

I. Multi-Object Tracking (MOT)

Multi-object tracking. Continuous detection and identification of multiple objects in a video.

Tracking-by-detection. For each frame, candidate each frame, candidate bounding boxes are given and associated with tracks.



2. Problem statement

Pose, depth maps, or thermal data can encode a deeper and robust understanding of the motion and the scene.



How to fuse costs from multiple domains into a single association metric?

Existing works [4, 3] such as Deep SORT exploit handwritten formulas or *heuristics*:

 $c_{i,j} = \lambda d^{(1)}(i,j) + (1-\lambda)d^{(2)}(i,j).$

The existing fusing strategies: assume the costs to be independent weigh the costs always the same, overlooking instance- and scene-related peculiarities (e.g., camera motion, lighting conditions, etc.)

Goal. Something like this:



TrackFlow: Multi-Object Tracking with Normalizing Flows Gianluca Mancusi¹, Aniello Panariello¹, Angelo Porrello¹, Matteo Fabbri², Simone Calderara¹, Rita Cucchiara^{1,3} ¹University of Modena and Reggio Emilia; ²GoatAl S.r.l.; ³IIT-CNR, Italy

3. DistSynth: estimating per-instance distance from a monocular image

Case study: refining costs related to 2D positions through 3D cues.

DistSynth (our proposal). An ANN giving, for each detected object, its distance from the camera.

It exploits the Feature Pyramid **Network** (FPN) layout to preserve the visual details of distant objects.

It takes the previous L = 5 frames as additional input to exploit temporal patterns and to handle temporary occlusions.



4. Fusing costs: formulation

Formulation. Given the track T, a candidate detection D and the resulting displacements Δ_p , $\Delta_{w,h}$, and Δ_d , we define the fusing cost $\Phi(T,D)$ as the **negative log-likelihood**:

 $\Phi(T,D) = -\log \mathcal{P}_{\theta}(D \in T | T).$

We apply Maximum Likelihood Estimation (MLE) and learn a deep generative model $f([\Delta_p, \Delta_{w,h}, \Delta_d] | T, \theta)$ promoting the likelihood of correct associations.

TrackFlow. The design of $f(\cdot | T, \theta)$ derives from normalizing flow models, which create an invertible mapping between a tractable base distribution and an arbitrary complex one.

<u>IV M</u>

Complex

High flexibility



Exact estimate of the likelihood of a sample





PARIS

6. Experiments: multi-object tracking

Results on MOTSynth [1] & MOTChallenge



(MOTSynth)	Easy		Moderate		Hard	
Metrics	HOTA ↑	IDF1 ↑	HOTA ↑	IDF1 ↑	HOTA ↑	IDF1 ↑
SORT	63.48	79.40	50.31	62.11	37.48	45.13
+ TrackFlow GT	+4.37	+7.41	+5.33	+9.09	+6.54	+10.88
+ TrackFlow	+0.31	+0.97	+0.81	+1.63	+0.74	+1.56
ByteTrack	63.22	80.84	49.91	62.46	37.61	46.15
+ TrackFlow GT	+3.76	+2.82	+5.47	+5.51	+5.08	+4.60
+ TrackFlow	+0.13	+1.80	+0.47	+1.21	+0.88	+1.81
OC-SORT	65.56	81.61	52.42	63.50	38.10	45.48
+ TrackFlow GT	+2.41	+3.76	+4.88	+7.70	+6.18	+9.55
+ TrackFlow	+0.44	+0.84	+0.60	+1.09	+1.17	+1.96

Improved identity accuracy and steady enhancement of IDF1.

MO	Г17	MOT20		
HOTA ↑	IDF1 ↑	HOTA ↑	IDF1 ↑	
64.17	72.98	60.56	74.30	
+1.78	+1.41	+0.15	+0.22	
67.73	79.81	58.94	74.89	
+0.40	+0.23	+0.54	+0.06	
66.22	77.74	55.18	71.22	
+0.35	+1.12	+0.53	+0.76	
	MO [¬] HOTA ↑ 64.17 +1.78 67.73 +0.40 66.22 +0.35	MOT17HOTA ↑IDF1 ↑64.1772.98+1.78+1.4167.7379.81+0.40+0.2366.2277.74+0.35+1.12	MOT17MOT $HOTA \uparrow$ $IDF1 \uparrow$ $HOTA \uparrow$ 64.17 72.98 60.56 $+1.78$ $+1.41$ $+0.15$ 67.73 79.81 58.94 $+0.40$ $+0.23$ $+0.54$ 66.22 77.74 55.18 $+0.35$ $+1.12$ $+0.53$	

7. Experiments: distance estimation

	AL(O)E. Avg. Localization (Occluded objs) Erro								
ALE	3 Zhu et al Ours 1 Ours 1 Ours 1 Thu et al Ours 1 Ours 1 Thu et al Ours 1 Thu et al Thu et al	[20, 30] istance [m]	[30, 100]	BOTH 2 	.3, 0.5] [0.5, 0.7 Occlusion %	5] [0.75, 1.0]			
	Metrics	$\delta_{<1.25}\uparrow$	RMSE ↓	ALP @0.5 <i>m</i>	ALP@1m	ALP@2m			
_	SVR	26.7%	12.5	3.4%	6.8%	13.8%			
	DisNet [2]	27.5%	12.1	3.8%	7.5%	14.6%			
	Zhu et al. [5]	94.7 %	2.15	34.5%	56.2%	78.5%			
	DistSynth	99.1%	1.91	48.0%	68.9%	86.1%			

[1] M. Fabbri et al. Motsynth: How can synthetic data help pedestrian detection and tracking? In *ICCV*, 2021

M. A. Haseeb et al. Disnet: a novel method for distance estimation from monocular camera. 10th Planning, Perception and Navigation for Intelligent Vehicles, IROS, 2018.

[3] J. Rajasegaran et al. Tracking people by predicting 3d appearance, location and pose. In *CVPR*, 2022.

N. Wojke et al. Simple online and realtime tracking with a deep association metric. In ICIP,

J. Zhu et al. Learning object-specific distance from a monocular image. In ICCV, 2019.