

## Decentralized multi-sensor fusion

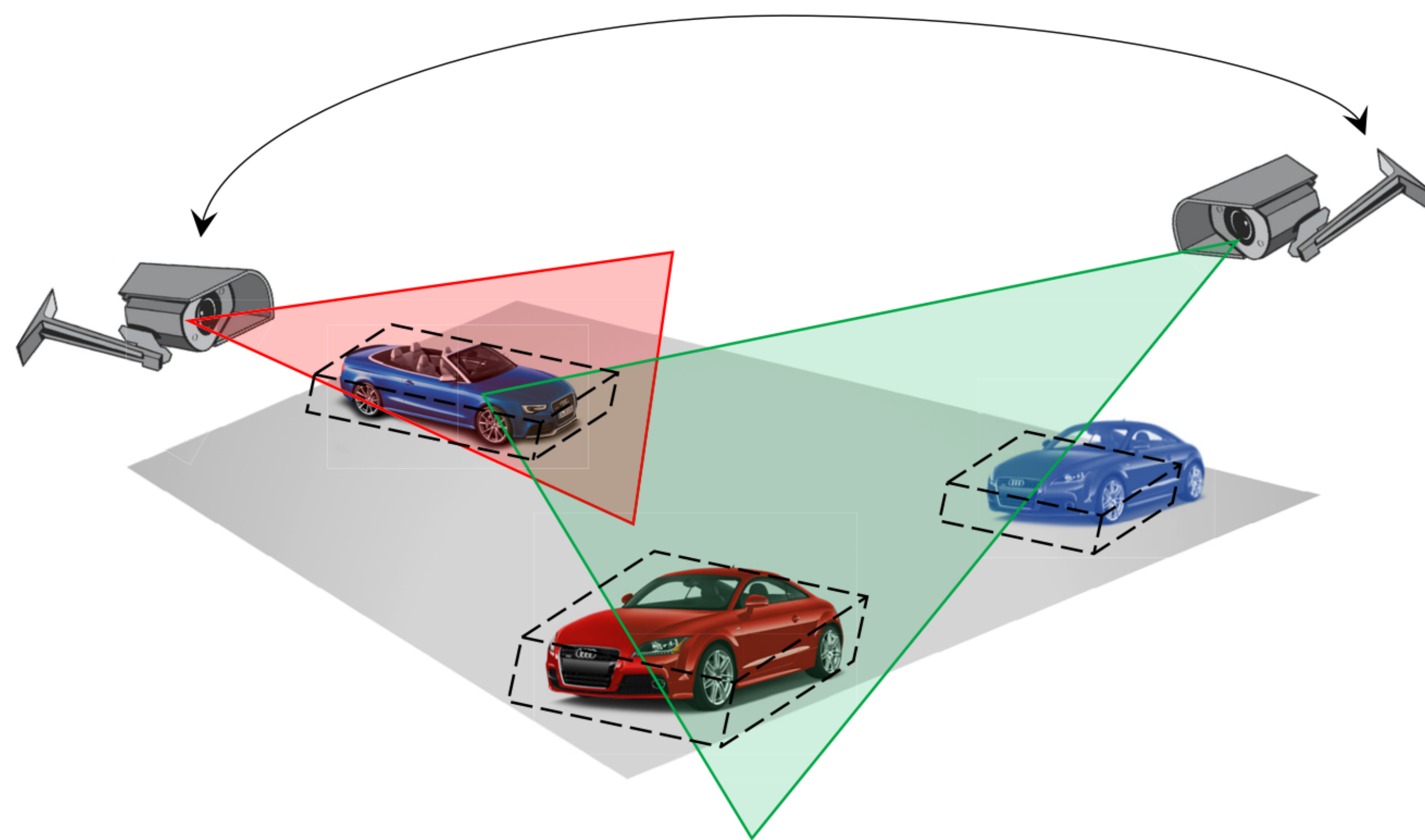


Figure 1. Scenario with three objects observed by two sensors with partially overlapping field-of-views. Image taken from [1].

- Multiple sensors with possibly different field-of-views (FoVs), where tracking is performed locally at each sensor, are utilized to estimate the moving objects in the region of interest.
- To leverage the information available at local sensors, the local multi-object densities (MODs) obtained at each sensor need to be fused to obtain a global MOD.

## Limitations of model-based methods

- Current model-based methods, in theory, do not provide Bayesian optimal fusion results.
- Current fusion methods only use the local multi-object posterior densities at the current time while ignoring information regarding the objects' previous states.
- It is difficult to develop a model-based method that leverages all the uncertainties captured in the MODs on sets of trajectories [2].

## Contributions

- The first deep learning-based solution to the fusion of MODs.
- MT3 [3], a transformer-based deep neural network, is used to fuse local MODs on sets of trajectories to obtain a global MOD that describes the set of current objects.

## Problem formulation

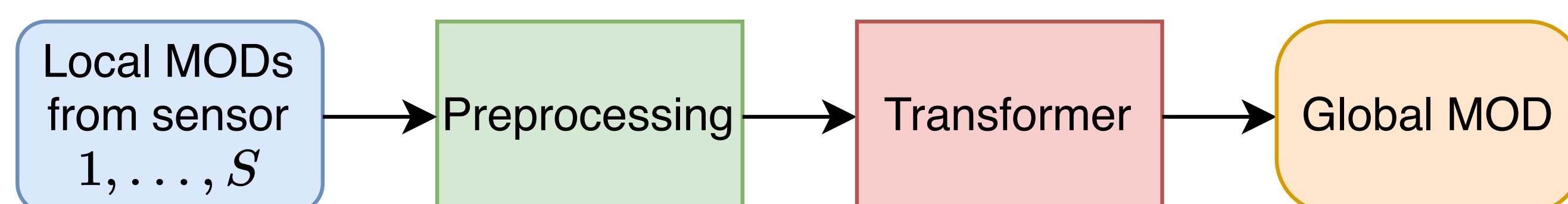


Figure 2. Diagram of fusion of MODs using transformer.

- The local MOD  $f^s(\mathbf{X}_T | \mathbf{z}_{1:T}^s)$  at the  $s$ -th sensor represents the posterior density of the set  $\mathbf{X}_T$  of trajectories in the time interval  $1 : T$  in the form of multi-Bernoulli (MB) [2].
- The global MOD  $f(\mathbf{x}_T | \mathbf{z}_{1:T}^1, \dots, \mathbf{z}_{1:T}^S)$  represents the set of detected objects  $\mathbf{x}_T$  at time step  $T$  in the form of MB, where an existence probability and a Gaussian single-object density characterize each Bernoulli component.

## Deep MOD fusion using transformer

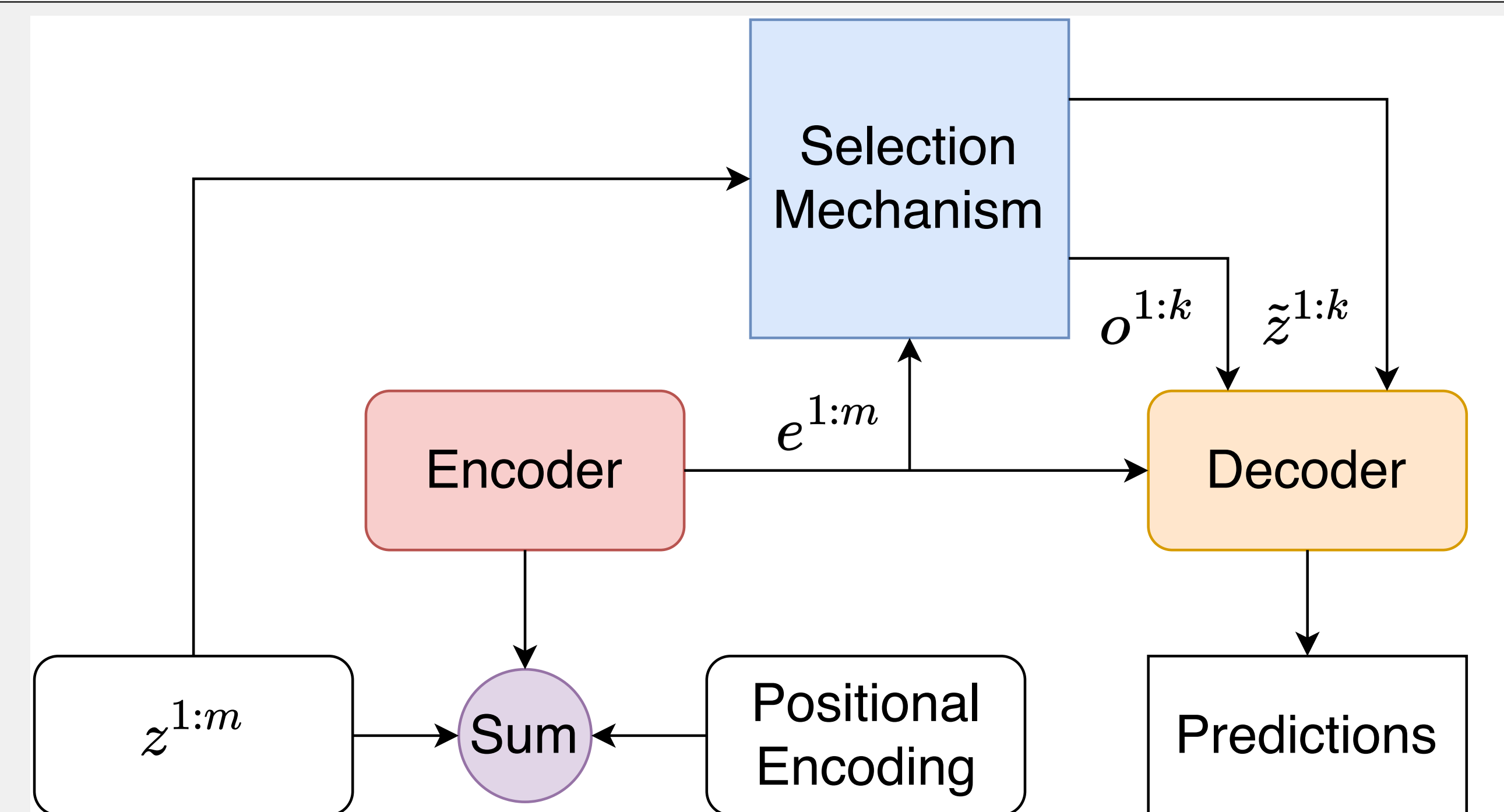


Figure 3. The structure of MT3 [3].

- Preprocessing** Concatenate parameters of the  $i$ -th Bernoulli at the  $s$ -th sensor:
  - Existence probability  $r^{s,i}$ .
  - Trajectory start time  $t^{s,i}$ .
  - Probability mass function of trajectory length  $w_j^{s,i}, j \in \{1, \dots, \ell^{s,i}\}$ .
  - Mean of state sequence  $x_{1:\ell^{s,i}}^{s,i} \in \mathbb{R}^{4\ell^{s,i}}$ .
  - (Reduced) covariance of state sequence  $P_{1:\ell^{s,i}}^{s,i} \in \mathbb{R}^{4\ell^{s,i} \times 4\ell^{s,i}}$ .
 Include sensor position information when sensors are mobile.
- Positional encoding** Encode time step  $t \in \{t^{s,i}, \dots, t^{s,i} + \ell^{s,i} - 1\}$ , trajectory index  $i \in \{1, \dots, n^s\}$  and sensor index  $s \in \{1, \dots, S\}$  into the input embedding.
- Loss** Negative log-likelihood (NLL) of the MB density evaluated at the ground truth [4].

## Scenario setup

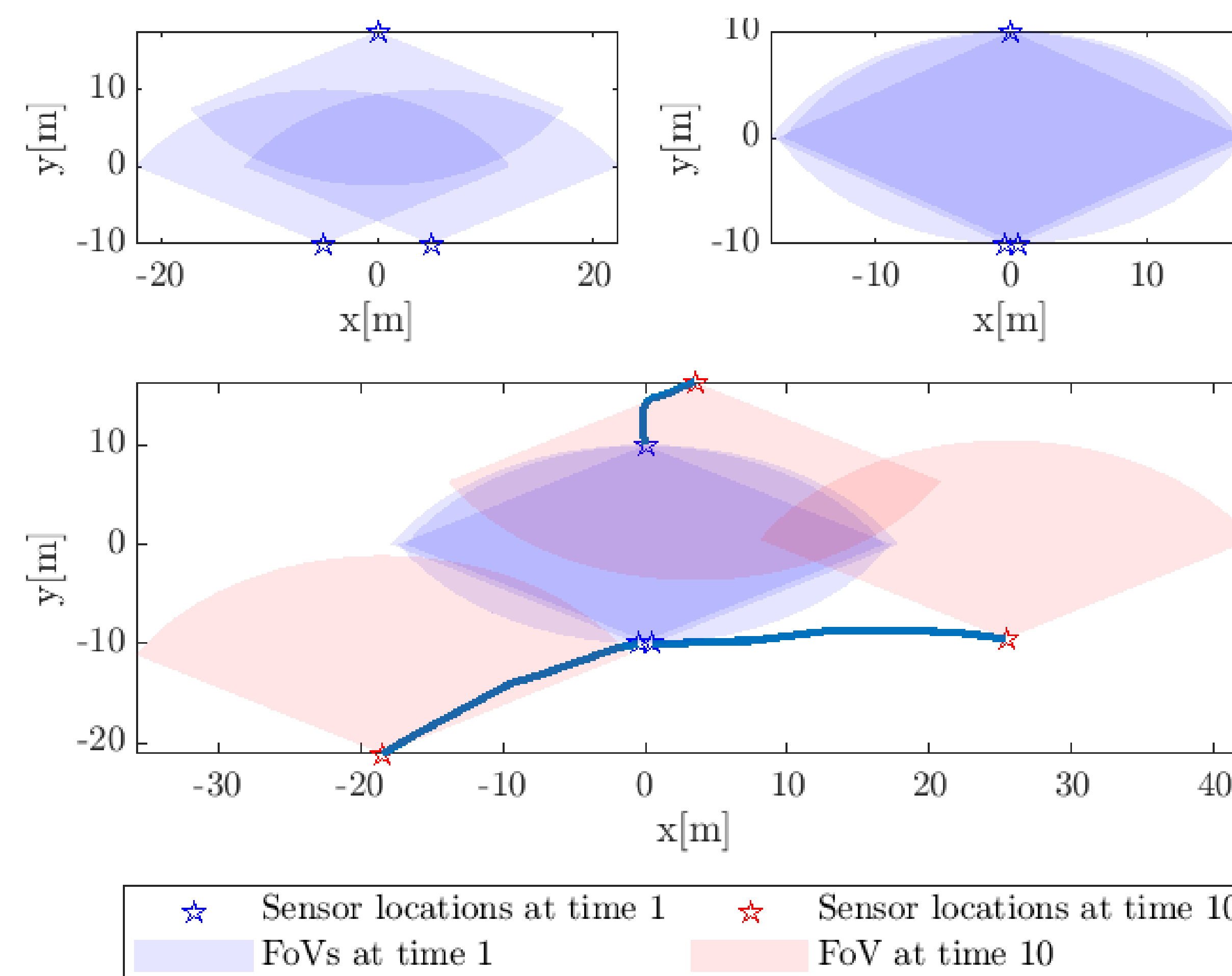


Figure 4. An illustration of scenario 1 (top left), scenario 2 (top right), and a sample plot of scenario 3 (bottom), sensor trajectories are shown in solid lines.

## Results on synthetic data

Scenario	1		2		3	
Method	Bayesian	MT3	Bayesian	MT3	Bayesian	MT3
GOSPA-total	1.618	<b>1.373</b>	1.440	<b>1.413</b>	1.294	<b>0.885</b>
GOSPA-loc	0.596	0.959	0.693	1.040	0.256	0.597
GOSPA-miss	0.058	0.051	0.076	0.031	0.003	0.025
GOSPA-false	0.964	0.363	0.671	0.342	1.035	0.263
NLL-total	12.207	<b>1.195</b>	2.695	<b>1.833</b>	13.105	<b>0.429</b>
NLL-loc	11.766	0.572	2.103	1.169	12.656	0.089
NLL-miss	0.320	0.304	0.253	0.411	0.428	0.137
NLL-false	0.121	0.319	0.340	0.253	0.021	0.202

Table 1. Performance comparison with [1] in terms of GOSPA [5] and NLL [4].

Results show that MT3 [3] outperforms the model-based Bayesian fusion method [1] in terms of both GOSPA and NLL in all the scenarios.

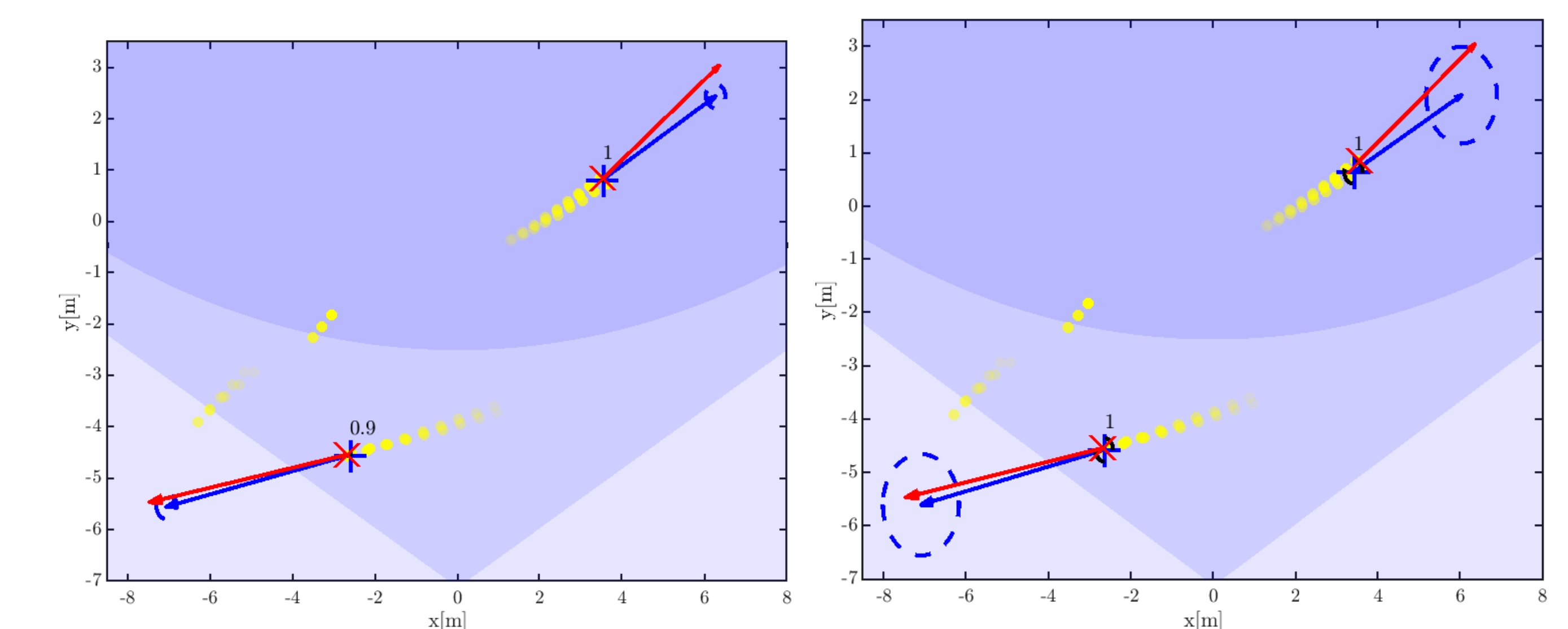


Figure 5. Sample plot of Bayesian method (Left) and MT3 (Right). Yellow-filled circles indicate estimated positions obtained from local filters. The ground truth positions/velocities at the current time are shown in red crosses/arrows, respectively, while predicted positions/velocities are shown in blue plus signs/arrows. The blue/black dashed ellipses represent the  $3\text{-}\sigma$  level of predicted velocities/positions.

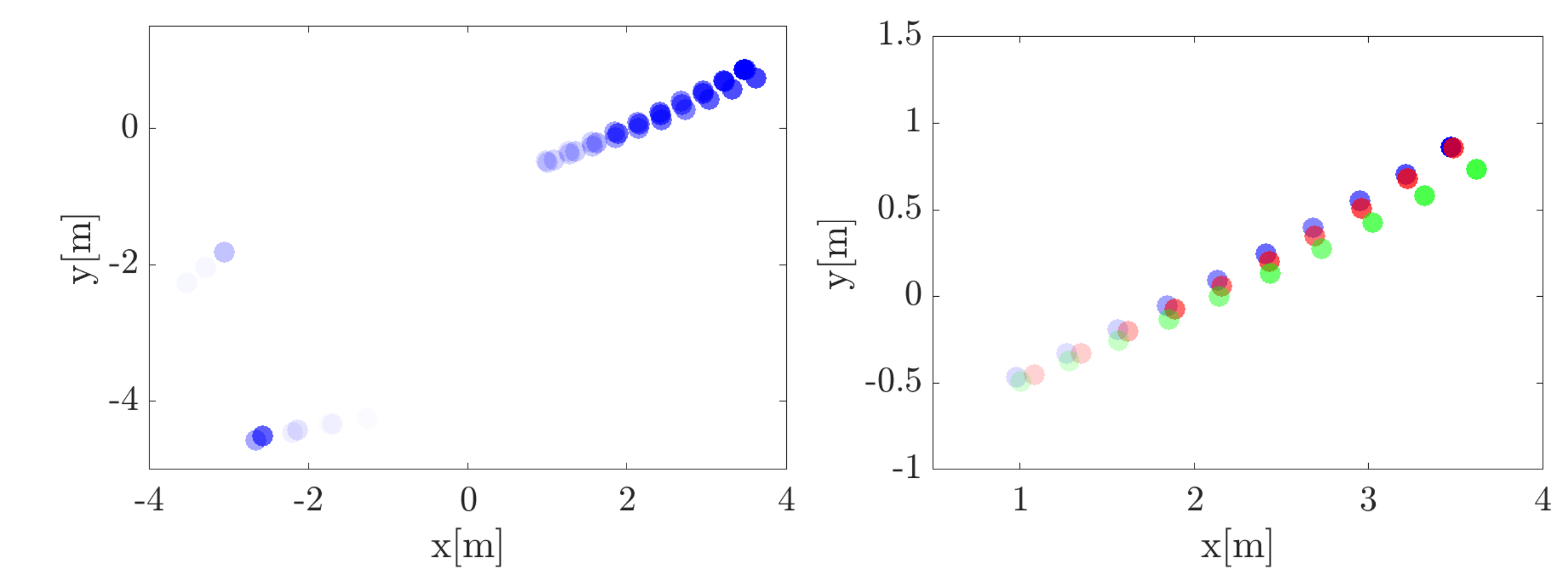


Figure 6. Left: attention maps of the predictions of MT3 in blue-filled circles. Right: corresponding trajectory estimates from different sensors, indicated using different colours.

- M. Fröhle, K. Granström, and H. Wymeersch, "Decentralized Poisson multi-Bernoulli filtering for vehicle tracking," *IEEE Access*, vol. 8, pp. 126414–126427, 2020.
- Á. F. García-Fernández, L. Svensson, J. L. Williams, Y. Xia, and K. Granström, "Trajectory Poisson multi-Bernoulli filters," *IEEE Trans. on Signal Process.*, vol. 68, pp. 4933–4945, 2020.
- J. Pinto, G. Hess, W. Ljungbergh, Y. Xia, L. Svensson, and H. Wymeersch, "Next generation multitarget trackers: Random finite set methods vs transformer-based deep learning," in *24th International Conference on Information Fusion*, IEEE, 2021.
- J. Pinto, Y. Xia, L. Svensson, and H. Wymeersch, "An uncertainty-aware performance measure for multi-object tracking," *IEEE Signal Process Lett.*, vol. 28, pp. 1689–1693, 2021.
- A. S. Rahmathullah, Á. F. García-Fernández, and L. Svensson, "Generalized optimal sub-pattern assignment metric," in *20th International Conference on Information Fusion*, IEEE, 2017.