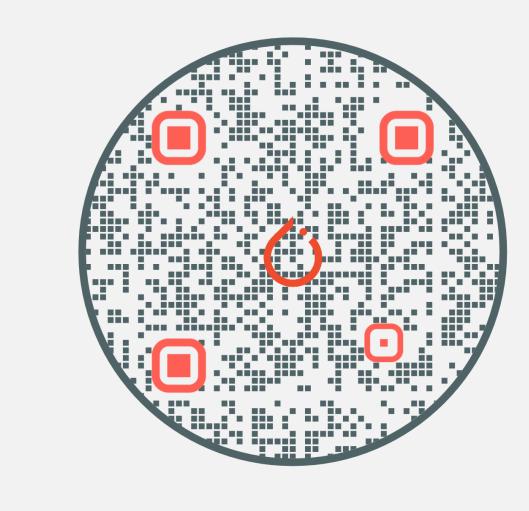


Tracking People by Predicting 3D Appearance, Location and Pose

Jathushan Rajasegaran, Georgios Pavlakos, Angjoo Kanazawa, Jitendra Malik

UC Berkeley





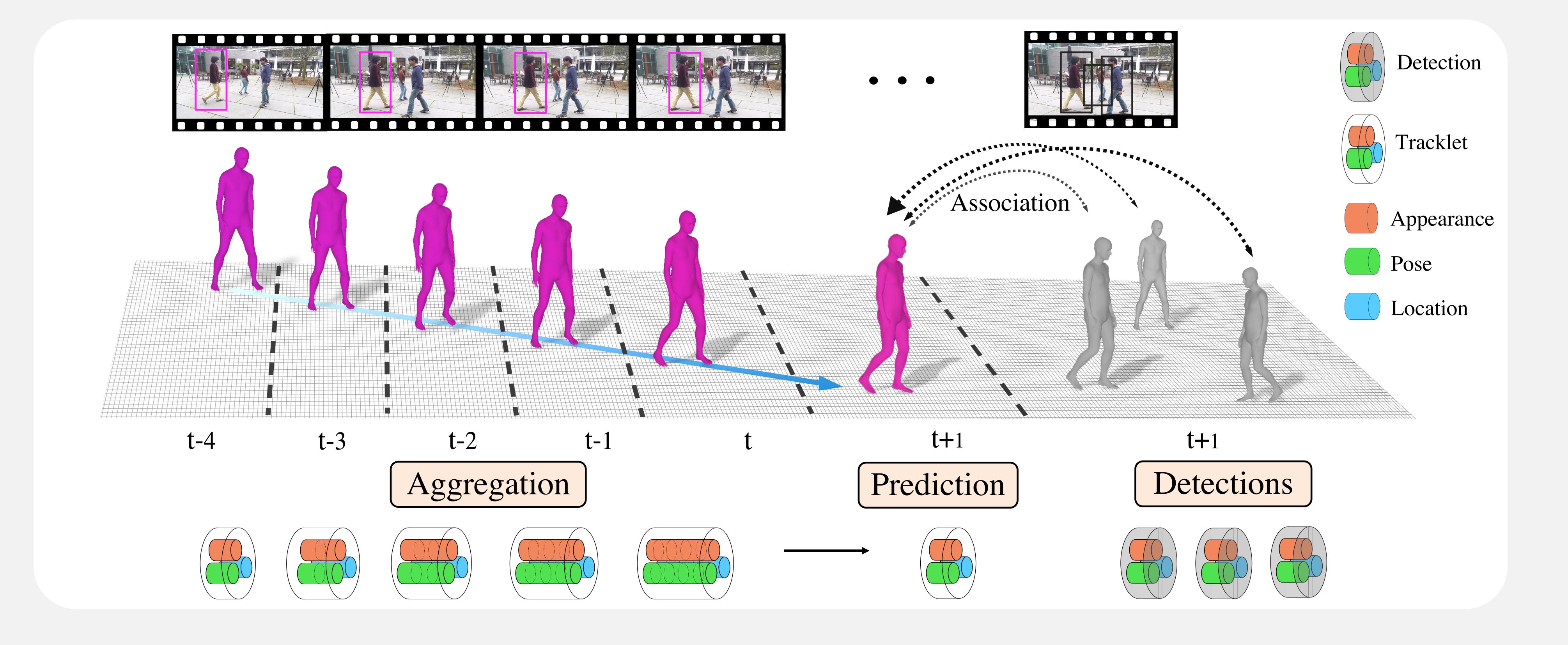
Introduction

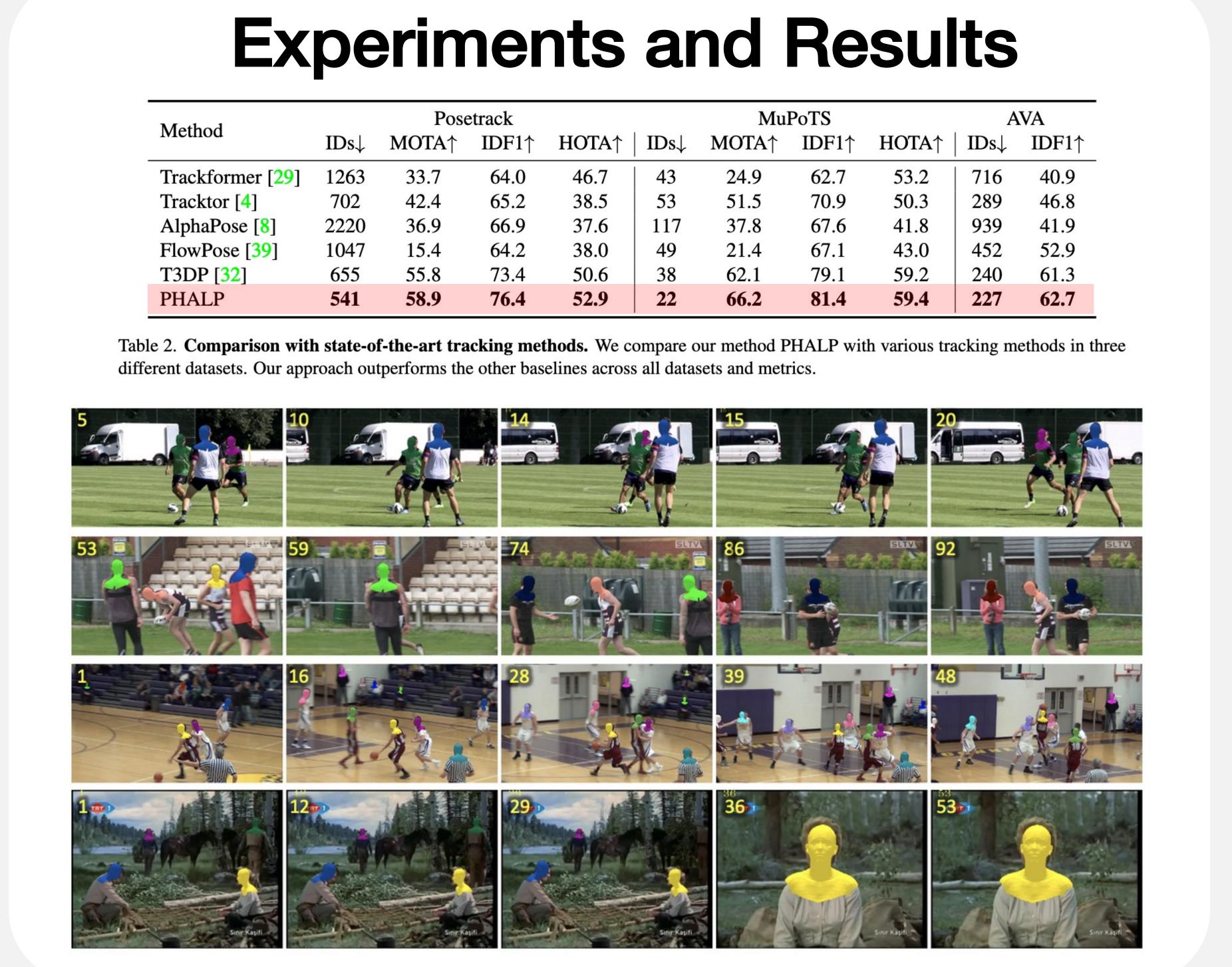
- We track people in monocular videos.
- we lift every people in a video to 3D.

 Rajasegaran et al.

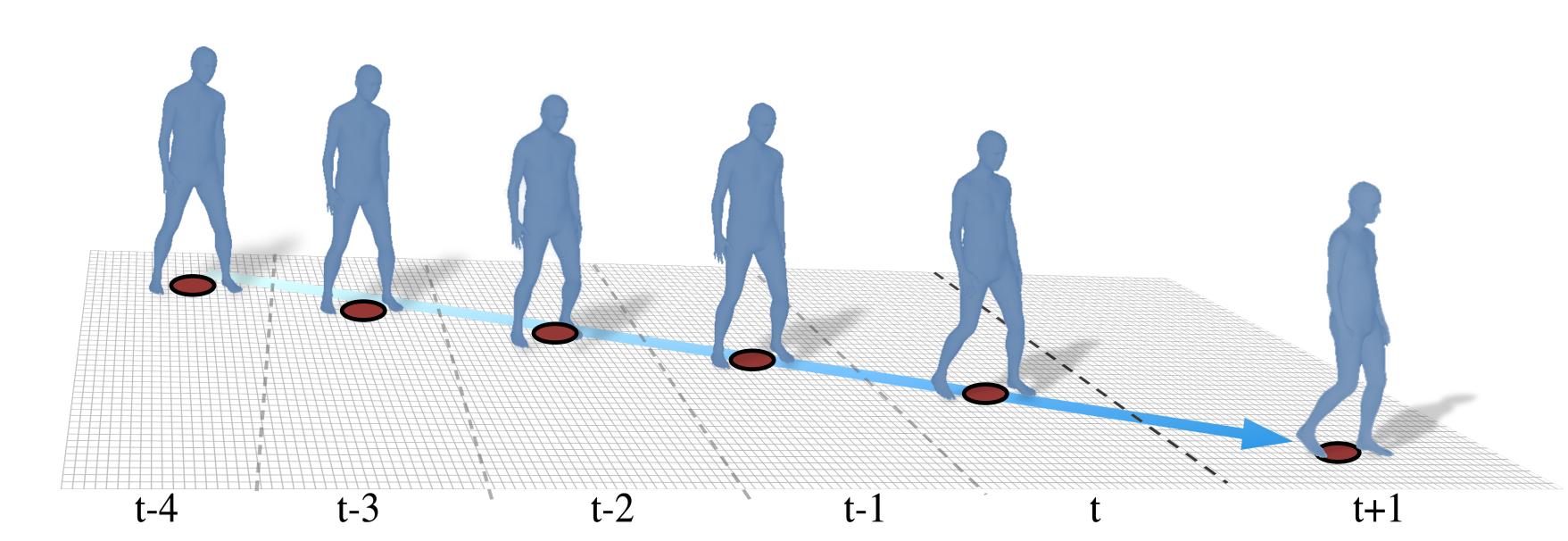
 NeurIPS 2021
- We aggregate the states and make predictions.







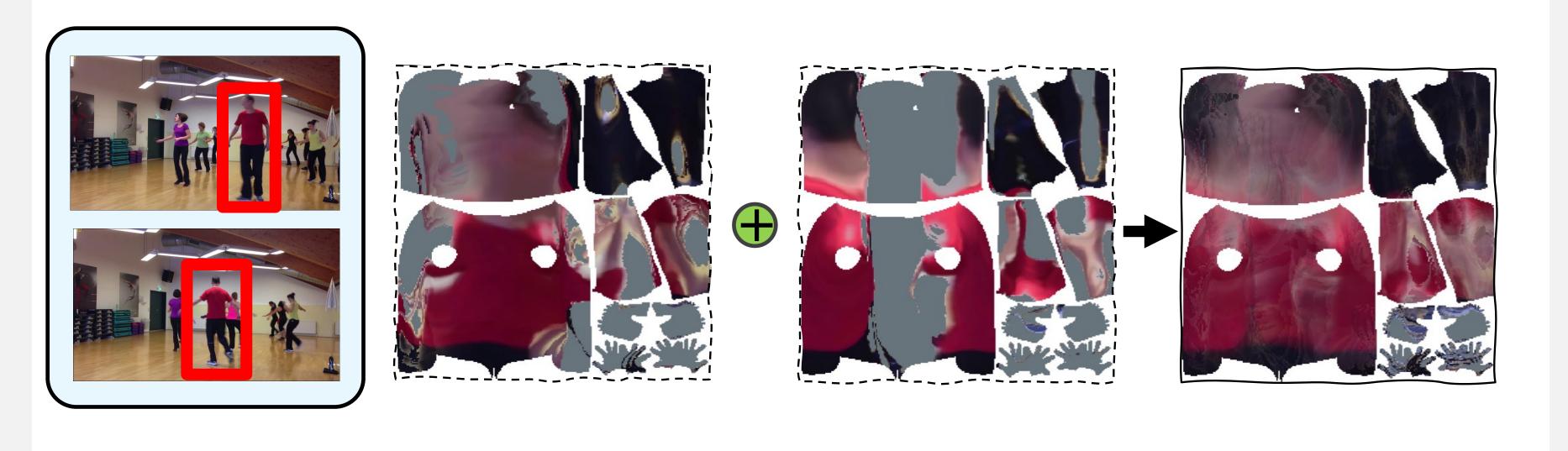
Location Prediction



- Human motion is linear.
- Place every human in X, Y and N.
- nearness = 1/scale.

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

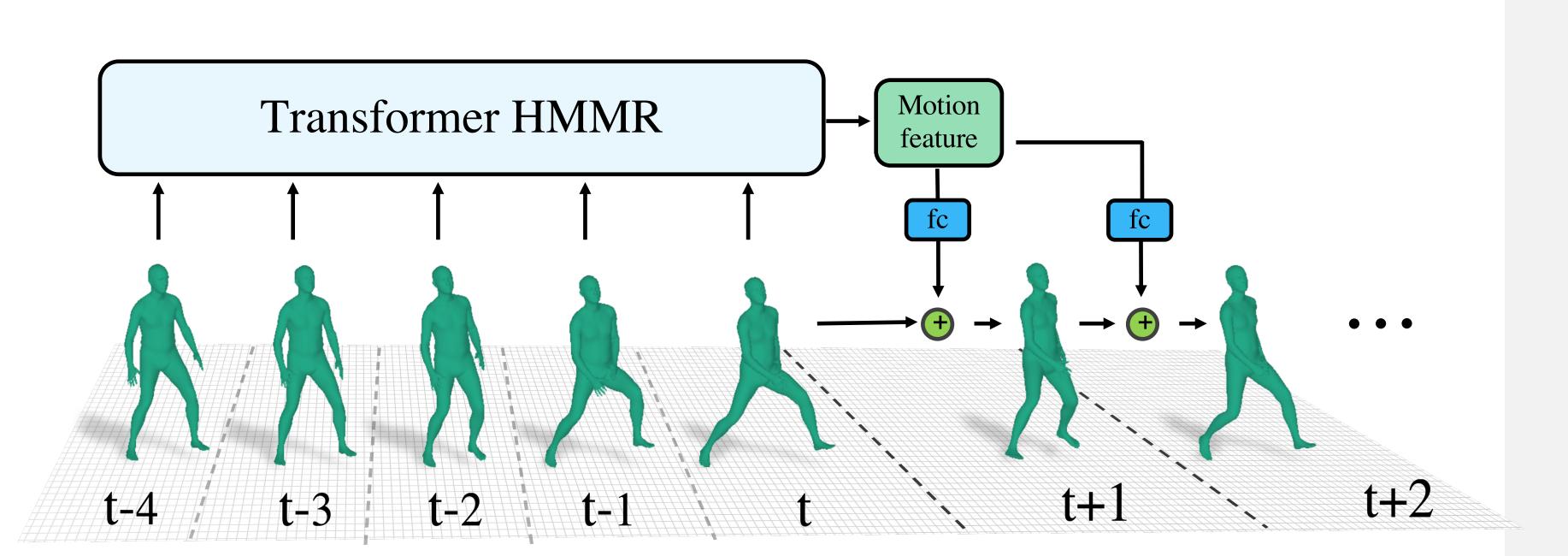
Appearance Prediction



- Appearance is constant over time
 Exponential moving average to aggregate.
- ullet Δ_{α} is the distance between prediction and detection.

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

Pose Prediction

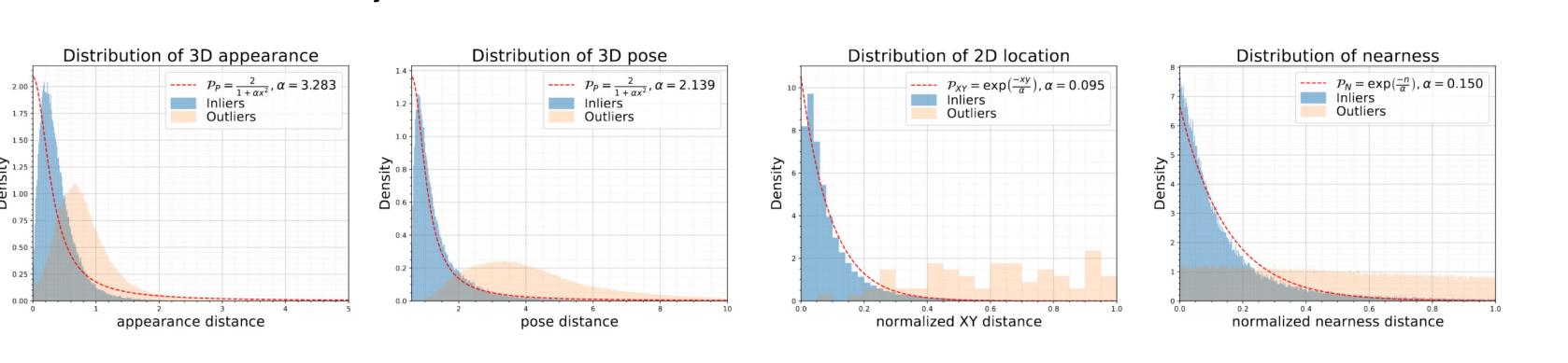


- Human Mesh and Motion Recovery to predict human poses. Kanazawa et al. CVPR 2019
- Pose distribution is modeled as a Cauchy distribution.

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

Association

Inlier, outlier distribution for appearance, pose,
2D location, and nearness.



 Total cost is the negative log likelihood of the joint distribution.

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

$$\Phi_C(D_j, T_i) = -\log(\mathcal{P}(D_j \in T_i))$$

$$= -\log(\mathcal{P}_A) - \log(\mathcal{P}_P) - \log(\mathcal{P}_{XY}) - \log(\mathcal{P}_N),$$

Hungarian to solve association.