

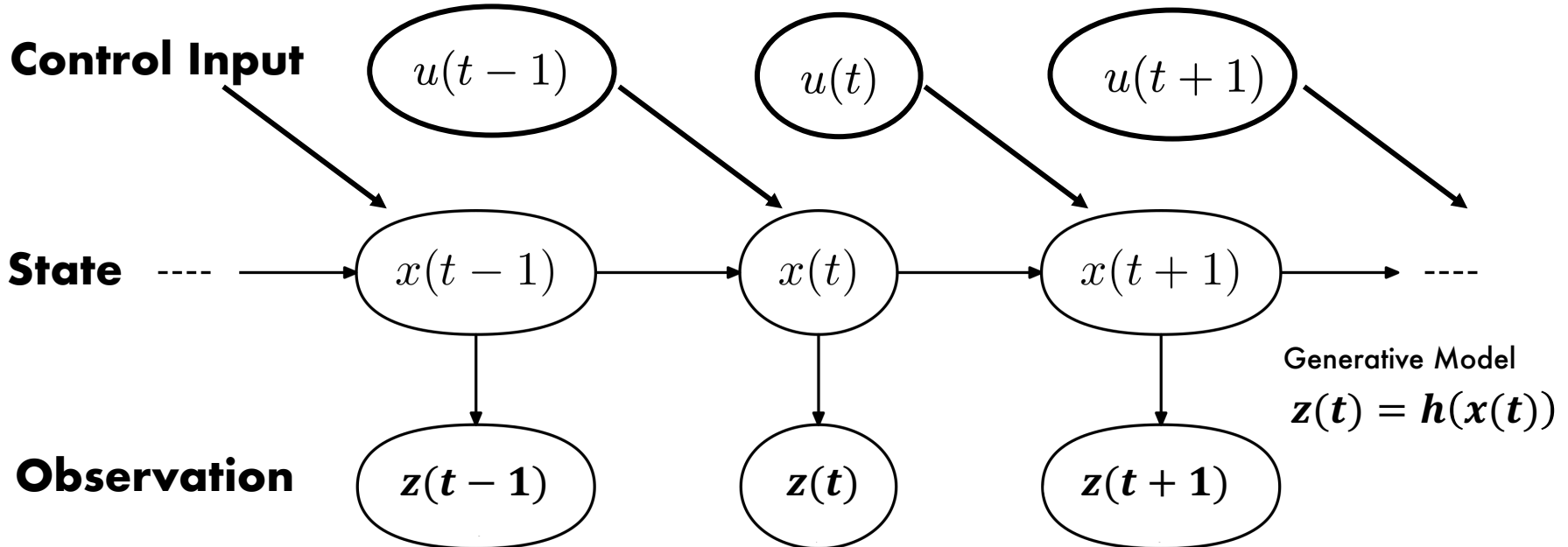
CS231A

**Computer Vision:
From 3D Reconstruction
to Recognition**

Optimal Estimation Cont'

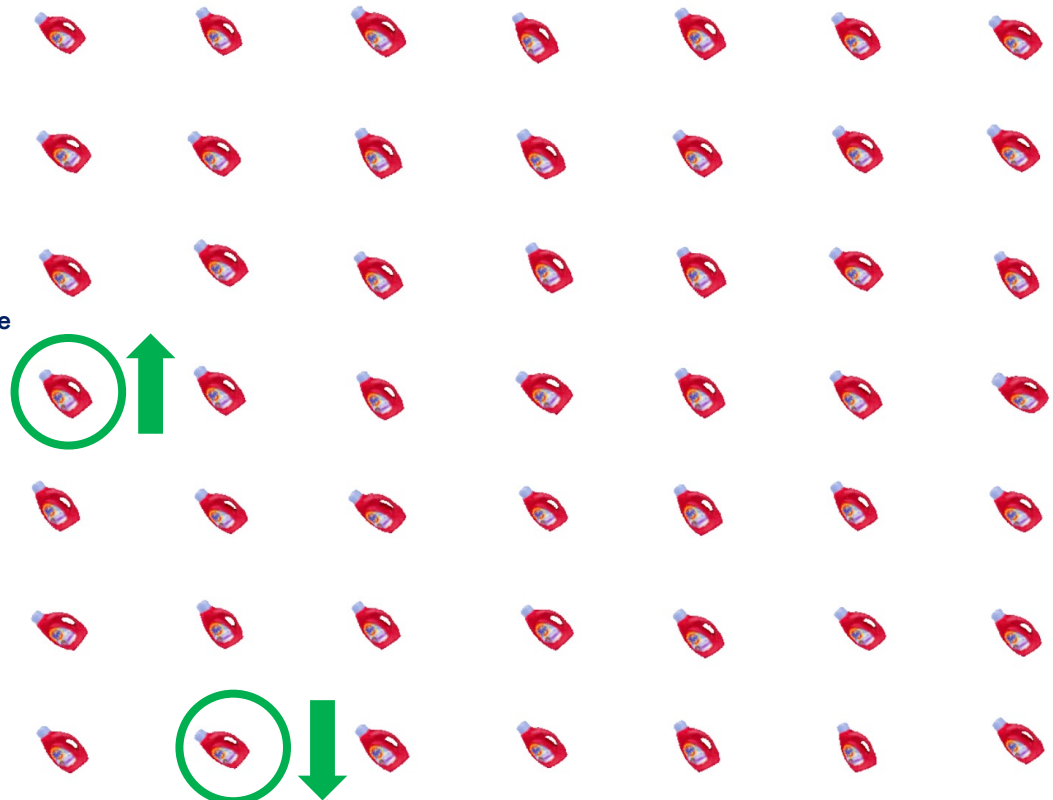
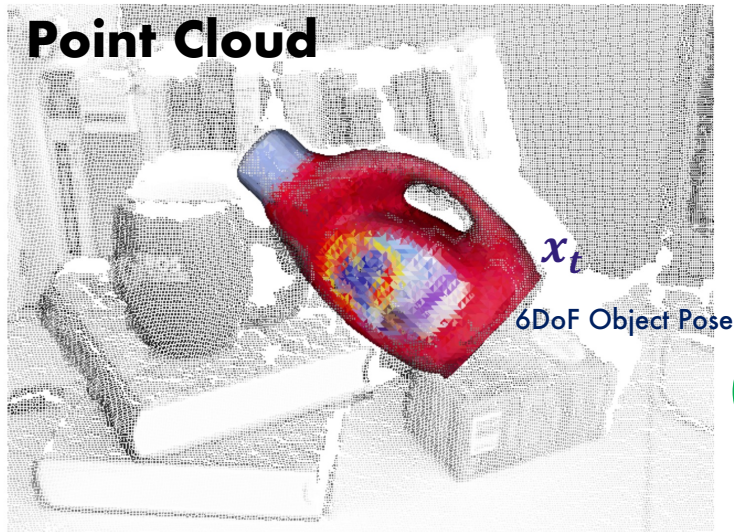


Graphical Model of System to Estimate



```
1: Algorithm Bayes filter( $bel(x_{t-1}), u_t, z_t$ ):  
2:   for all  $x_t$  do  
3:      $\overline{bel}(x_t) = \int p(x_t | u_t, x_{t-1}) bel(x_{t-1}) dx$   
4:      $bel(x_t) = \eta p(z_t | x_t) \overline{bel}(x_t)$   
5:   endfor  
6:   return  $bel(x_t)$ 
```

Example Observation model for 3D object



Algorithm Particle filter($\mathcal{X}_{t-1}, u_t, z_t$):

```

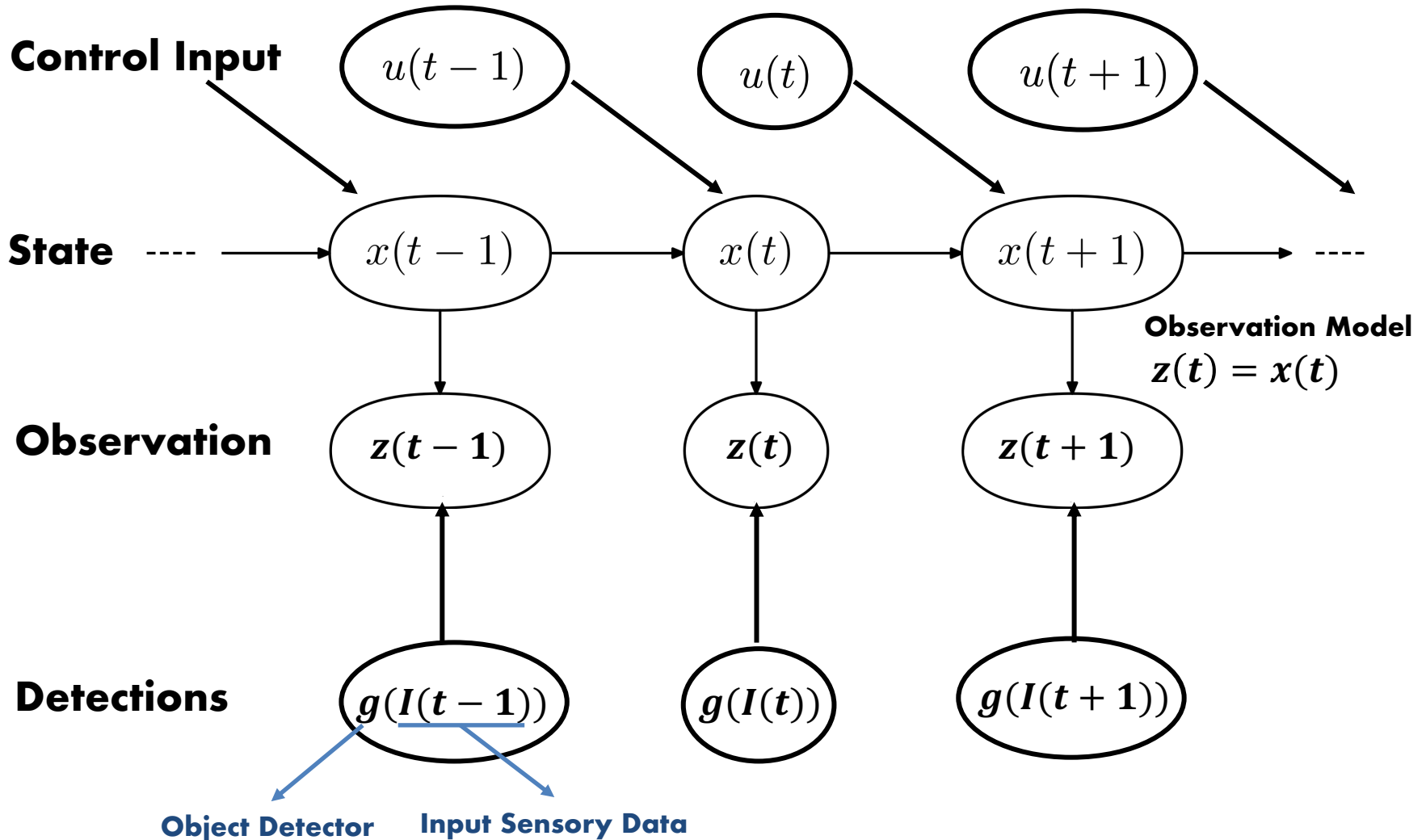
 $\mathcal{X}_t = \mathcal{X}_{t-1} = \emptyset$ 
for  $m = 1$  to  $M$  do
  sample  $x_t^{[m]} \sim p(x_t | u_t, x_{t-1}^{[m]})$ 
   $w_t^{[m]} = p(z_t | x_t^{[m]})$ 
   $\bar{x}_t = \bar{x}_t + \langle x_t^{[m]}, w_t^{[m]} \rangle$ 
endfor
for  $m = 1$  to  $M$  do
  draw  $i$  with probability  $\propto w_t^{[i]}$ 
  add  $x_t^{[i]}$  to  $\mathcal{X}_t$ 
endfor
return  $\mathcal{X}_t$ 

```

Importance Sampling

Rendered Particles

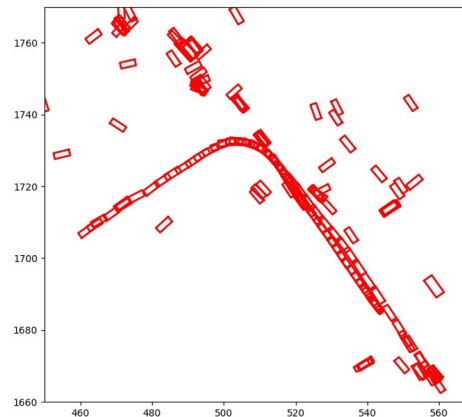
Tracking by Detection



Problem Statement: Input

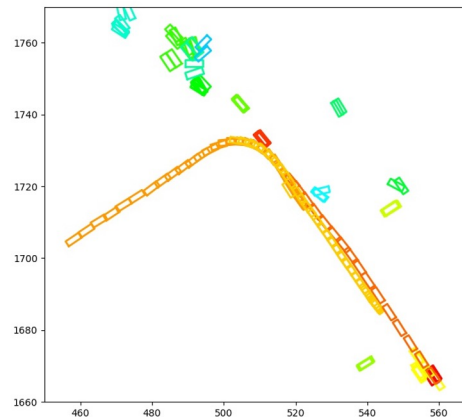
Probabilistic 3d multi-object tracking for autonomous driving. H Chiu, A Prioletti, J Li, J Bohg
arXiv preprint arXiv:2001.05673

- **Object detections** at each frame in a sequence
- Each **detection bounding box** is represented by:
 - center position (x, y, z) , rotation angle along the z -axis (a) , and the scale (l, w, h)
 - category label (car, pedestrian, ...), confidence score (c)



Problem Statement: Output

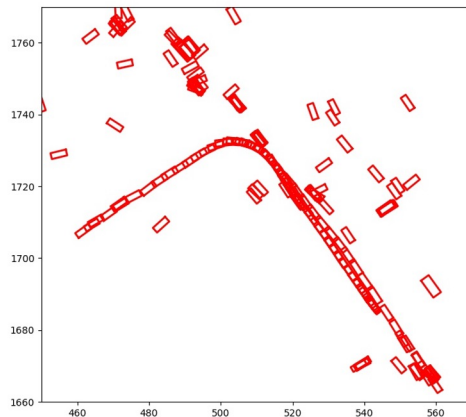
- Tracking object bounding boxes at each frame in a sequence
- Each tracking bounding box is represented by:
 - center position (x, y, z) , rotation angle along the z -axis (a) , and the scale (l, w, h)
 - category label (car, pedestrian, ...), confidence score (c)
 - **tracking id**: one unique tracking id for each object instance across frames



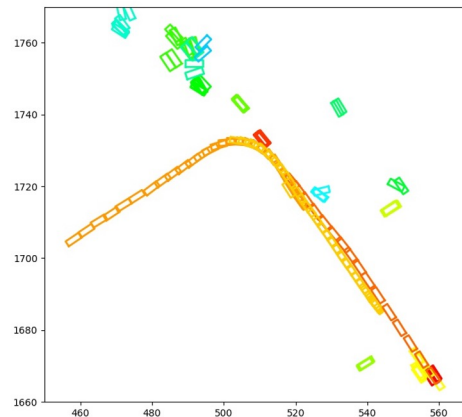
Why Tracking?

- Filter out the out-liners from the detection results
- Continue estimating object states even if occluded
- Forecast the future based on past trajectories and motion patterns
- Make autonomous driving decisions

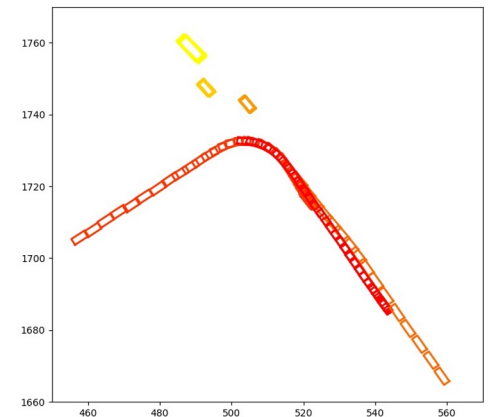
Detection

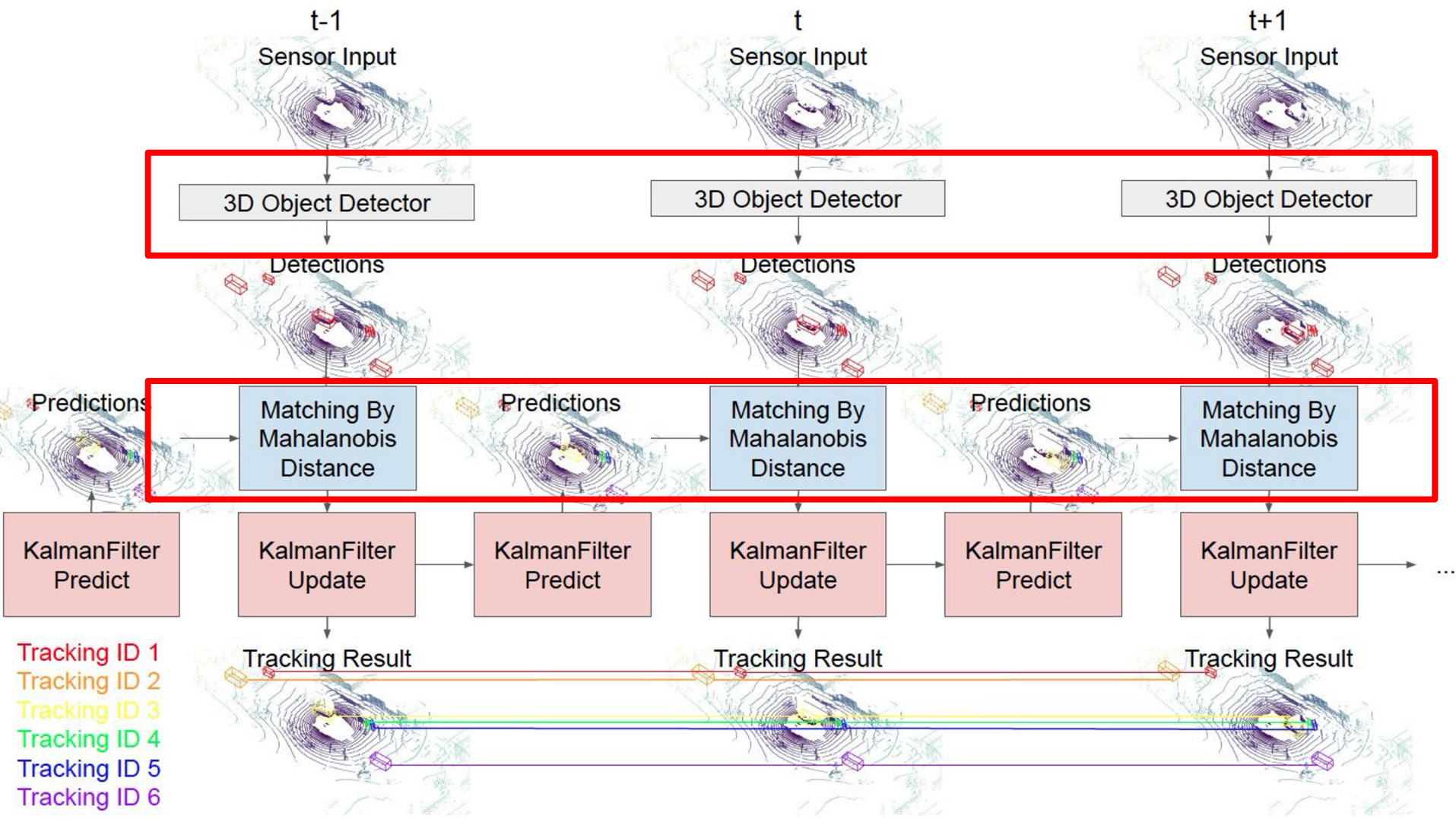


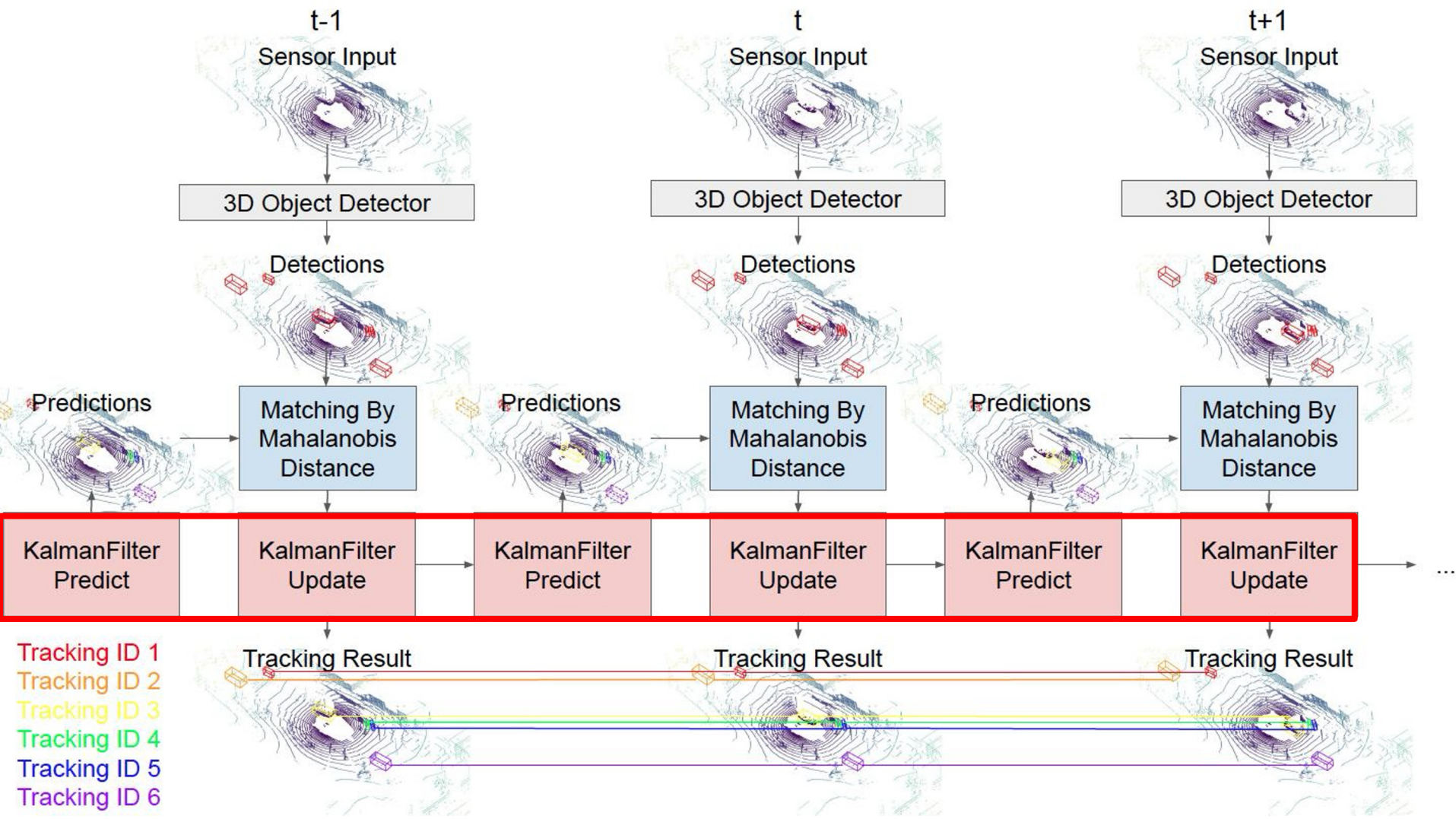
Tracking



Ground-truth







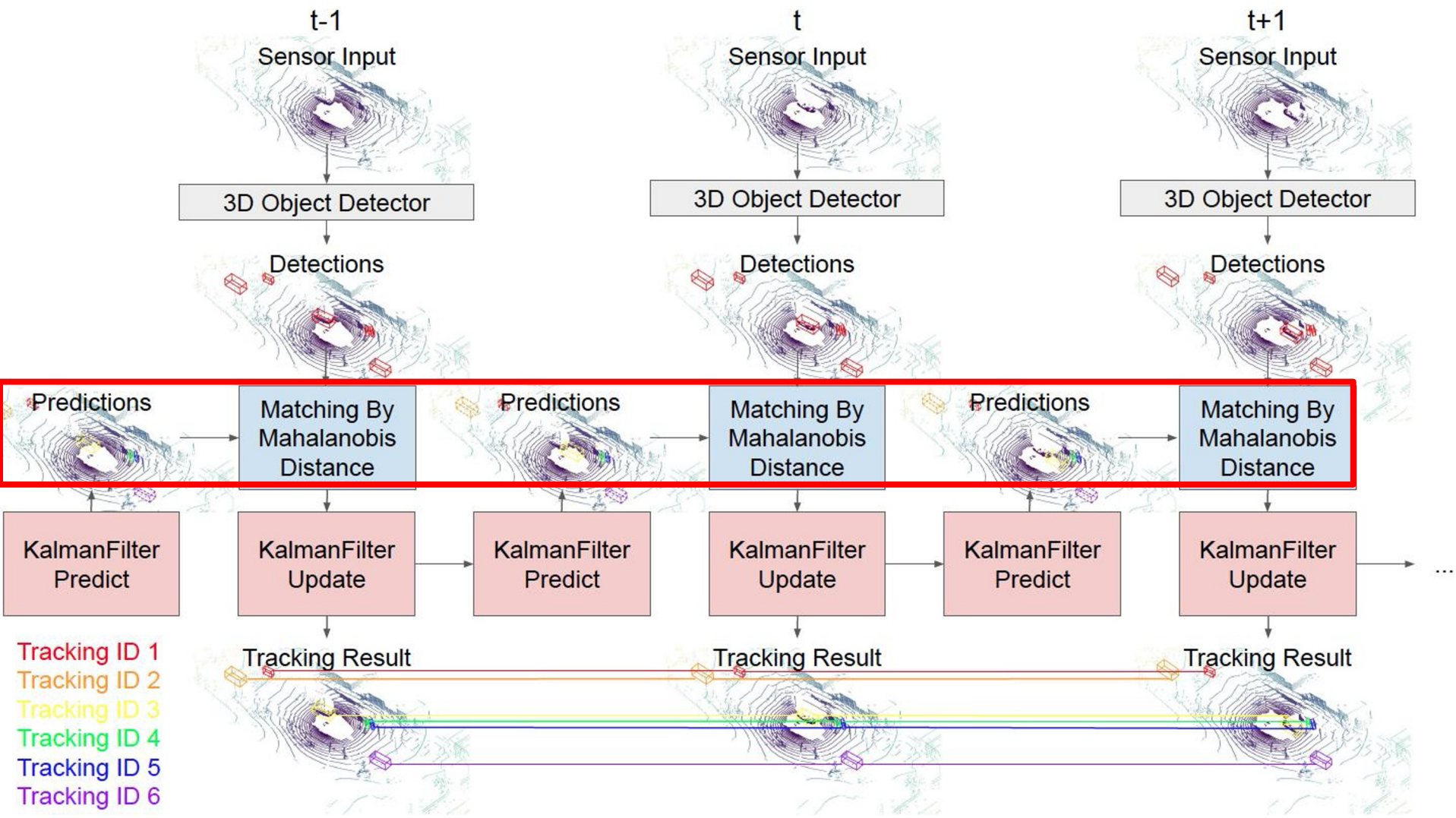
Kalman Filter for Tracking

Define the object **state** using a vector of random variables including the position, the rotation, the scale, linear velocity, and the angular velocity.

$$\mathbf{s}_t = (x, y, z, a, l, w, h, d_x, d_y, d_z, d_a)^T$$

Define the **Process Model** for prediction based on the constant velocity motion:

$$\begin{array}{l|l|l} \hat{x}_{t+1} = x_t + d_{x_t} + q_{x_t}, & \hat{d}_{x_{t+1}} = d_{x_t} + q_{d_{x_t}} & \hat{l}_{t+1} = l_t \\ \hat{y}_{t+1} = y_t + d_{y_t} + q_{y_t}, & \hat{d}_{y_{t+1}} = d_{y_t} + q_{d_{y_t}} & \hat{w}_{t+1} = w_t \\ \hat{z}_{t+1} = z_t + d_{z_t} + q_{z_t}, & \hat{d}_{z_{t+1}} = d_{z_t} + q_{d_{z_t}} & \hat{h}_{t+1} = h_t \\ \hat{a}_{t+1} = a_t + d_{a_t} + q_{a_t}, & \hat{d}_{a_{t+1}} = d_{a_t} + q_{d_{a_t}} & \end{array}$$

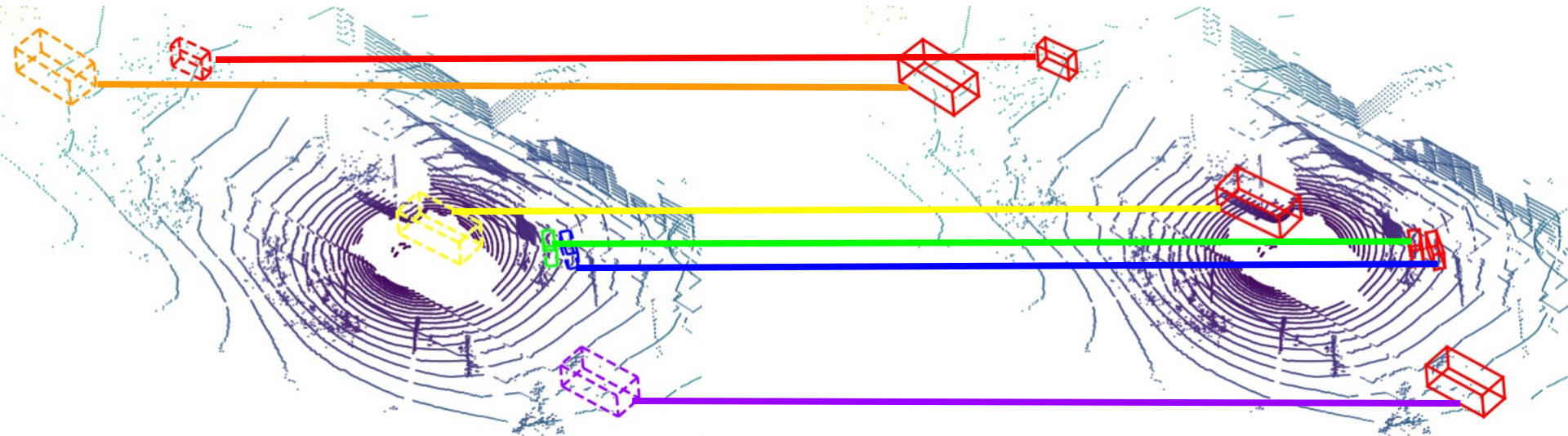


Data Association

$$\text{Mahalanobis Distance } m = \sqrt{(z_t - C\bar{\mu}_t)^T S_t^{-1} (z_t - C\bar{\mu}_t)}$$

S = Innovation Covariance

$z_t - C\bar{\mu}_t = \text{innovation}$



**Kalman Filter
Predictions**

Object Detections

Kalman Filter

- 1: **Algorithm Kalman filter**($\mu_{t-1}, \Sigma_{t-1}, u_t, z_t$):
- 2: $\bar{\mu}_t = A_t \mu_{t-1} + B_t u_t$
- 3: $\bar{\Sigma}_t = A_t \Sigma_{t-1} A_t^T + R_t$
- 4: $K_t = \bar{\Sigma}_t C_t^T (C_t \bar{\Sigma}_t C_t^T + Q_t)^{-1} = S_t^{-1}$
- 5: $\mu_t = \bar{\mu}_t + K_t (z_t - C_t \bar{\mu}_t)$
- 6: $\Sigma_t = (I - K_t C_t) \bar{\Sigma}_t$
- 7: return μ_t, Σ_t

Data Association

Mahalanobis Distance $m = \sqrt{(z_t - C\bar{\mu}_t)^T S_t^{-1} (z_t - C\bar{\mu}_t)}$

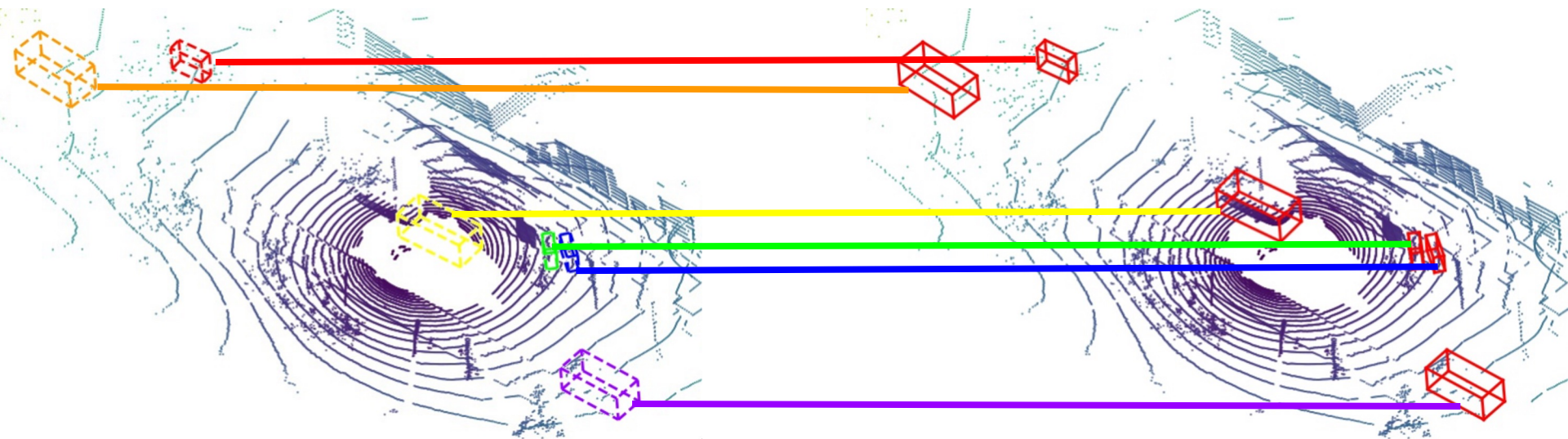
If $m > 3 * \sigma$ then reject as outlier. 99.7% of values lie within 3* standard deviation.

Measuring the distance between the observation and the distribution of the predicted state.

Providing distance measurement **when there is no overlap** between the prediction and detection.

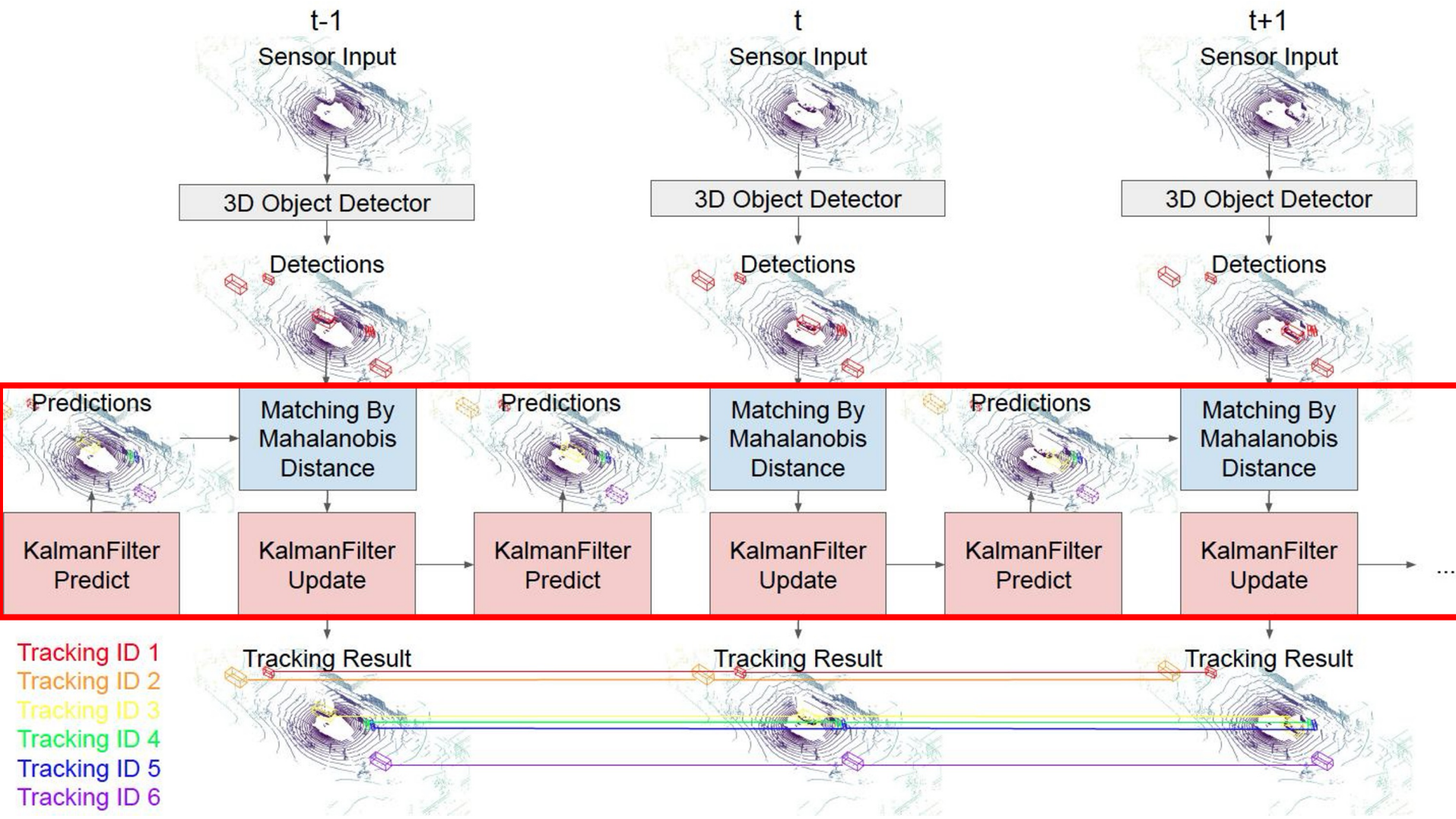
Taking the **uncertainty** information from the prediction into account.

Data Association - Greedy



**Kalman Filter
Predictions**

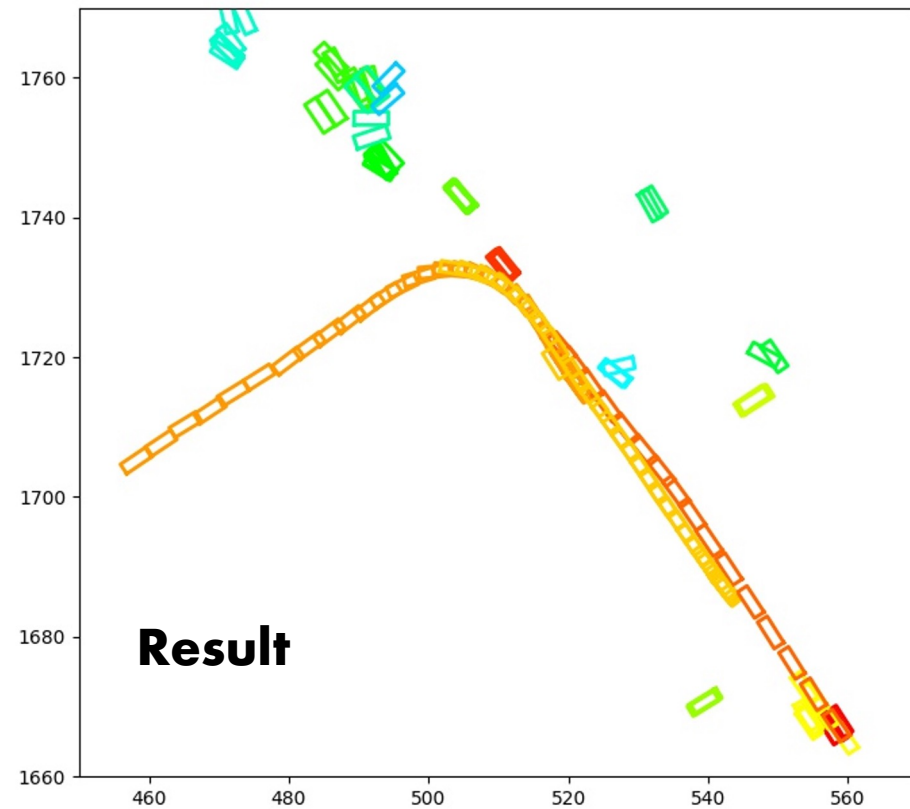
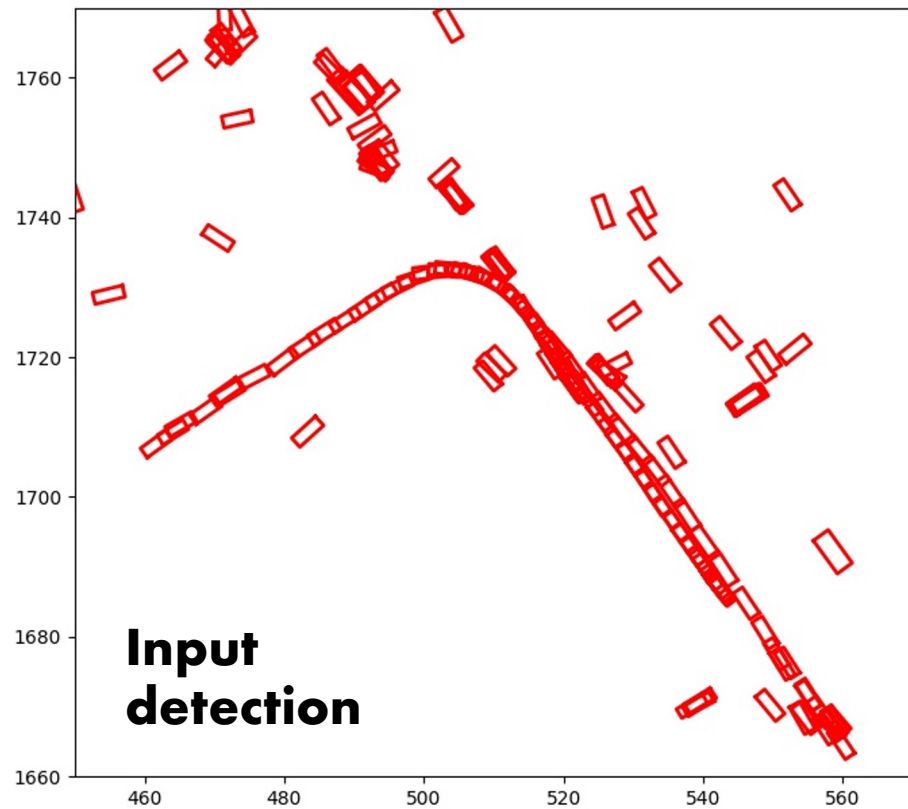
Detections



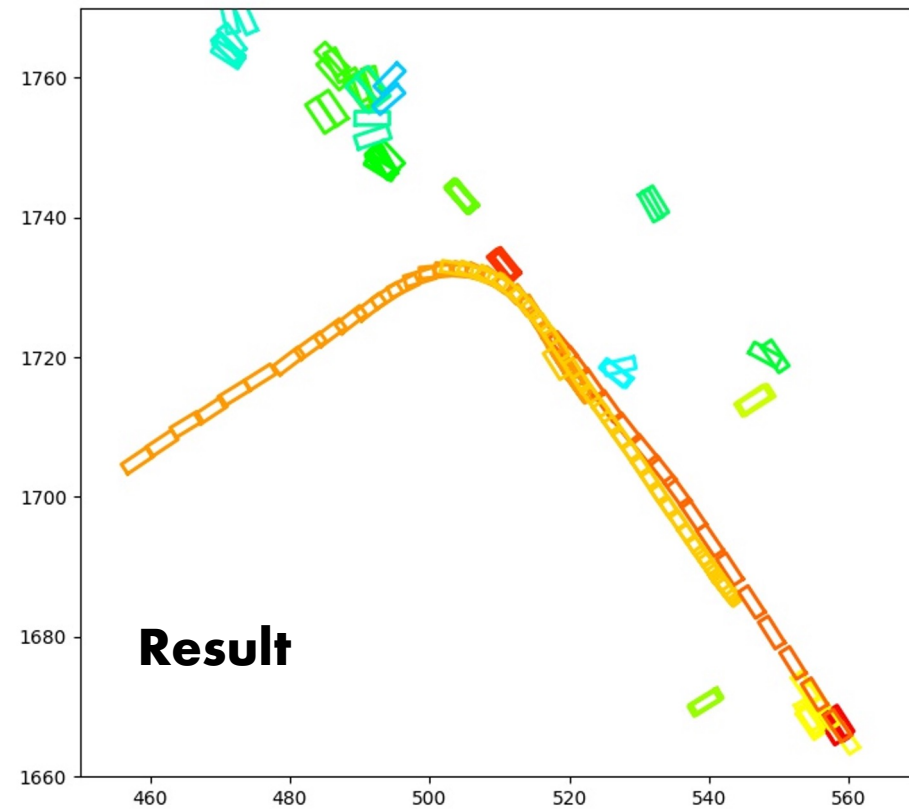
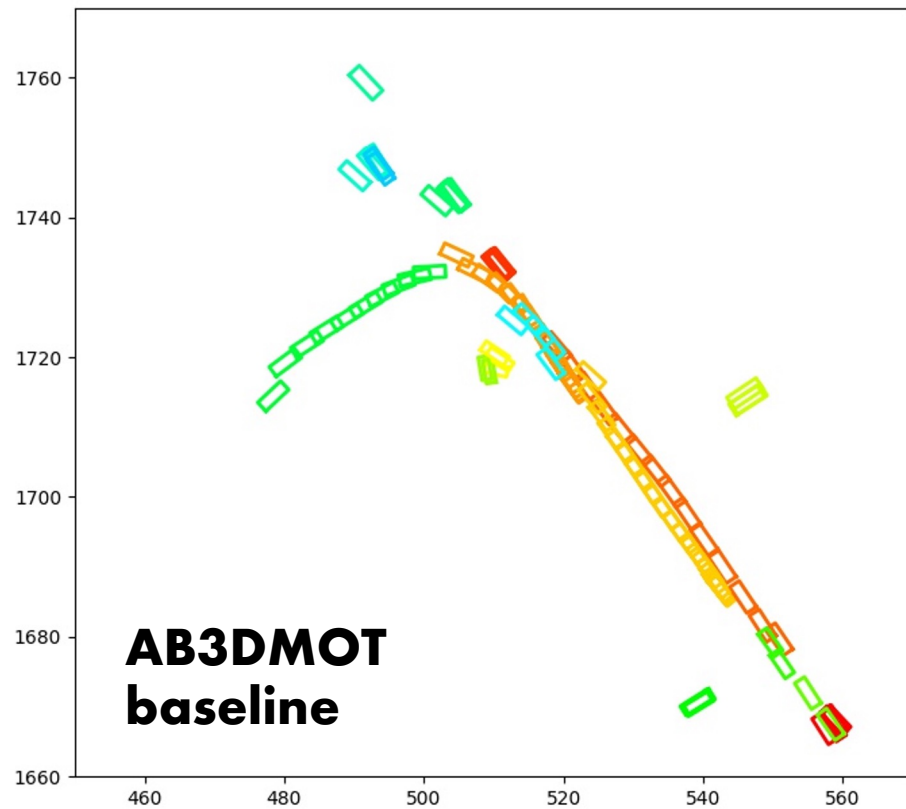
Kalman Filter

1: **Algorithm Kalman filter**($\mu_{t-1}, \Sigma_{t-1}, u_t, z_t$):
2: $\bar{\mu}_t = A_t \mu_{t-1} + B_t u_t$
3: $\bar{\Sigma}_t = A_t \Sigma_{t-1} A_t^T + R_t$
4: $K_t = \bar{\Sigma}_t C_t^T (C_t \bar{\Sigma}_t C_t^T + Q_t)^{-1}$
5: $\mu_t = \bar{\mu}_t + K_t (z_t - C_t \bar{\mu}_t)$
6: $\Sigma_t = (I - K_t C_t) \bar{\Sigma}_t$
7: return μ_t, Σ_t

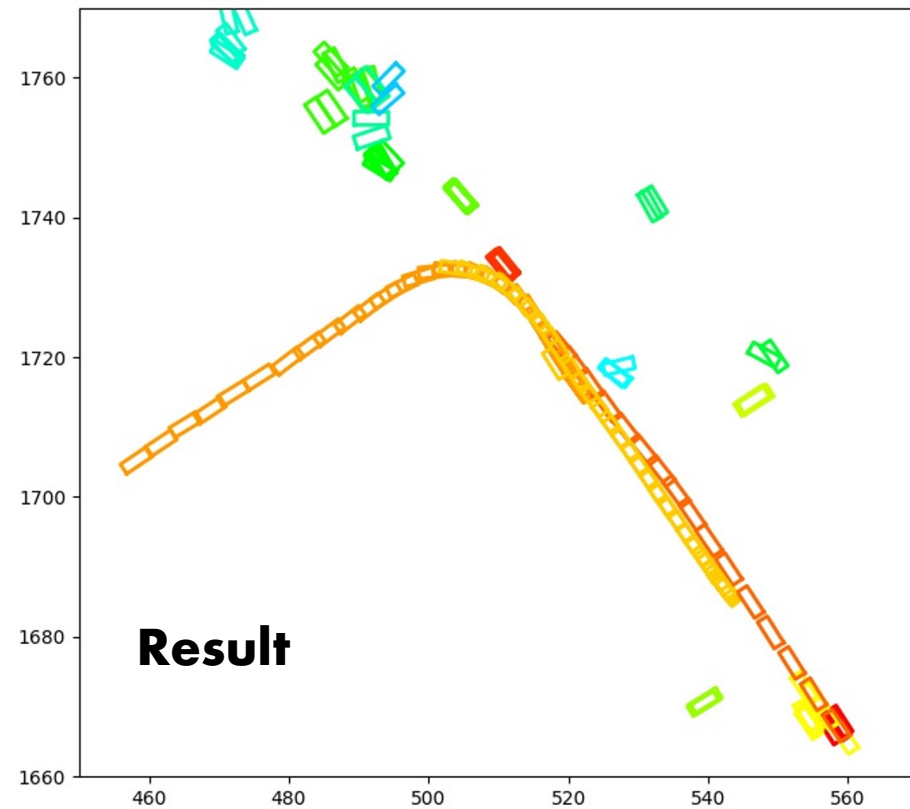
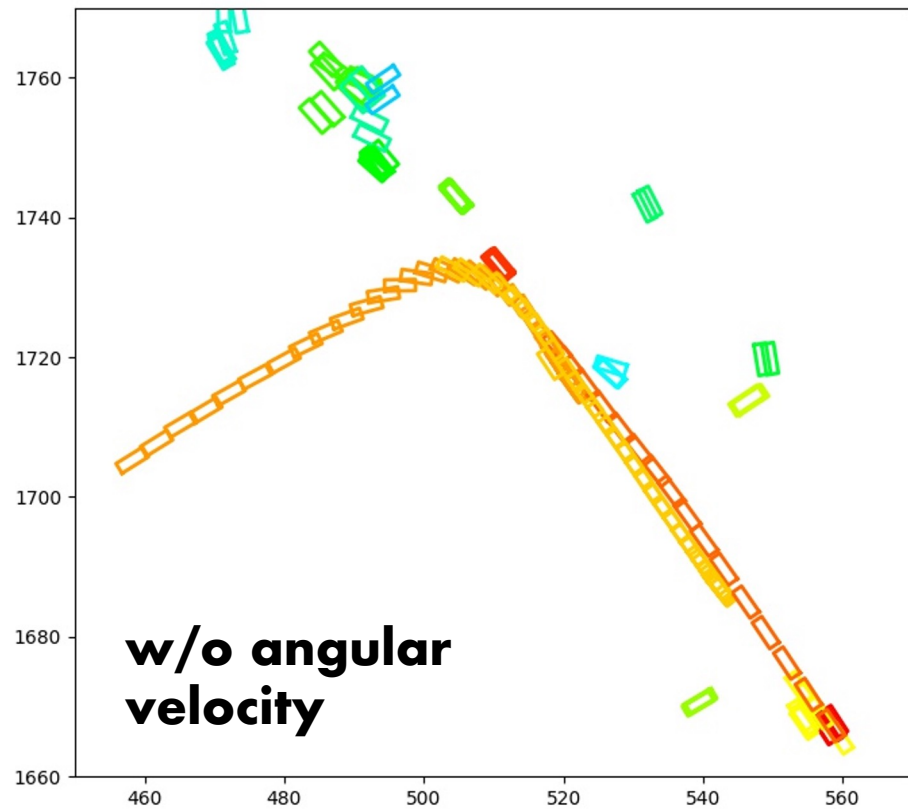
Qualitative Results



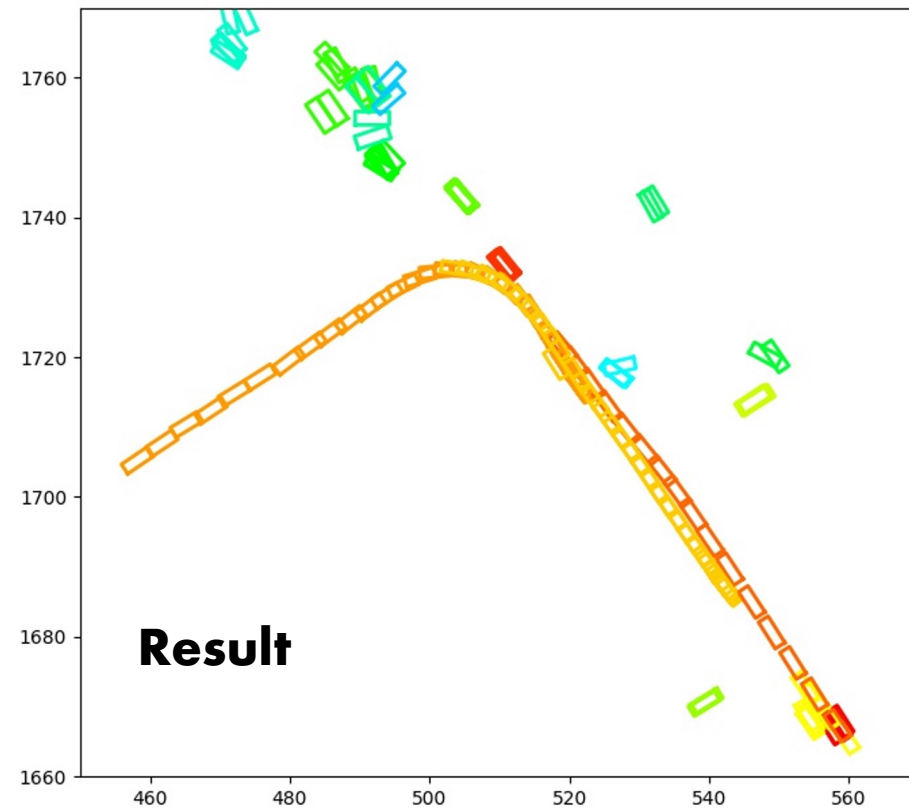
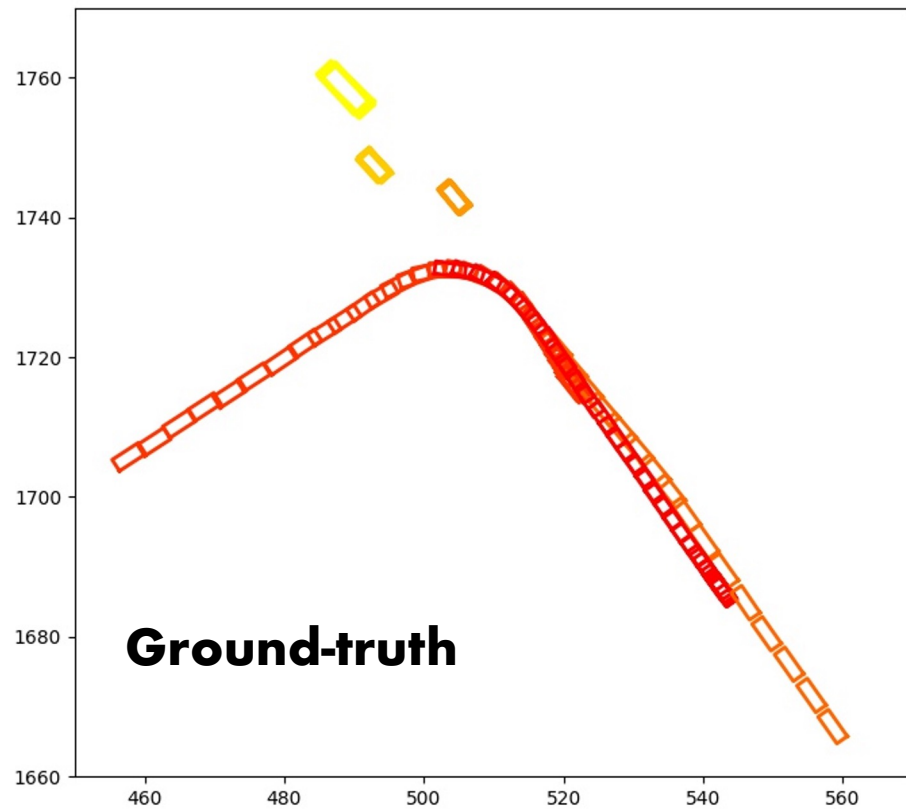
Qualitative Results



Qualitative Results



Qualitative Results



Priors and Hyperparameters

A lot of hardcoded knowledge!

- State Representation
- Models
 - Forward Model
 - State to next state
 - Action to next state
 - Measurement Model
- Probabilistic Properties
 - Process Noise
 - Measurement Noise



Differentiable filters

Can we learn models and hyperparameters from data?

Approach: Embed algorithmic structure of Bayesian Filtering into a recurrent neural network.

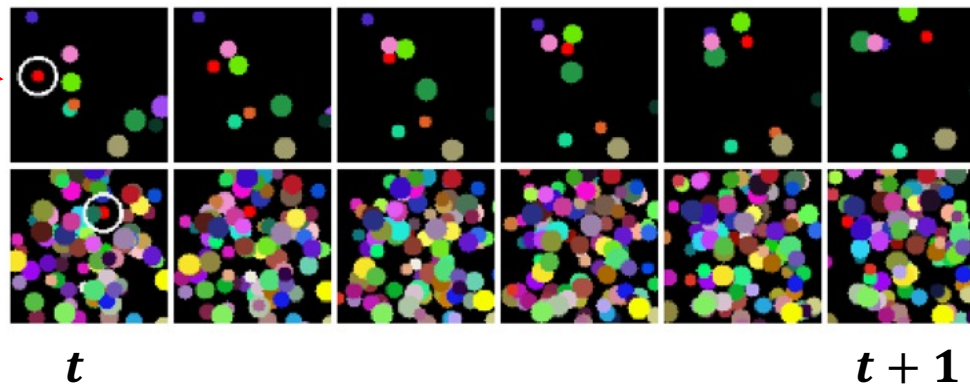
- prevents overfitting through regularization
- Avoids manual tuning and modeling

Estimators. Haarnoja et al. NeurIPS 2016

- Differentiable version of the Kalman Filter
- Uses Images as observations; learns a sensors that outputs state directly

$$g(I_t) = \mathbf{z}_t \approx \mathbf{x}_t$$

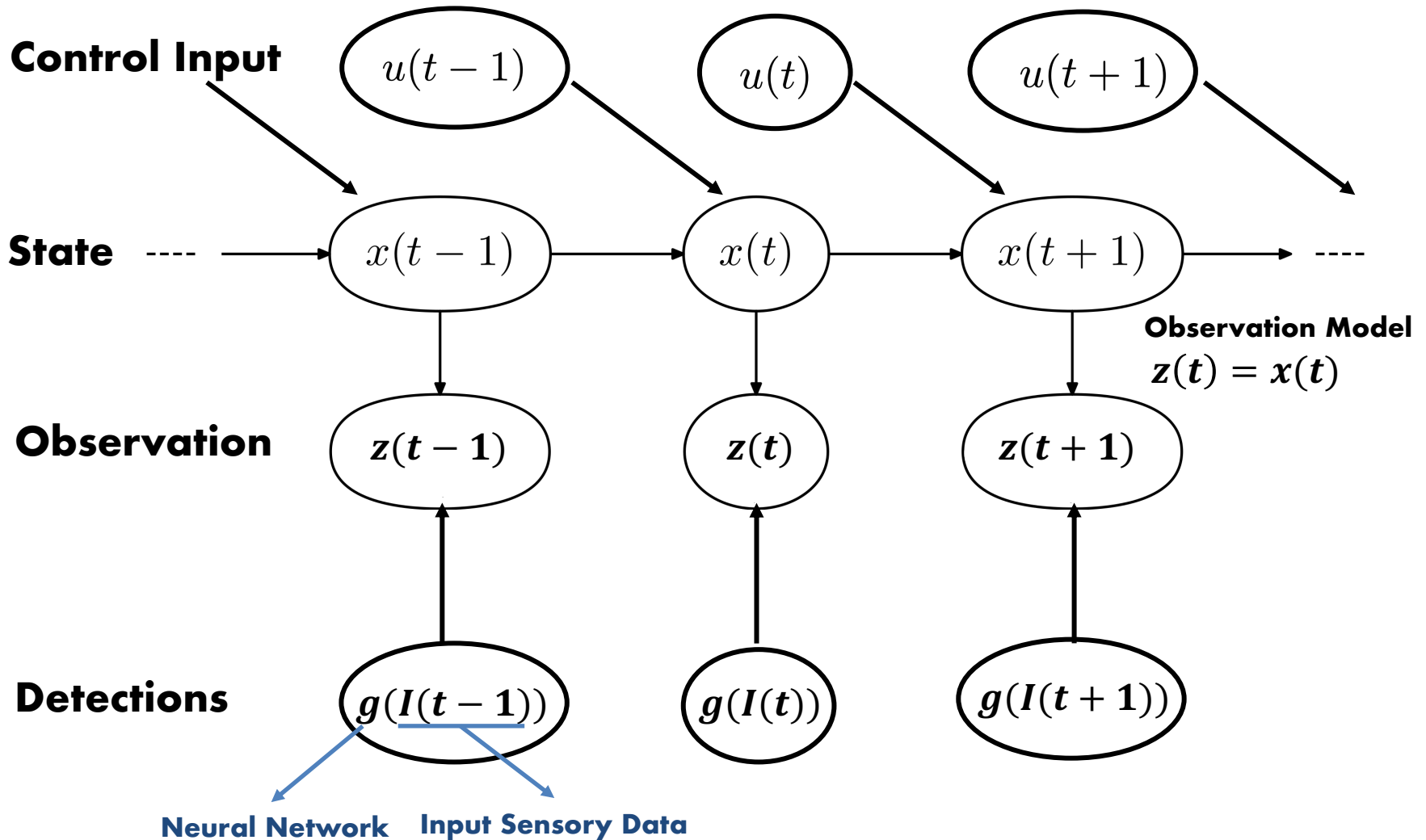
Track red
disk position



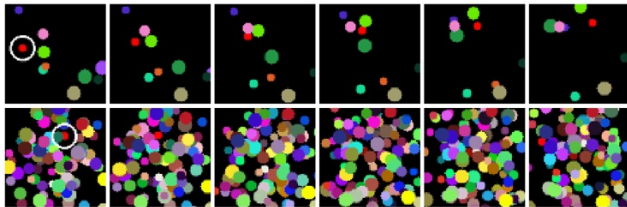
Example Sequence w/ few occlusions

Example Sequence w/ many occlusions

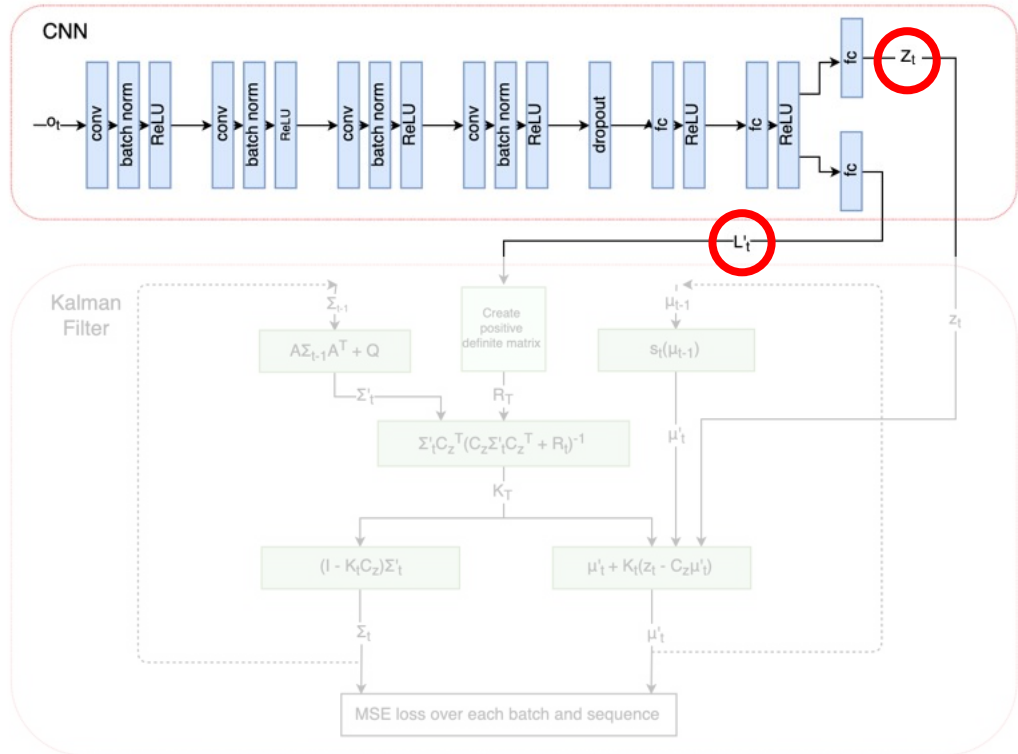
Tracking by Detection



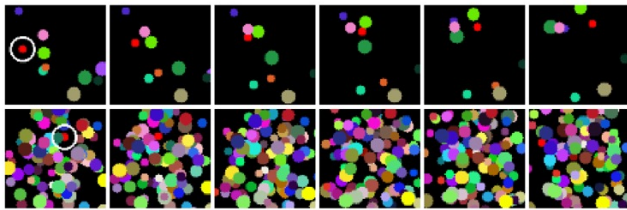
Differentiable Kalman Filter - Structure



$$g(I_t) = z_t \approx x_t$$

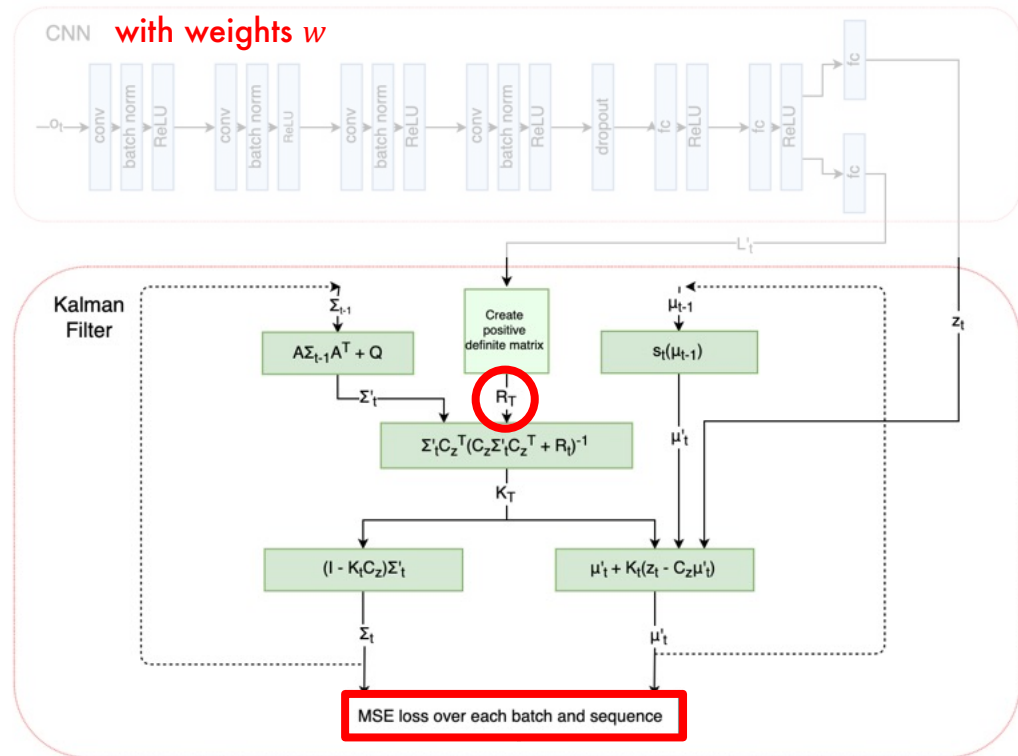


Differentiable Kalman Filter - Structure



↑
R is high if red disk is occluded

$$L'L^T = R$$



$$\frac{\delta Loss}{\delta w}$$

Differentiable Kalman Filter – Loss Function

Ground truth state Network weights

Belief

$$L(\mathbf{l}_{0\dots T}, \mu_{0\dots T}, \Sigma_{0\dots T}, \mathbf{w}) =$$

$$\lambda_1 \sum_{t=0}^T \frac{1}{2} \underbrace{((\mathbf{l}_t - \mu_t)^T \Sigma_t^{-1} (\mathbf{l}_t - \mu_t) + \log(|\Sigma_t|))}_{\text{Negative log likelihood of ground truth given current belief}} + \lambda_2 \sum_{t=0}^T \underbrace{\| (\mathbf{l}_t - \mu_t) \|_2}_{\text{Mean-Squared Error}} + \lambda_3 \underbrace{\| \mathbf{w} \|_2}_{\text{Regularization}}$$

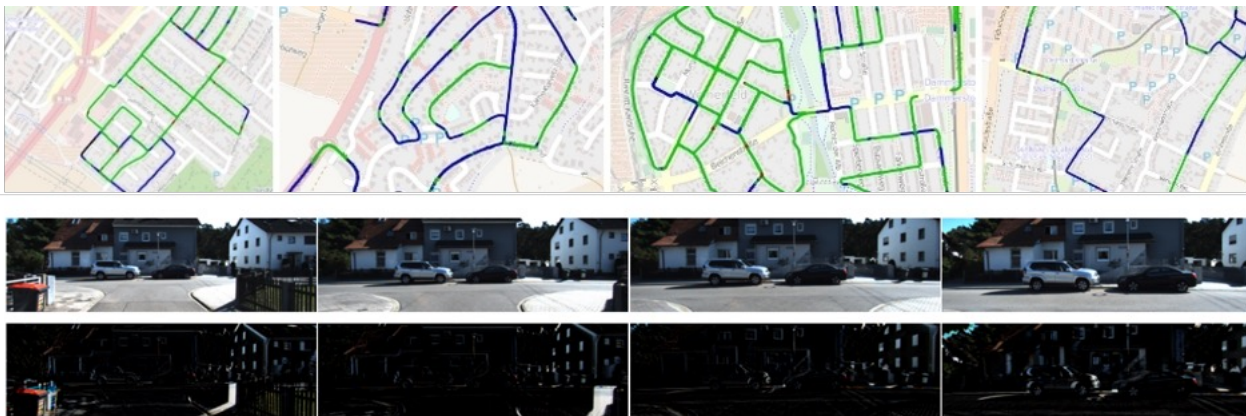
Differentiable Kalman Filter – Experiments and Baselines

Table 1: Benchmark Results

State Estimation Model	# Parameters	RMS test error $\pm \sigma$
feedforward model	7394	0.2322 \pm 0.1316
piecewise KF	7397	0.1160 \pm 0.0330
LSTM model (64 units)	33506	0.1407 \pm 0.1154
LSTM model (128 units)	92450	0.1423 \pm 0.1352
BKF (ours)	7493	0.0537 \pm 0.1235

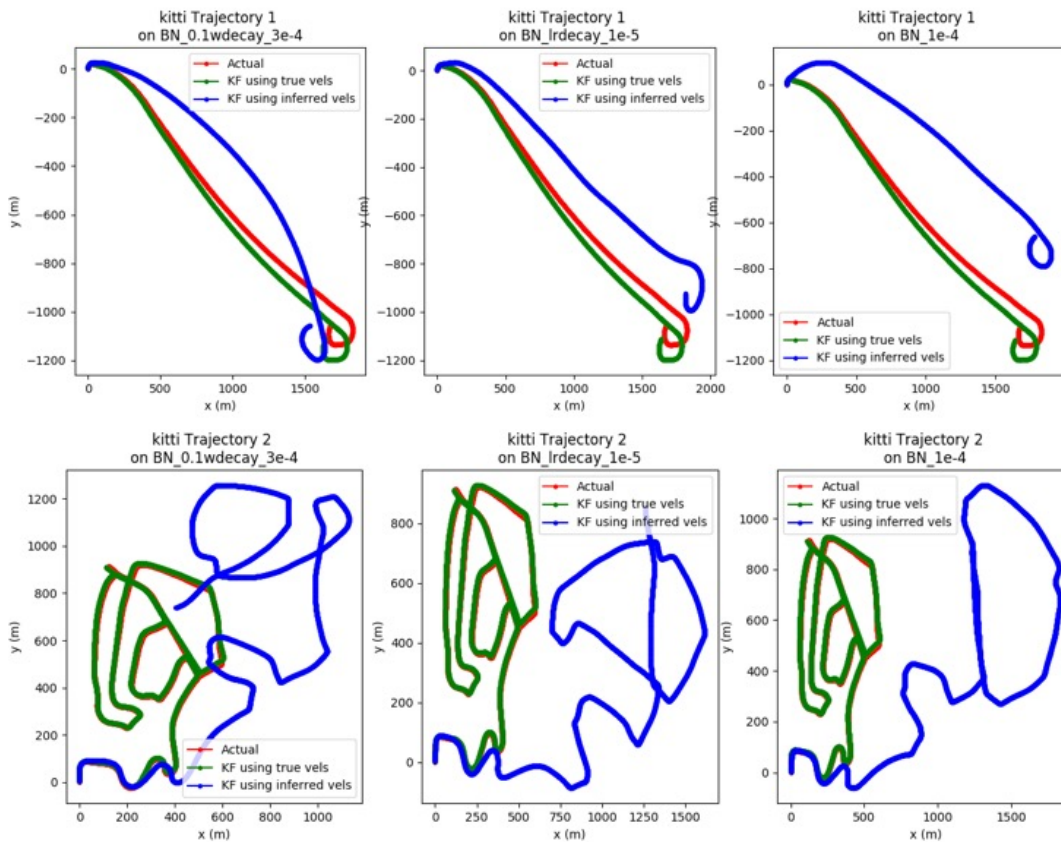
Differentiable Kalman Filter – Experiments and Baselines

- **Kitti – Visual Odometry Dataset**
- **22 stereo sequences with LIDAR**
 - **11 sequences with ground truth (GPS/IMU data)**
 - **11 sequences without ground truth (for evaluation)**



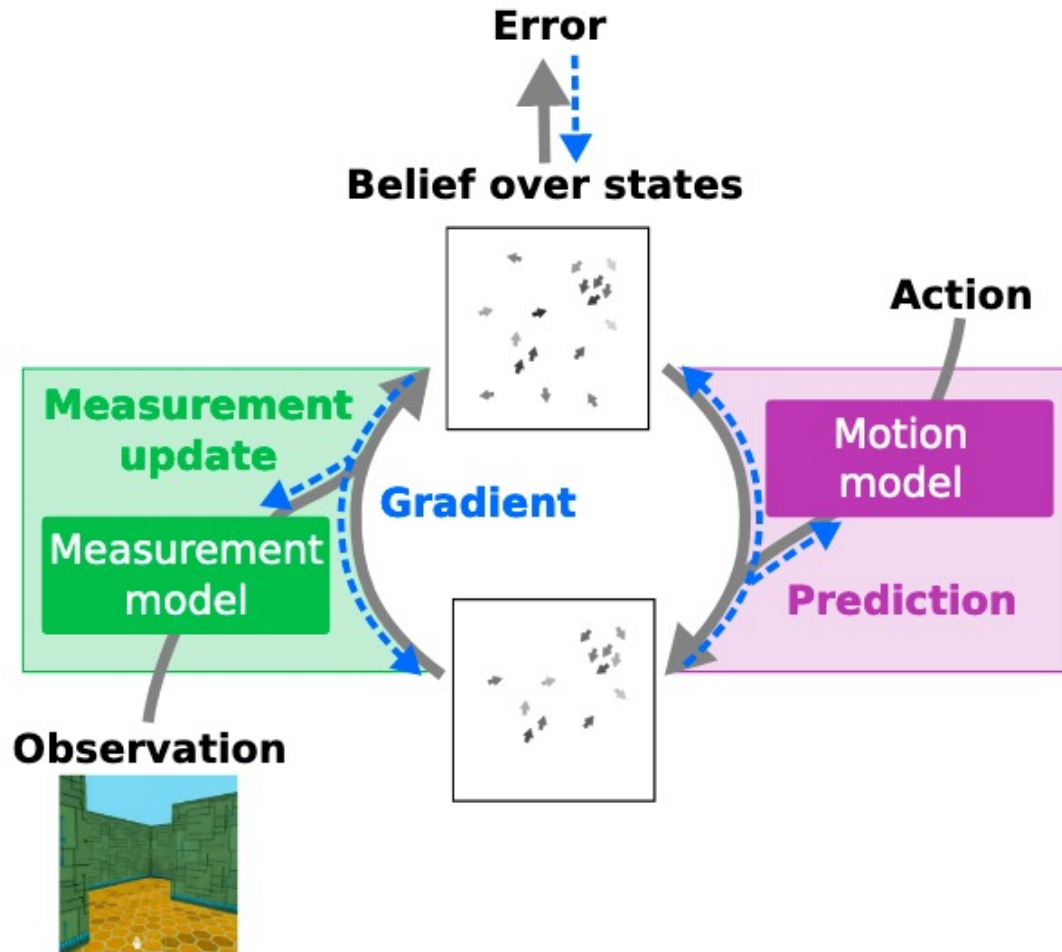
Differentiable Kalman Filter – Experiments and Baselines

Results reproduced by Claire Chen

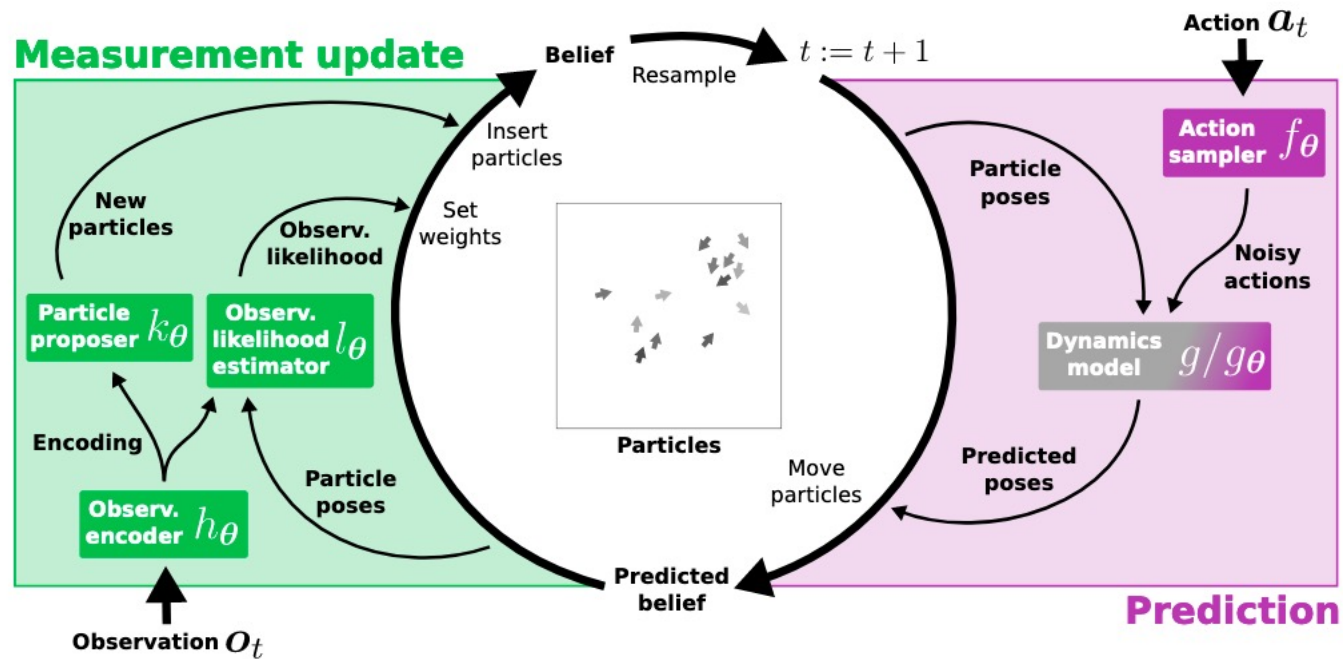


Blue – Result of BackpropKF
Red – Ground truth
Green – w/ Ground truth velocities

Differentiable Particle Filters: End-to-End Learning with Algorithmic Priors. Jonschkowski et al. RSS 2018.

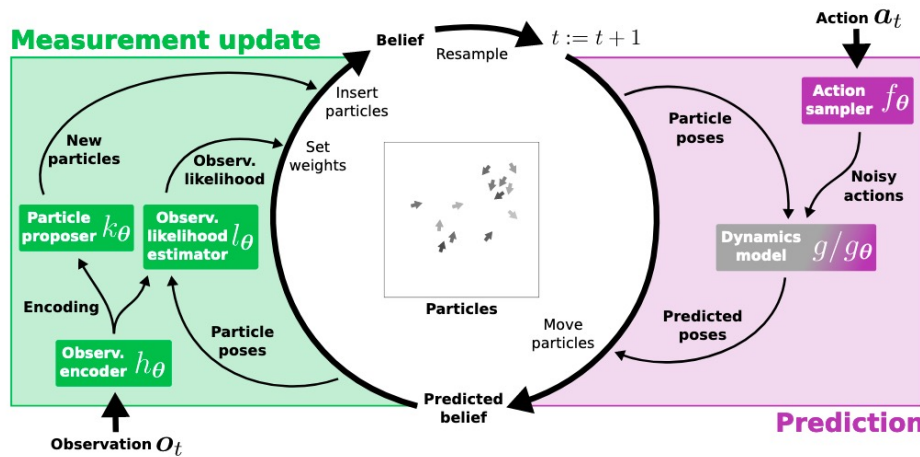


Differentiable Particle Filters: End-to-End Learning with Algorithmic Priors. Jonschkowski et al. RSS 2018.



Differentiable Particle Filters: End-to-End Learning with Algorithmic Priors. Jonschkowski et al. RSS 2018.

- Prediction Step



w_t weights

a_t actions

s_t states

o_t observations

$$\text{bel}(s_t) = (S_t, w_t)$$

$$\hat{a}_t^{[i]} = a_t + f_\theta(a_t, \epsilon^{[i]} \sim \mathcal{N}),$$

$$s_t^{[i]} = s_{t-1}^{[i]} + g(s_{t-1}^{[i]}, \hat{a}_t^{[i]}),$$

Action sampler

Dynamics model

Differentiable Particle Filters: End-to-End Learning with Algorithmic Priors. Jonschkowski et al. RSS 2018.

- Measurement Update

w_t weights

a_t actions

s_t states

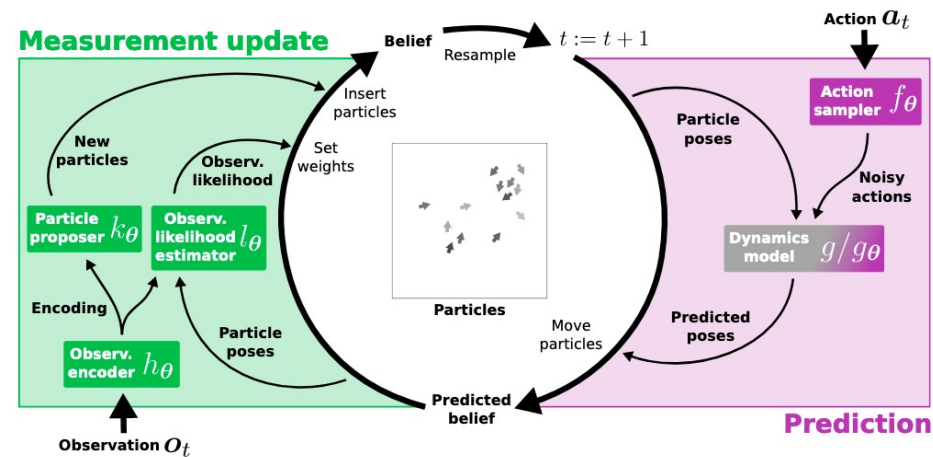
o_t observations

e_t encoding

$$e_t = h_{\theta}(o_t),$$

$$s_t^{[i]} = k_{\theta}(e_t, \delta^{[i]} \sim B),$$

$$w_t^{[i]} = l_{\theta}(e_t, s_t^{[i]}),$$



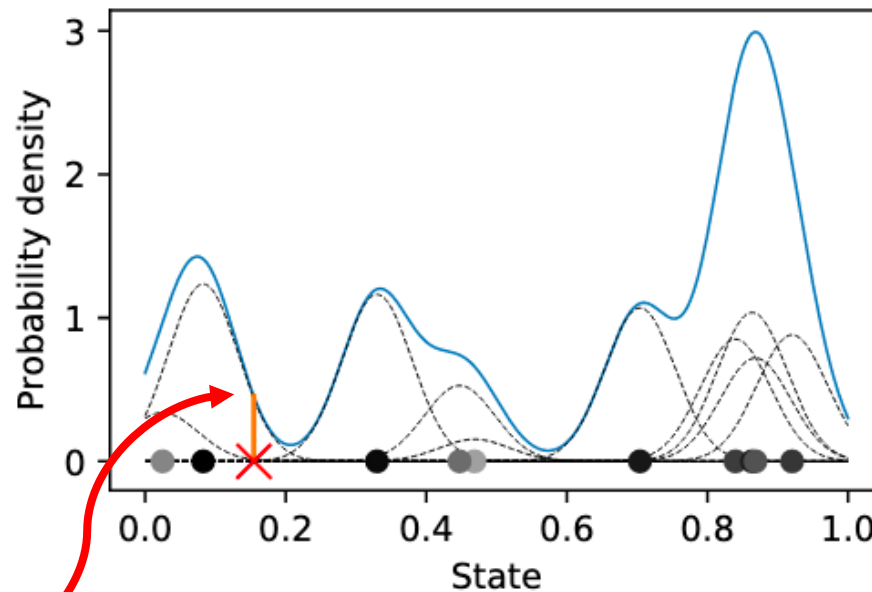
Observation Encoder

Particle Proposer

Observation likelihood estimator

Differentiable Particle Filter – Loss Function

- Supervised learning given data $\mathbf{o}_{1:T}, \mathbf{a}_{1:T}, \mathbf{s}_{1:T}^*$



$\text{bel}(\mathbf{s}_t^*)$

$$\theta^* = \operatorname{argmin}_{\theta} -\log E_t[\text{bel}(\mathbf{s}_t^*; \theta)].$$

Maximizing the belief at the true state

Differentiable Particle Filter – Experiments and Baselines

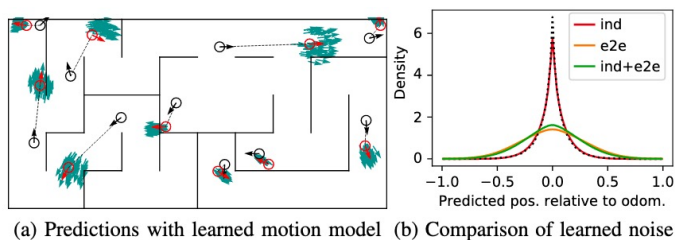


Fig. 5: **Learned motion model.** (a) shows predictions (cyan) of the state (red) from the previous state (black). (b) compares prediction uncertainty in x to true odometry noise (dotted line).

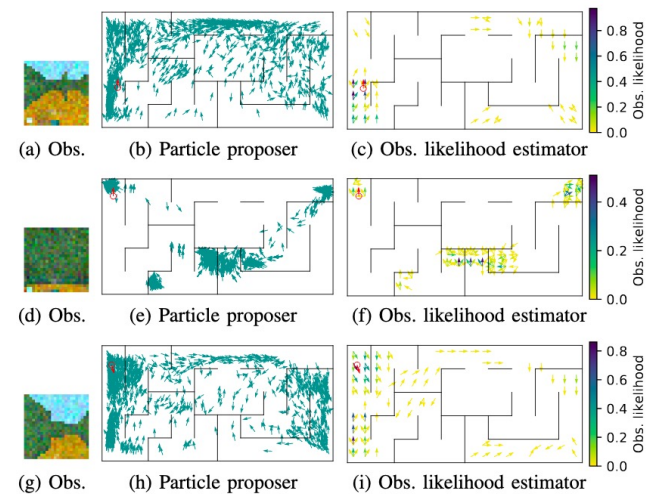


Fig. 6: **Learned measurement model.** Observations, corresponding model output, and true state (red). To remove clutter, the observation likelihood only shows above average states.

Differentiable Particle Filter – Experiments and Baselines

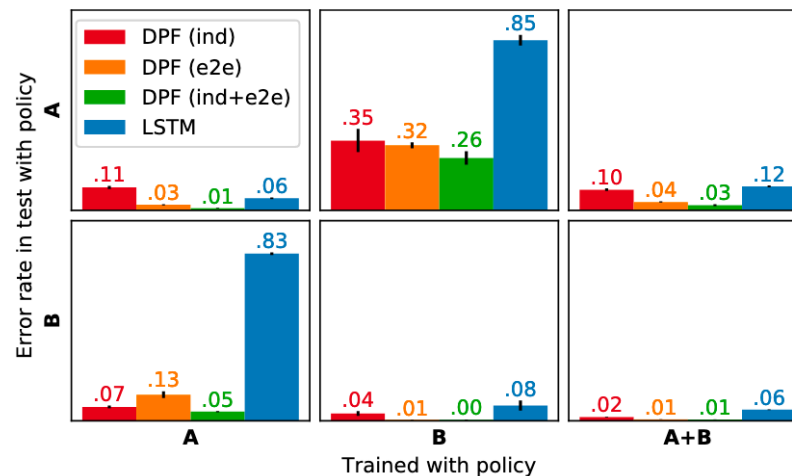
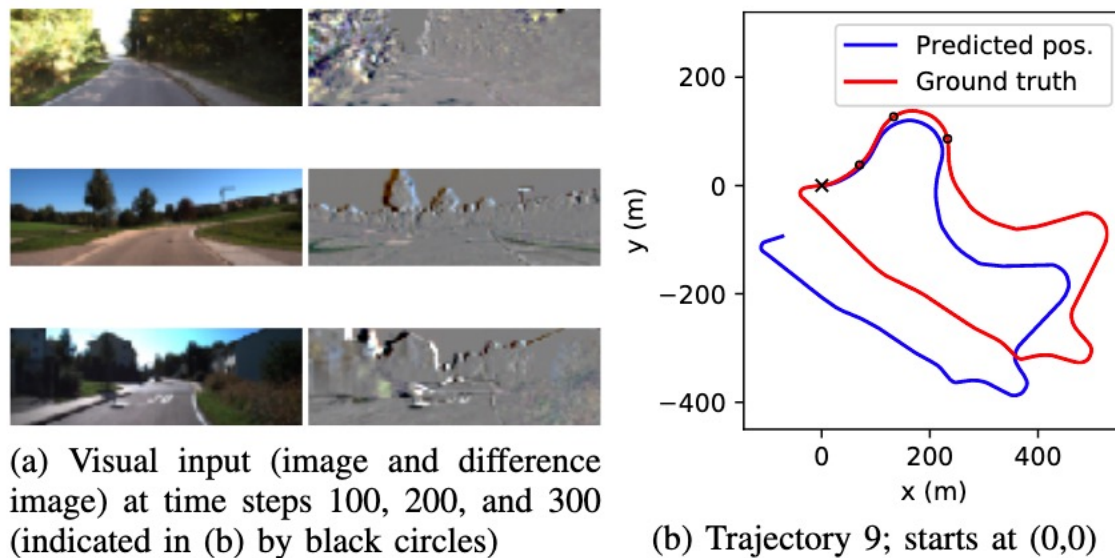


Fig. 9: **Generalization between policies** in maze 2. A: heuristic exploration policy, B: shortest path policy. Methods were trained using 1000 trajectories from A, B, or an equal mix of A and B, and then tested with policy A or B.

Differentiable Particle Filter – Experiments and Baselines



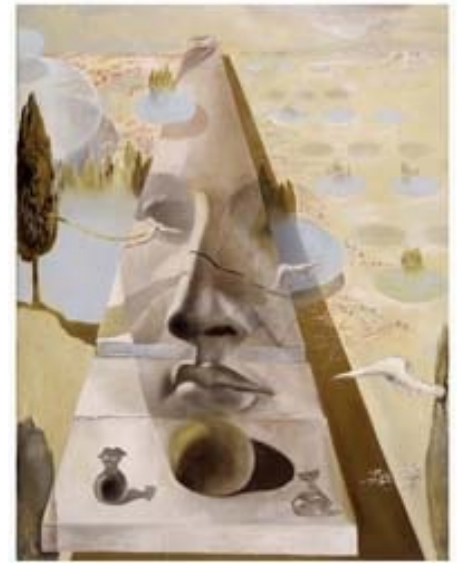
(a) Visual input (image and difference image) at time steps 100, 200, and 300 (indicated in (b) by black circles)

(b) Trajectory 9; starts at (0,0)

Fig. 10: Visual odometry with DPFs. Example test trajectory

CS231A

Computer Vision: From 3D Reconstruction to Recognition



Next lecture:

Neural Radiance Fields for Novel View
Synthesis