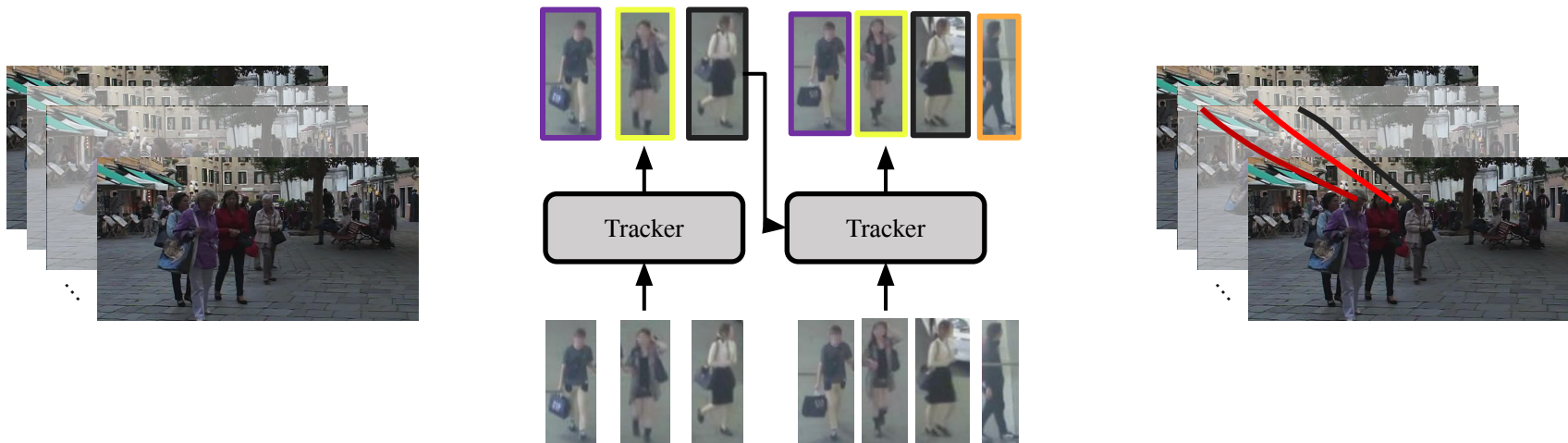


# Simple Cues Lead to a Strong Multi-Object Tracker

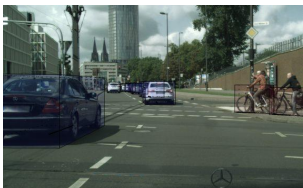
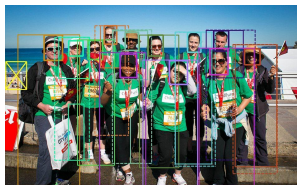
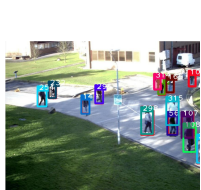
Jenny Seidenschwarz, Guillem Brasó, Victor Castro  
Serrano, Ismail Elezi, Laura Leal-Taixé

# Haunting Completely End-to-End Trackers with GHOST

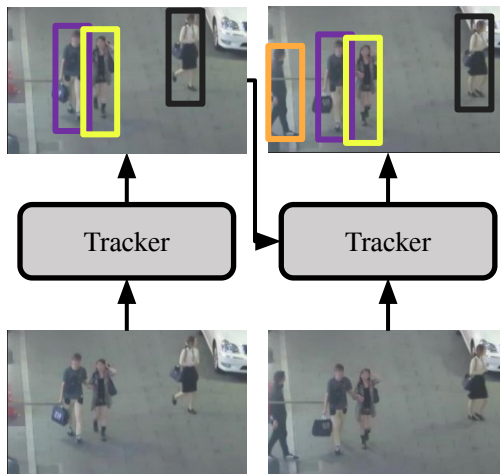
Tracking-by-Detection



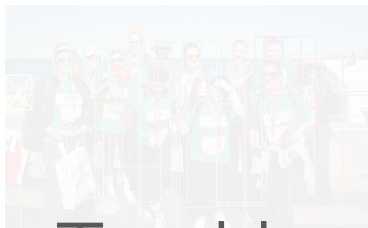
# Haunting Completely End-to-End Trackers with GHOST



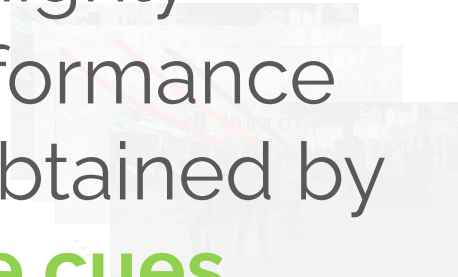
End-to-End Trainable



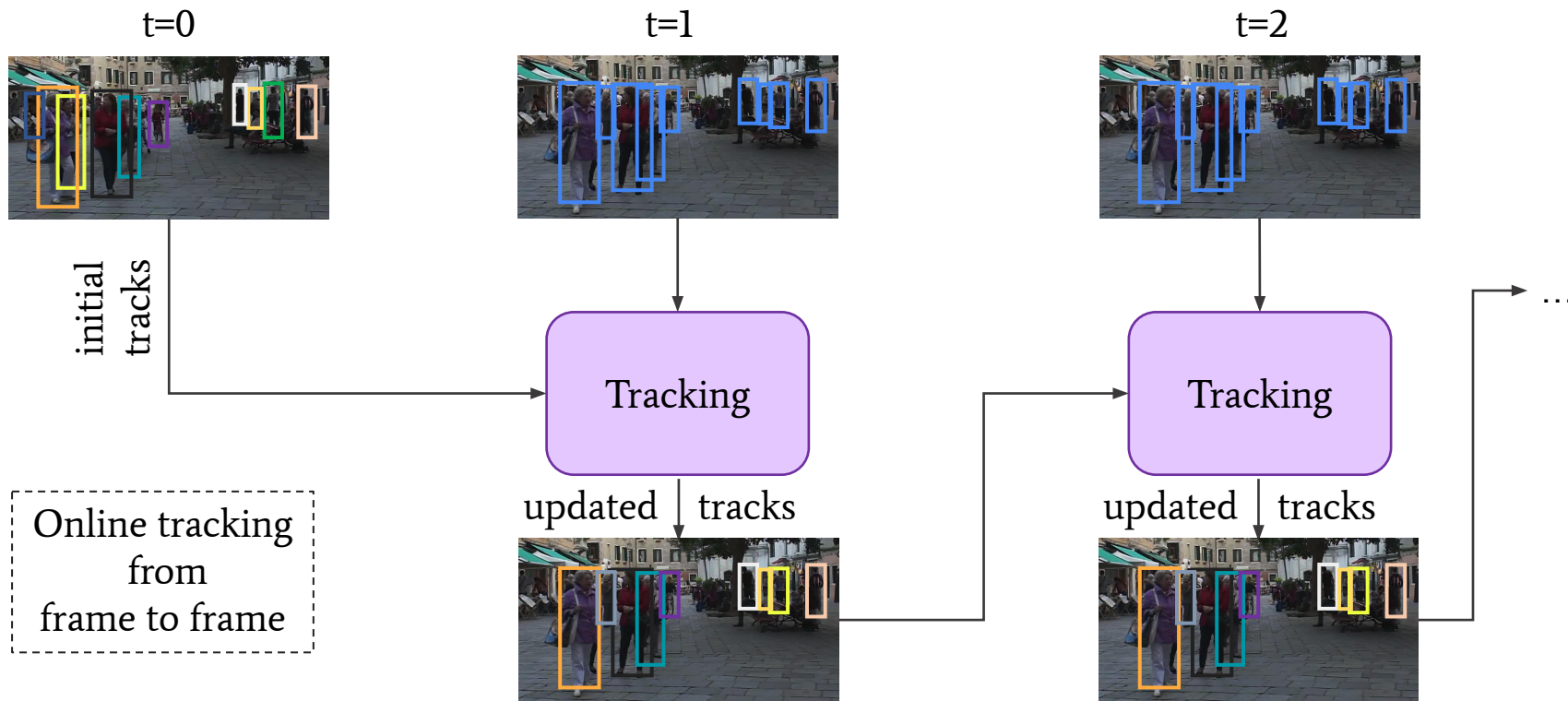
# Haunting Completely End-to-End Trackers with GHOST



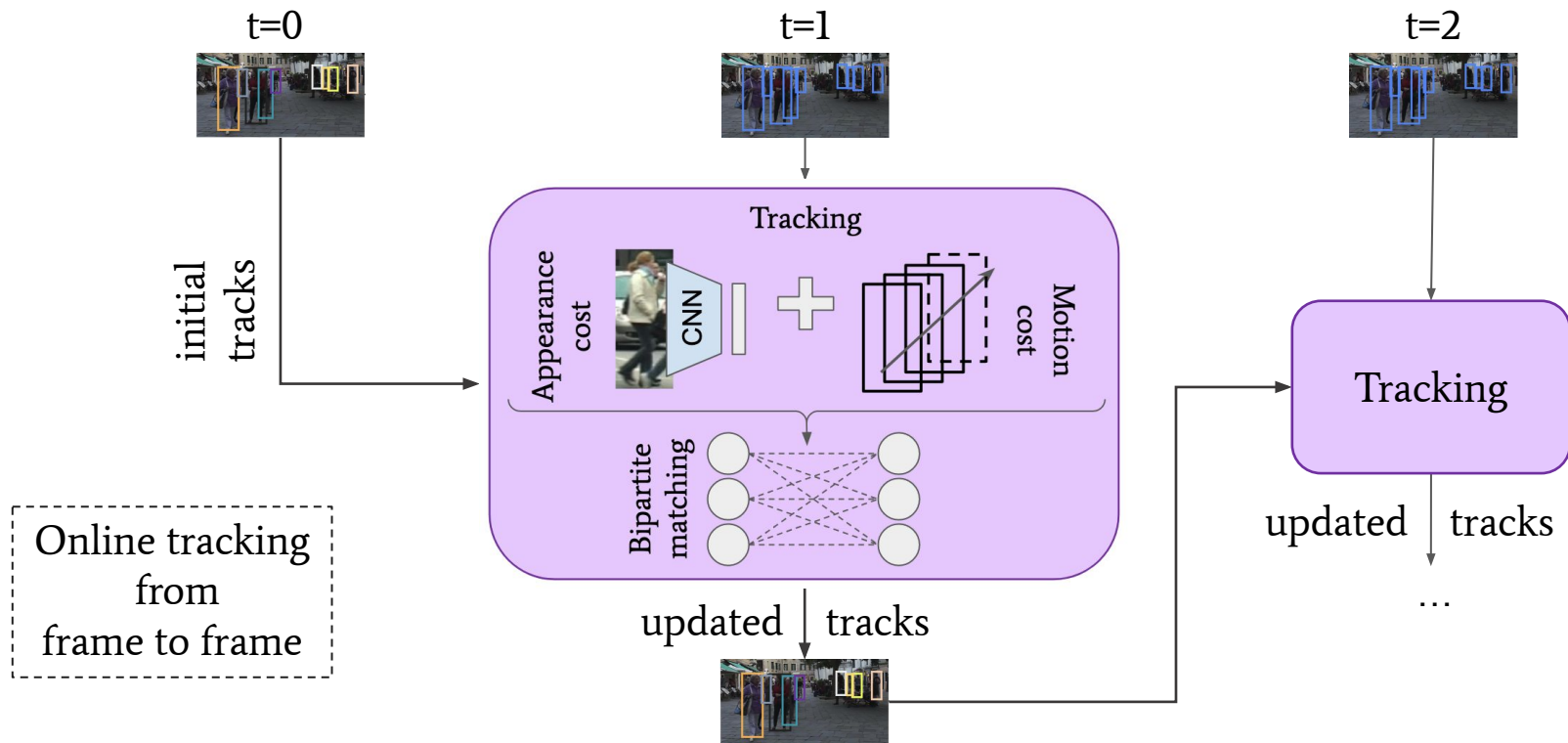
Tracking-by-Detection is highly **general** and shows **SOTA** performance if we follow **key observations** obtained by in-depth **analysis of simple cues**



# Simple Online Tracking



# Simple Online Tracking



# Observations Appearance Model

1. different challenges **active** and **inactive** tracks



person gets  
...  
occludes



track inactive

person  
...  
reappears



# Observations Appearance Model

1. different challenges **active** and **inactive** tracks

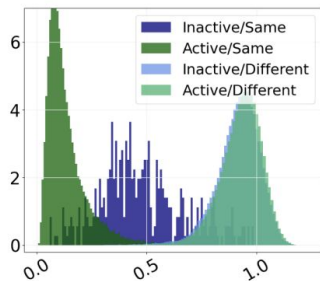


person gets  
...  
occludes



track inactive

person  
...  
reappears



(a) Appearance



# Observations Appearance Model

1. different challenges **active** and **inactive** tracks

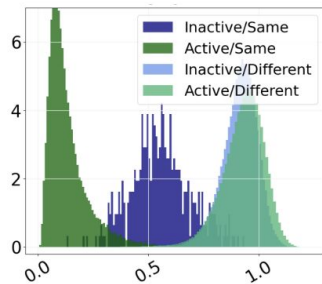


person gets  
...  
occludes



track inactive

person  
...  
reappears



(b) Proxy Appearance

$$d_{i,k} = \frac{1}{N_k} \sum_{n=1}^{N_k} d(f_i, f_k^n)$$

# Observations Appearance Model

1. different challenges **active** and **inactive** tracks

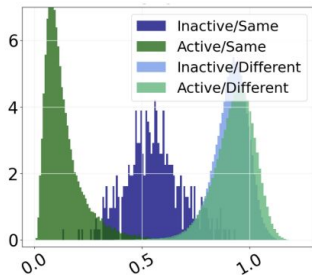


person gets  
...  
occludes



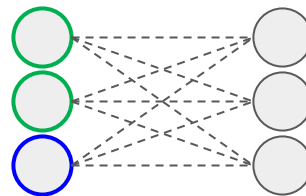
track inactive

person  
...  
reappears

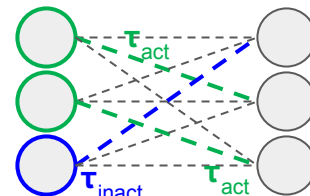


(b) Proxy Appearance

$$d_{i,k} = \frac{1}{N_k} \sum_{n=1}^{N_k} d(f_i, f_k^n)$$



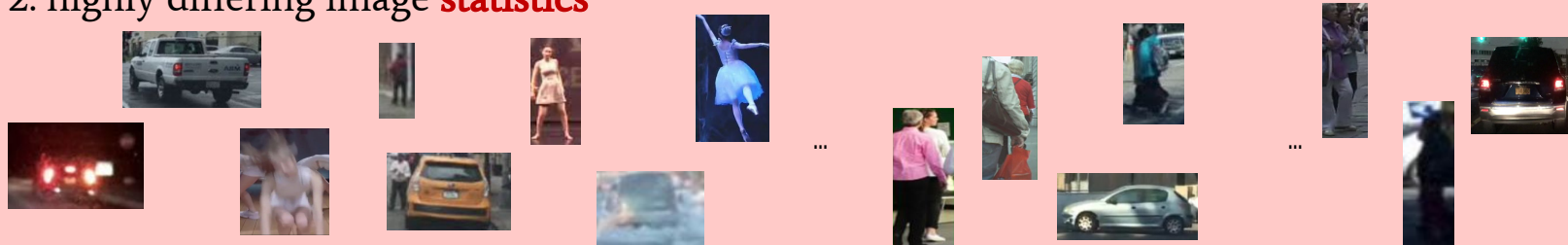
matching



thresholding

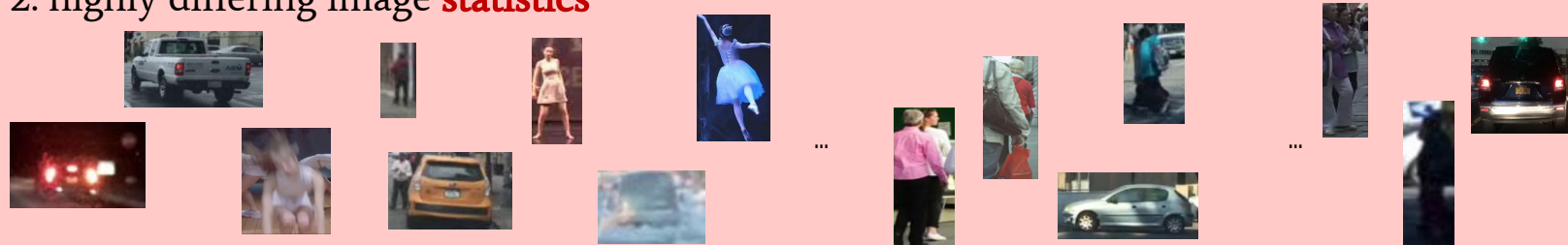
# Observations Appearance Model

2. highly differing image **statistics**



# Observations Appearance Model

2. highly differing image **statistics**



Frame-wise statistics



D detections

$$\mu_b = \frac{1}{D} \sum_{i=0}^D x_i$$
$$\sigma_b = \frac{1}{D} \sum_{i=0}^D (x_i - \mu_b)^2$$

Features adapted to specific frame

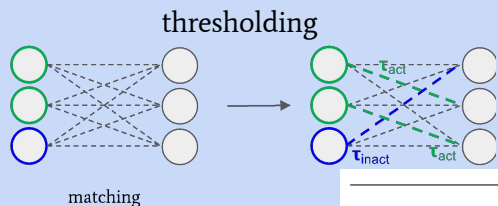
$$\hat{x}_i = \gamma \frac{x_i - \mu_b}{\sqrt{\sigma_b + \epsilon}} + \beta$$

# Observations Appearance Model

1. different challenges **active** and **inactive** tracks

proxy distance

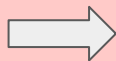
$$d_{i,k} = \frac{1}{N_k} \sum_{n=0}^{N_k} d(f_i, f_k^n)$$



Spiced up appearance model  
better suited for MOT

diff $\tau$	IP	DA	MOT 17			BDD		
			HOTA $\uparrow$	IDF1 $\uparrow$	MOTA $\uparrow$	HOTA $\uparrow$	IDF1 $\uparrow$	MOTA $\uparrow$
			61.6	72.2	69.5	41.3	47.7	42.0
✓			61.8	72.7	69.6	41.9	48.6	41.9
✓	✓		62.4	73.6	69.6	42.9	50.4	43.7
✓	✓	✓	63.3	75.3	69.6	43.7	51.5	43.9

2. highly differing image **statistics**

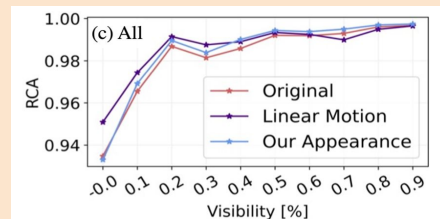
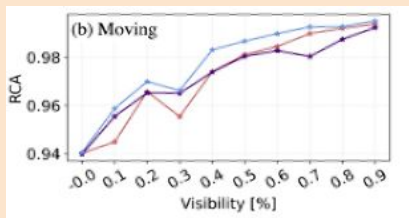
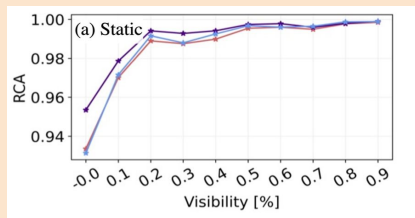


Features adapted to  
specific frame

$$\hat{x}_i = \gamma \frac{x_i - \mu_b}{\sqrt{\sigma_b + \epsilon}} + \beta$$

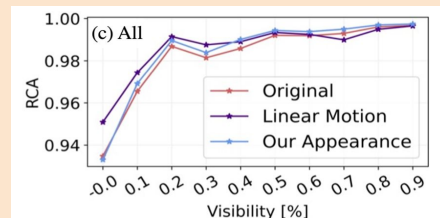
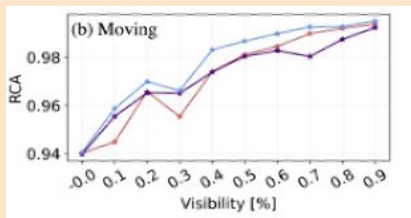
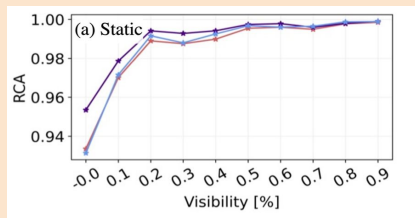
# Observations Motion Model

## 3. Linear motion model complements appearance well

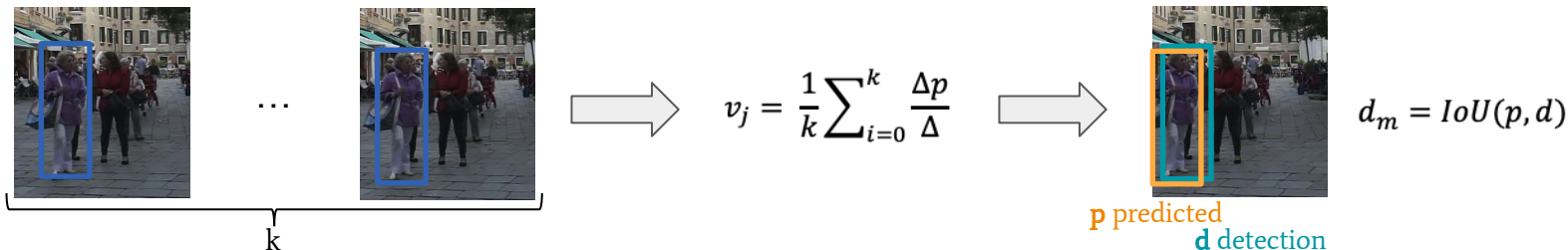


# Observations Motion Model

## 3. Linear motion model complements appearance well

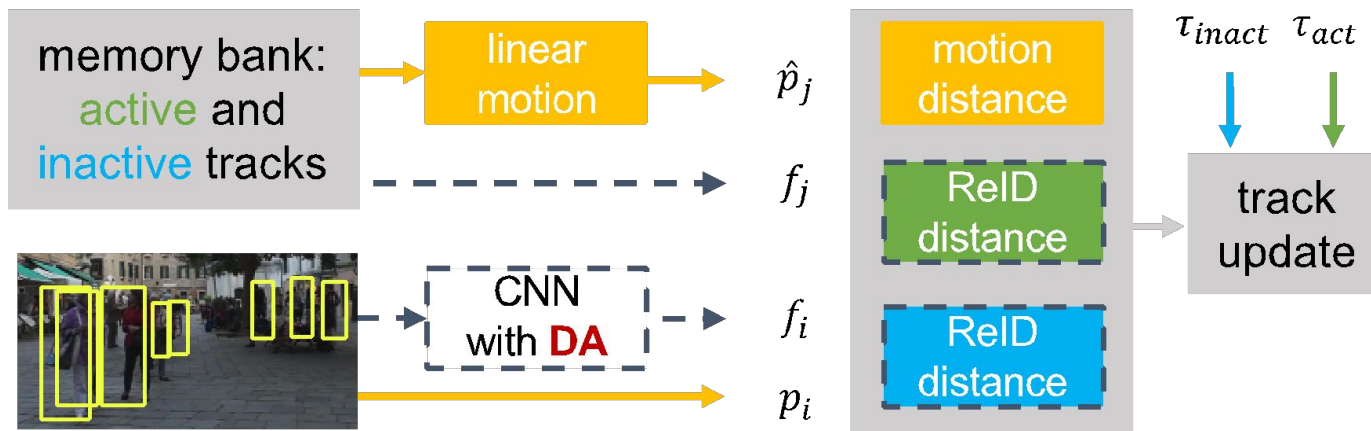


Motion still struggles with moving camera → adaptive frames  $k$  for velocity computation



# Good Old Hungarian Simple Tracker or GHOST\*

\*the order of the letters of the acronym does not change the product





# SOTA on highly differing datasets w/o training on tracking data!



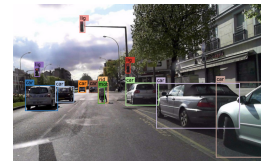
MOT17

	HOTA ↑	IDF1 ↑	MOTA ↑	IDSW ↓
<i>Private MOT17</i>				
CenterTrack [16]	52.2	64.7	67.8	3039
TraDeS [12]	52.7	63.9	69.1	3555
QDTrack [7]	53.9	66.3	68.7	3378
FairMOT [15]	<b>59.3</b>	<b>72.3</b>	73.7	3303
MeMOT [3]	56.9	<b>72.5</b>	69.0	2724
GTR [17]	<b>59.1</b>	71.5	<b>75.3</b>	<b>2859</b>
MOTR [13]	57.8	68.6	73.4	2439
ByteTrack* [14]	<b>62.8</b>	<b>77.1</b>	<b>78.9</b>	<b>2363</b>
ByteTrack* [14]	63.1	77.3	80.3	2196
GHOST	<b>62.8</b>	<b>77.1</b>	<b>78.7</b>	<b>2325</b>



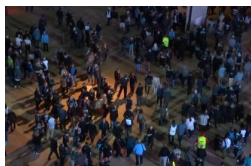
BDD100k

	mHOTA ↑	mIDF1 ↑	mMOTA ↑	HOTA ↑	IDF1 ↑	MOTA ↑
<i>validation</i>						
Yu et. al. [69]	-	44.5	25.9	-	66.8	56.9
ByteTrack [73]	<b>45.4</b>	<b>54.6</b>	<b>45.2</b>	<b>61.6</b>	<b>70.2</b>	<b>68.7</b>
QDTrack [41]	<b>41.7</b>	51.5	36.3	<b>60.9</b>	<b>71.4</b>	<b>63.7</b>
MOTR [71]	-	43.5	32.0	-	-	-
TETer [29]	-	<b>53.3</b>	<b>39.1</b>	-	-	-
GHOST	<b>45.7</b>	<b>55.6</b>	<b>44.9</b>	<b>61.7</b>	<b>70.9</b>	<b>68.1</b>
<i>test</i>						
Yu et. al. [69]	-	44.7	26.3	-	68.2	58.3
ByteTrack [73]	-	<b>55.8</b>	<b>40.1</b>	-	<b>71.3</b>	<b>69.6</b>
QDTrack [41]	-	<b>52.3</b>	<b>35.5</b>	-	<b>72.3</b>	<b>64.3</b>
GHOST	<b>46.8</b>	<b>57.0</b>	<b>39.5</b>	<b>62.2</b>	<b>72.0</b>	<b>68.9</b>



MOT20

	HOTA ↑	IDF1 ↑	MOTA ↑	IDSW ↓
<i>Private MOT20</i>				
GSDT [10]	53.6	67.5	67.1	3230
FairMOT [15]	<b>54.6</b>	67.3	61.8	5243
MeMOT [3]	54.1	66.1	<b>63.7</b>	<b>1938</b>
MTrack [13]	-	<b>69.2</b>	63.5	6031
ByteTrack* [14]	<b>60.4</b>	<b>74.5</b>	<b>74.2</b>	<b>925</b>
ByteTrack* [14]	61.3	75.2	77.8	1223
GHOST	<b>61.2</b>	<b>75.2</b>	<b>73.7</b>	<b>1264</b>



DanceTrack

	HOTA ↑	IDF1 ↑	MOTA ↑	DetA ↑	AssA ↑
CenterTrack [79]	41.8	35.7	86.8	<b>78.1</b>	22.6
FairMOT [74]	39.7	40.8	82.2	66.7	23.8
QDTrack [41]	<b>54.2</b>	50.4	<b>87.7</b>	<b>80.1</b>	<b>36.8</b>
TraDeS [64]	43.3	41.2	86.2	74.5	25.4
MOTR [71]	<b>54.2</b>	<b>51.5</b>	79.7	73.5	<b>40.2</b>
GTR [80]	<b>48.0</b>	50.3	84.7	72.5	31.9
ByteTrack [73]	47.7	<b>53.9</b>	<b>89.6</b>	71.0	32.1
GHOST	<b>56.7</b>	<b>57.7</b>	<b>91.3</b>	<b>81.1</b>	<b>39.8</b>



# SOTA on highly differing datasets w/o training on tracking data! 🤩

