



LesseTrack: End-to-End Learnable Multi-Person Articulated 3D Pose Tracking | Pose | Po





¹Carnegie Mellon University

*Equal Contribution ²Amazon

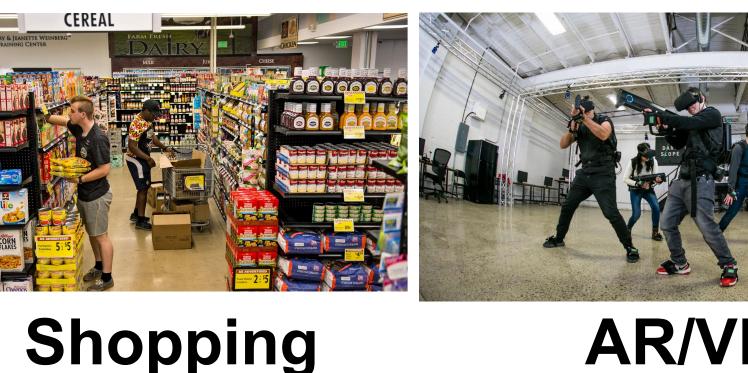
Project: http://www.cs.cmu.edu/~ILIM/projects/IM/TesseTrack/

Introduction

Goal: Multi-Person 3D Pose Tracking using End-to-End Learning









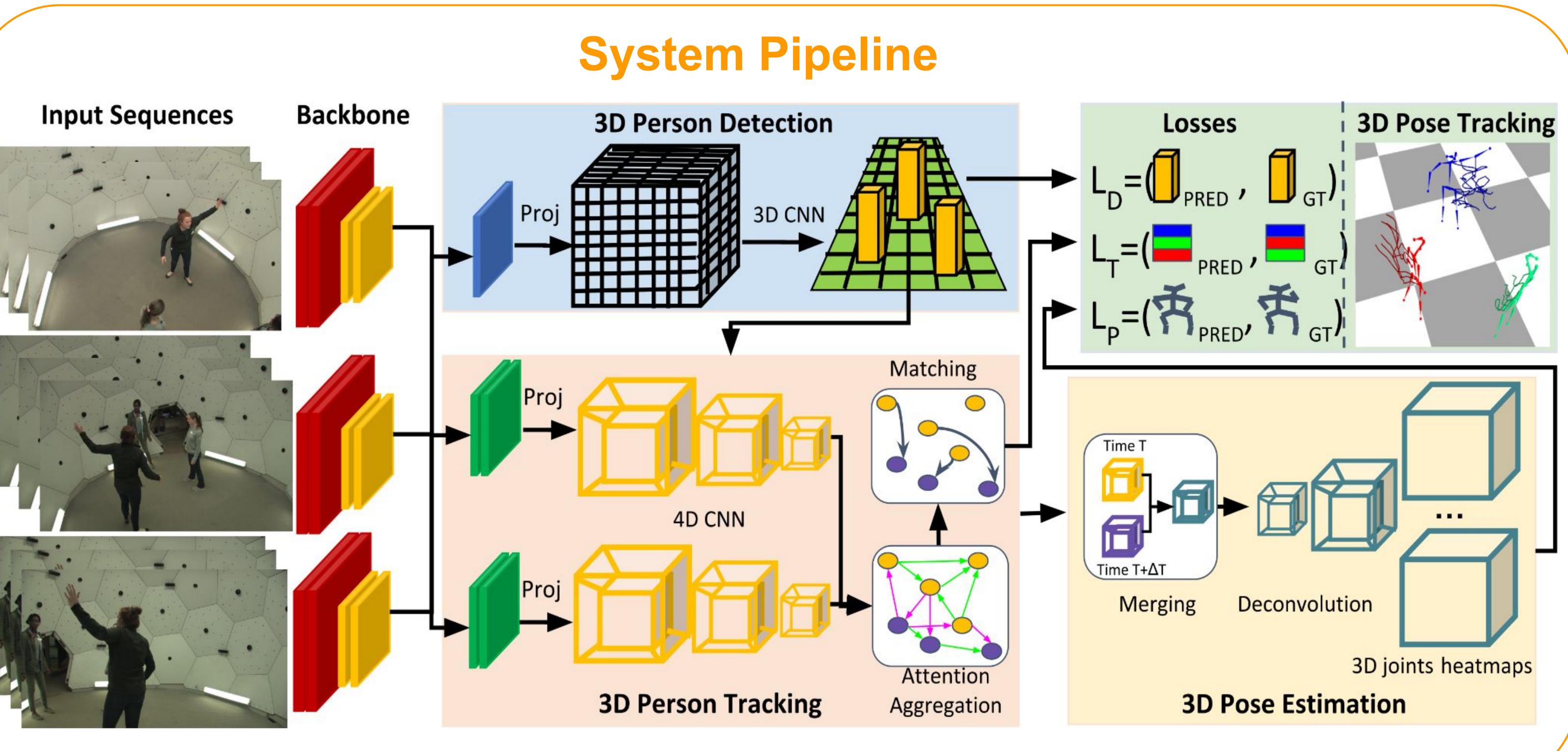
Input image + estimated 2D pose

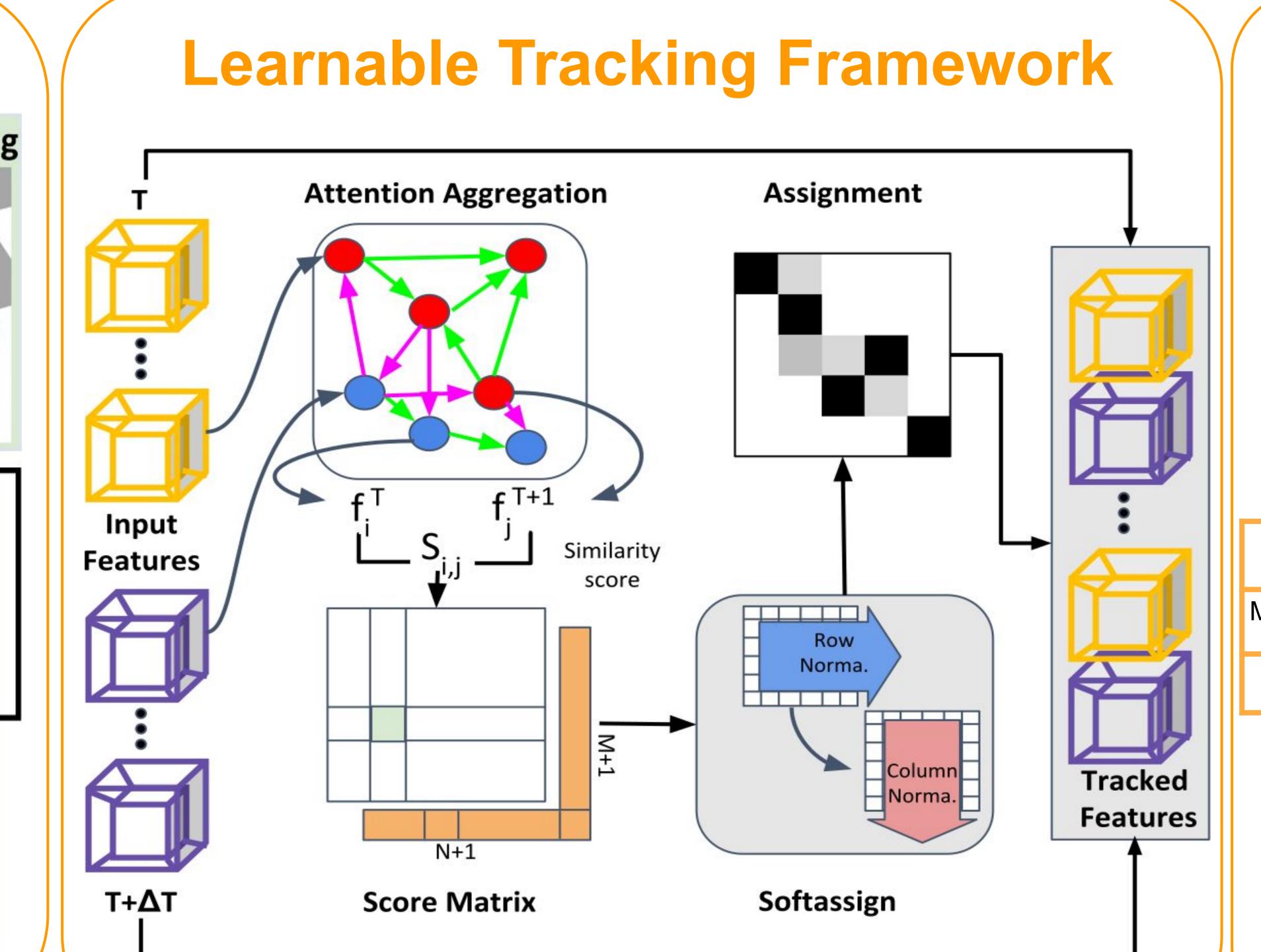
AR/VR **Automation**

Contributions:

Sports

- A novel spatio-temporal formulation that allows simultaneous 3D body joint reconstruction and tracking of multiple individuals.
- A novel learnable tracking formulation that allows extending person-specific spatio-temporal representation learning.
- A novel evaluation framework for multi-person 3D pose tracking.
- An in-depth ablation study of the proposed approach and thorough comparisons to current methods on several standard benchmarks.





SOTA Comparisions

ethod	Actor-1	Actor-2	Actor-3	Total	Method	Actor-1	Actor-2	Actor-3	Total
giannis	93.5	75.7	84.4	84.5	Belagiannis	93.5	75.7	84.4	84.5
shadi	94.2	92.9	84.6	90.6	Ershadi	94.2	92.9	84.6	90.6
ong	97.2	93.3	98.0	96.3	Dong	97.2	93.3	98.0	96.3
Et al.	97.6	93.8	98.8	96.7	Tu Et al.	97.6	93.8	98.8	96.7
seTrack	97.9	95.2	99.1	97.4	TesseTrack	97.9	95.2	99.1	97.4

Campus Dataset

Shelf Dataset

	Monocular Method(MPJPE,mm)					Multi-View Method(MPJPE,mm)					
Method	Martinez	Iskakov	Pavollo	Cheng	Tessetrack	Martinez	Pavlakos	Padoy	Iskakov	TesseTrack	
All	62.9	49.9	46.8	40.1	44.6	57.0	56.7	49.1	20.8	18.7	

Human3.6M Dataset

	Multi-Vie	ew(5 Views)	Monocular		
Method	Tu et al.	TesseTrack	Tu Et al.	TesseTrack	
MPJPE(mm)	17.7	7.3	51.1	18.9	

Panoptic Dataset

State-of-the Art

Single person 3D Pose Estimation:

Good Reconstruction accuracy

Francesc Moreno-Noguer. 3d human pose estimation from a single image via distance matrix regression.. In CVPR 2017. Julieta Martinez, Rayat Hossain, et al.. A simple yet effective baseline for 3d hu-man pose estimation. In ICCV, 2017.

Fails with Multi-person or occlusions

Multi-Person 3D Pose Estimation:

Good 3D learning Frameworks

Fails Accurate Tracking

Hanyue Tu, et al. Voxelpose:Towards multi-camera 3d human pose estimation in wild environment. ECCV, 2020. Je Sun, Bin Xiao, et al.. Deephigh-resolution representation learning for human pose esti-mation. In CVPR, 2019

Multi-Person 3D Pose Tracking:

Working tracking Frameworks

PieceWise Tracking modules

Andrei Zanfir et al.. Monocular 3d pose and shape estimation of multiple people in natural scenes.in cvpr 2018

Mvkhavlo Andriluka, Et al.. Monocular 3d pose estimation and tracking

Losses

Detection Loss:

$$L_D^t = \sum_{w=1}^{N} \sum_{h=1}^{M} \sum_{d=1}^{N} ||V_{Pred}^{w,h,d} - V_{GT}^{w,h,d}||$$

Tracking Loss:

$$L_T^t = -\sum_{(i,j)\in G} \log P_{i,j}$$

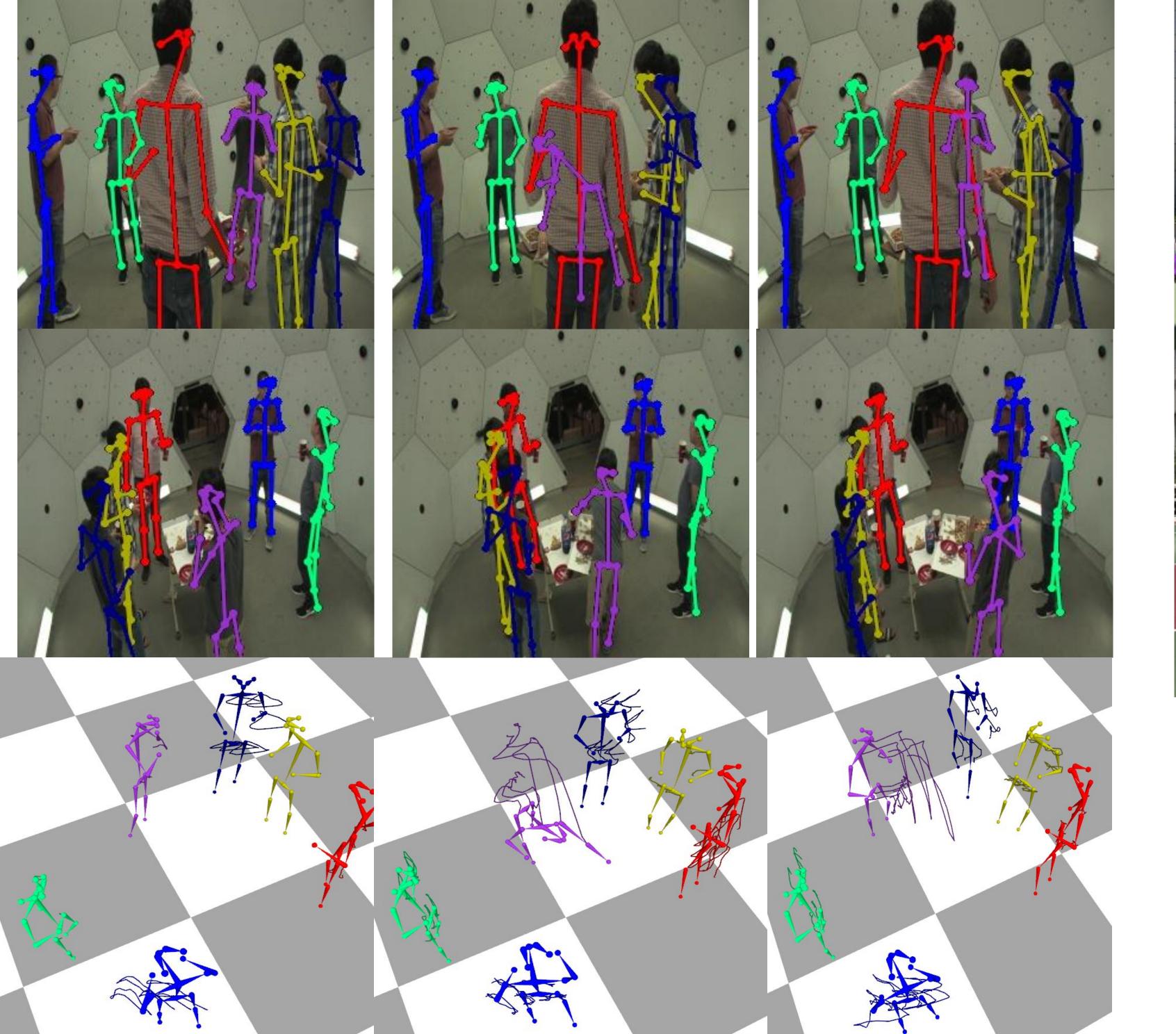
3D Pose Estimation Loss:

$$L_P^{t,d} = \sum_{q=1}^{Q} [||k_{Pred}^q - k_{GT}^q||_1 - \beta.\log(T_{Pred}^q(k_{GT}^q))]$$

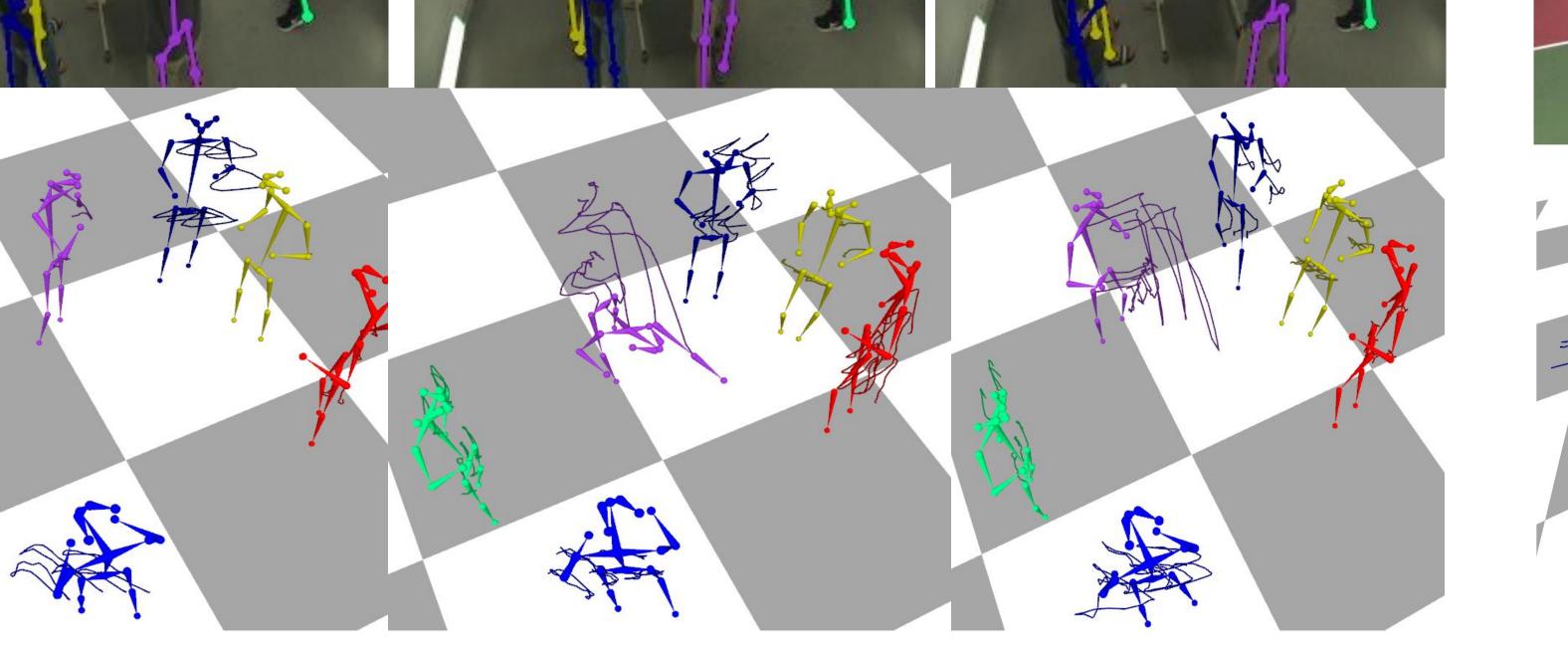
Total Loss:

$$L = \sum_{t \in D} \left[L_D^t + \alpha L_T^t + \gamma \sum_{p \in TP(t)} L_P^{t,p} \right]$$

Qualitative Results



Panoptic Dataset



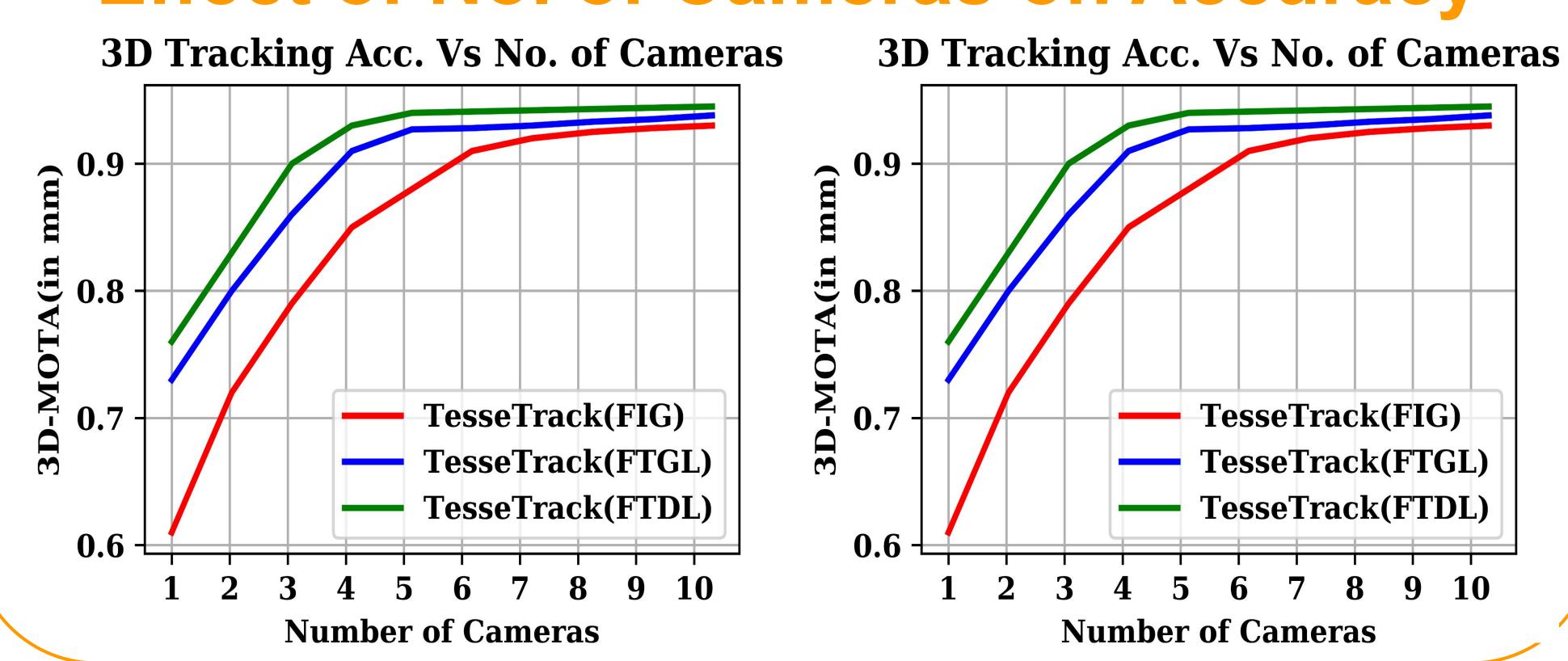
Wild Dataset

3D Pose Tracking Accuracy(3D-MOTA)

Method	Neck	head	Shou.	Elbow	Wrist	Hip	Knee	Ankle	Avg
FIG	89.7	87.4	90.8	0.88	82.2	92.7	89.1	92.4	87.6
FTGL	93.9	91.7	93.0	92.1	87.4	94.4	93.9	94.6	92.1
FTDL	94.6	93.6	93.4	92.7	88.2	94.7	93.8	95.0	94.1

Panoptic Dataset

Effect of No. of Cameras on Accuracy



This paper was supported in parts by NSF Grants IIS1900821 and CNS-2038612, DOT RITA Mobility-21 Grant 69A3551747111, and a PhD fellowship from Amazon.