

Multi-Object Tracking with Machine Learning based trajectories prediction and soft track-association

Christian Alexander Holz, Christian Bader, Matthias Drüppel

Abstract—This study introduces Machine Learning (ML)-based methodologies for Multi-Object Tracking (MOT) tailored to Advanced Driver Assistance Systems (ADAS), addressing the growing complexity and precision requirements in the automotive sector. We propose and evaluate three novel neural networks (NN) designed for MOT: (i) the Single Prediction Network (SPENT) for trajectory forecasting, (ii) the Single Association Network (SANT) for mapping single sensor inputs to existing tracks, and (iii) the Multi-Association Network (MANTA) for associating multiple sensor inputs to multiple tracks. These ML models are integrated with a traditional Kalman filter framework to enhance MOT performance while preserving the system's modularity and interpretability. Our evaluation demonstrates significant improvements: SPENT reduces the Root Mean Square Error (RMSE) by 50% on the KITTI tracking dataset compared to a standard Kalman Filter; SANT and MANTA achieve a 95% validation accuracy in sensor object assignment to tracks. These results highlight the efficacy of ML integration in enhancing the robustness, modularity, and maintainability of ADAS tracking systems.

Index Terms—Article submission, IEEE, IEEEtran, journal, LATEX, paper, template, typesetting.

I. INTRODUCTION

THE ongoing evolution of Advanced Driver Assistance Systems has brought the need for precise and reliable Multi Object Tracking into the spotlight [1]–[8]. In complex and dynamic environments, as encountered in urban traffic, it is crucial to simultaneously and accurately capture the positions and movements of multiple objects - a key challenge in computer vision (CV) for automated driving.

In the commonly used Tracking-by-Detection (TbD) paradigm, a tracker fuses detected sensor objects (SO) to create consistent object tracks over time. A crucial step within this paradigm is the association of the incoming measured SO with their corresponding existing object tracks to update their properties. If no association can be made, new object tracks must be initialized.

Tracking frameworks form the heart of ADAS systems that are used in millions of vehicles around the globe. The vast majority of these frameworks rely on classical approaches such as the Kalman filter (KF) or its variants [1]–[3]. These classical tracking theories have the great benefit of being modular and interpretable. The task is split into clearly separated subtasks

such as the prediction of currently tracked objects and the association with newly measured ones. Furthermore, the math and the theory itself is often clean and comprehensible. However, in the automotive industry tracking systems are not developed for a single model, but are usually used as a platform and are deployed to a variety of different car models with different sensor sets, different installation heights of the sensors, used in different countries and must always perform for a wide range of driving scenarios. This usually leads to poor performance in certain scenes for a specific system. With classical systems that not learn directly from data, these situations are often solved by the implementation of heuristics and by tweaking parameters of the tracking system for the troublesome scenes. But from a software engineering point of view hand-engineered parameters and heuristics are extremely hard to maintain and develop further for new scenarios and new system configurations. Moreover, hand-engineered classical approaches can reach its limits for the overall performance and in challenging situations that would require a more automated and reproducible development strategy.

In this work, we propose a data driven approach to tracking frameworks, which would allow the same system to be fine tuned for specific configurations relying only on data, thus increasing maintainability and adaptability. We do this preserving one of the biggest strengths of classical approaches: its modularity, by replacing only single tracking components with ML models.

In contrast to the Kalman filter [3], our Prediction Network is capable of predicting the state of individual objects without the need for a predefined state or prediction model at runtime. The self-learning, data driven approach enables adaptability to various scenarios and the ability to effectively handle non-linearities. Many conventional tracking systems rely on static methods for data association. Commonly used algorithms like the Hungarian algorithm (HA) [9] require heuristics and fixed thresholds. Our Single Association Network replaces the calculation of a distance metric for the corresponding assignment by employing machine learning. Both, the Prediction Network and SANT can be developed and evaluated as stand-alone models. For a throughout evaluation, we proceed by integrating them into an existing tracking system and demonstrate their performance through multiple tests and comparisons with established methods.

This work provides new insights and advancement in the development of ADAS, contributing to the further evolution of technologies for autonomous driving (AD).

C. Holz is with Daimler Truck AG, Research and Advanced Development, Stuttgart, Germany

C. Bader is with Daimler Truck AG, Research and Advanced Development, Stuttgart, Germany

M. Drüppel is with the Center for Artificial Intelligence, Duale Hochschule Baden-Württemberg (DHBW), Stuttgart, Germany

II. RELATED WORK

A. *Tracked object prediction*

One fundamental problem in Tracking-by-Detection frameworks is the prediction of the states of the already tracked objects. In many approaches Kalman filters and their variants have proven to be effective for state prediction [1], [2], [10]. However, they reach their limits in more complex scenarios, particularly in the presence of non-linear motion patterns and interactions among multiple objects [11]–[13]. [11] discusses the limitations of Kalman Filters in nonlinear and non-Gaussian scenarios and introduces Particle Filters as an alternative solution. [12] proposes the Unscented Kalman Filter (UKF) as an extension to the standard Kalman Filter to handle non-linear motion models more effectively. [13] introduces the Gaussian Mixture Model (GMM) for tracking multiple objects in a cluttered environment.

In this work, we introduce a novel MOT approach that leverages Machine Learning to overcome these challenges. We specifically focus on the development and implementation of Neural Networks, which can enable more precise and flexible data-driven object tracking.

B. *Association*

Another fundamental problem for a TbD tracker is the data association. For this, some approaches only consider the current state of the tracks, while others integrate temporal information such as the track history. This aggregate of temporal information can be done for example using attention mechanisms as used in [8], [14] or recurrent neural networks (RNNs) as developed in [4], [5]. The latter is what we are also pursuing in this work. Similar to the problem statement in Mertz et al [5], the aim of our work is to develop a data-based approach that can learn to completely solve the combinatorial non deterministic polynomial time (NP) hard optimization problem of data association. Mertz et al [5] use a distance matrix based on the Euclidean distance measure as input for the developed association network, thus replacing an association algorithm such as the Hungarian Algorithm. In the context of this work, we put forward the hypothesis that a Gated Recurrent Unit (GRU)-based association network can be designed and trained using an undefined distance measure.

C. *Real time applications*

ADAS are embedded real-time applications where it is crucial to predict the state of objects immediately after their detection as described in [1], [15]. This rules out offline tracking methods as presented in [16] that process the entire video material at once in a batch process. Therefore, most recent approaches for tracking multiple objects rely on online methods that do not depend on future image information. Online methods use various features to estimate the similarity between the recognized objects and the existing tracks. This can be done on the basis of their predicted positions or even similarities in appearance.

a) *Kalman Filter based tracker:* [1], [1]–[3], [10]–[13] propose a Kalman Filter based TbD multi-object tracker. [1] presents a simple and efficient multi-object tracking approach based on the KF and the Hungarian Algorithm. [2] presents a MOT approach based on simple visual cues. The authors contend that many existing multi-object trackers are too complex and require a large amount of computational resources. Instead, they propose a simpler approach based on basic visual features such as color, shape and motion. These visual cues are used to track objects at the image level and make associations between frames.

b) *Recurrent Neural Network based Tracker:* Similar to our methodologies, studies by [4], [5], [17] introduce RNN-based approaches for online multi-target tracking. Specifically, the work by Mertz et al. [5] focuses on data association within a TbD framework. Their proposed DeepDA model, an LSTM-based Deep Data Association Network, is designed to learn and execute the task of associating objects across frames. This model's ability to discern association patterns directly from data enables the achievement of robust and reliable tracking outcomes, even in environments with significant disturbances. Mertz et al. employ a distance matrix, derived from the Euclidean distance measure, as the input for the DeepDA network. This innovative approach effectively supersedes traditional association algorithms, such as the Hungarian Algorithm. It is inferred that the Euclidean distance measure served as a foundation not only for generating the ground truth (GT) training data (i.e., distance matrices) but also for the subsequent evaluation process. However, this is not explicitly stated. It can therefore be argued that this preprocessing step deprives the network opportunity to follow a different association logic or to learn it data based.

c) *Attention Mechanism based Tracker:* Each of the papers [6]–[8], [14], [18] present approaches for a tracker that utilize the attention mechanism [19], for example to compute soft data association [8]. The main research focus of the paper [8] is on soft data association, which enables the tracker to make probabilistic associations between objects and account for uncertainties in the associations. Soft data association in the SoDA model works by using attention mechanisms to aggregate information from all detections in a given temporal window. This allows the model to learn long-term and highly interactive relationships between detections and tracks from large datasets without using complex heuristics and hyperparameters.

III. OVERVIEW OF OUR PROPOSED MODELS

Our primary contribution is the development and evaluation of three NN that we label: (i) the Single Prediction Network (SPENT), (ii) the Single Association Network (SANT), and (iii) the Multi-Association Network (MANTa). Figure 1 provides an overview of our approach in which the association network can be implemented using either SANT or MANTa. The input for the proposed Prediction Network (i) are the sensor objects in every time step. If new objects are detected, these are stored in a fixed-dimensional state vector that contains information such as object position, orientation

and dimension. The Prediction Network predicts all **vectors** to the next **time step** where they are used as input to either the SANT or MANTa association networks. These **provide** the association matrix, that is used to update the tracks using the corresponding sensor objects or create new tracked objects. The Track Management can then decide to delete tracks, that were not updated for a specific amount of time and **send out** tracks to the next higher software component that have been confirmed by sensor objects.

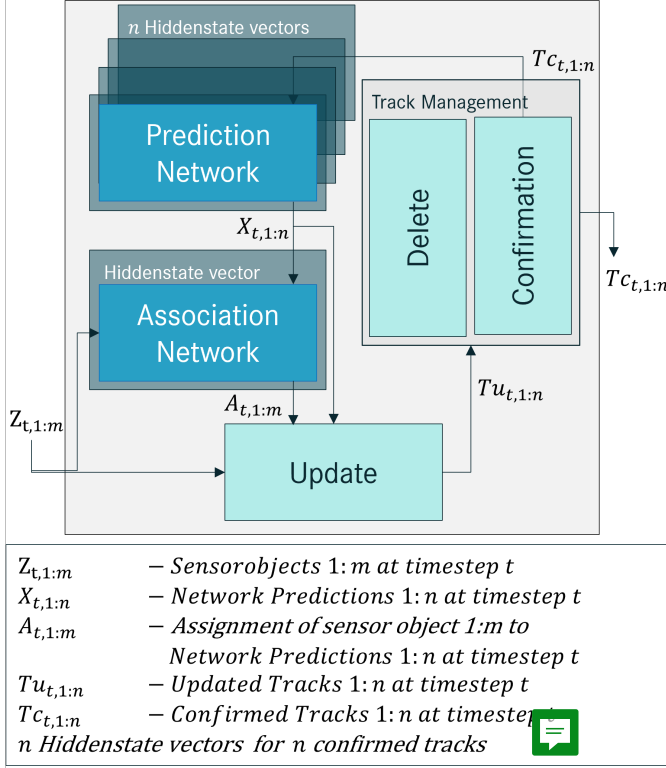


Fig. 1. Schematic representation of the two integrated networks (blue) in a tracking-by-detection framework. The Association Network can either be filled with the our SANT (single) or MANTa (multi) association model.

IV. TRACKING WITH PREDICTION AND ASSOCIATION NETWORKS

We apply the tracking-by-detection paradigm, in which a tracker fuses object detections (**sensor objects**) to generate object tracks that are consistent over time. [10], for example, provides a framework for analyzing tracking approaches that follow the TbD paradigm. This was used accordingly in this study. To enable our models to incorporate temporal information, we use LSTM and bidirectional Long Short-Term Memory (BiLSTM) networks. Both for the prediction and the association (SANT) of the sensor objects to the existing tracks.

A. Single Prediction Network (SPENT)

In our methodology, we leverage the hidden states of the LSTM layer as a dedicated information repository for each object. The Prediction Network is architected for open-loop operation, signifying that it forecasts future state values based on data previously received. Within the context of an online

MOT process, this enables the network to prognosticate the most probable subsequent state values by utilizing the state values received for a specific sensor object's state vector.

As depicted in Figure 2, the schematic illustration of the generic structure of the Prediction Network elucidates how the architecture is adeptly designed to address the challenges of real-time state prediction. This representation highlights the strategic deployment of the LSTM layer for storing and processing object-specific information, facilitating accurate and timely predictions of object states.

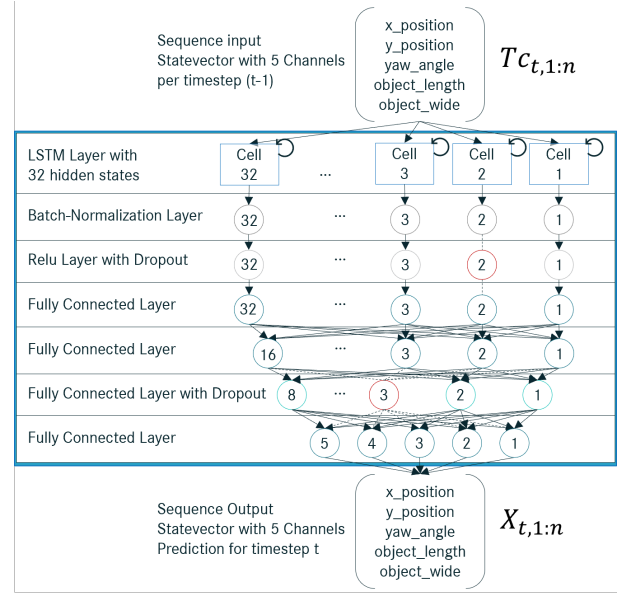


Fig. 2. Schematic representation of the generic prediction network structure.

a) *Network architecture:* The foundational layer of our Prediction Network is constituted by an LSTM layer, wherein the internal states — known as hidden states — are dynamically updated at each timestep in response to incoming measurement data. This iterative updating mechanism facilitates a continuous correction throughout the sequence of each track, thereby enhancing the predictive accuracy of the network. The quantity of hidden units within this layer directly correlates to the volume of information retained across timesteps, as illustrated in Figure 2. These hidden states are capable of encapsulating information from all preceding timesteps, independent of the sequence's length, ensuring a comprehensive temporal understanding. Subsequent to the LSTM layer, our architecture incorporates a batch normalization layer, which standardizes the LSTM output prior to its transition to the following Relu layer. This normalization significantly expedites the training process and fosters better convergence by mitigating internal covariate shift, as supported by [20]. Following this, a Relu layer is employed to apply a non-linear threshold operation, setting any input value below zero to zero. During training, a dropout layer is introduced to randomly nullify input elements with a specified probability, thereby imposing a regularization effect and preventing overfitting, as detailed in [21]. The architecture culminates in a Fully Connected (FC) Layer, which amalgamates the localized insights garnered by preceding layers. The dimensionality of the final FC layer is

meticulously aligned with the number of response variables required by the output layer, as elucidated in [22].

Our model's loss function is predicated on the mean-squared-error (MSE) metric, calculated for each state value prediction. The MSE quantifies the average squared discrepancy between the predicted and actual target values, serving as a critical measure of our model's predictive fidelity.

$$MSE = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2, \quad (1)$$

where m is the size of the predicted state vector, y_i is the ground truth value (KITTI car tracks) and \hat{y}_i is the prediction of the network for training set sample i . The loss is evaluated for several sequences, each with numerous time stamps. For our Prediction Network the loss function is half the mean-square-error of the predictions added up for each time step in the training set, normalized by the total number of all time stamps in the used sequences T :

$$loss = \frac{1}{2T} \sum_{j=1}^T \sum_{i=1}^m (y_{ij} - \hat{y}_{ij})^2. \quad (2)$$

During training, the average loss is calculated using the observations in the mini-batch, so T equals the mini-batch length.

Within the context of the KITTI dataset, encompassing both cars and vans, our model attained a Root Mean Square Error of 0.029 for the positional prediction of all objects within the validation dataset. This performance metric translates to an average deviation of 0.38 meters for predictions pertaining to the X coordinate and 0.21 meters for those related to the Y coordinate. Using an inhouse implemented Kalman Filter carried by Daimler Truck Research Group following [1], [3], an RMSE of 0.066 was achieved on the same data set.

b) Data preprocessing: In the development of our model, ground truth data comprising vehicle tracks from the KITTI dataset was utilized. A GT track encapsulates the temporal evolution of an object, delineating its trajectory from the moment it enters until it exits the sensor's detection range. To enhance model generalization and foster convergence during training, we normalized the state values of tracks at time t (predictors) and at time $t + 1$ (targets) in accordance with the methodology outlined in [23]. This normalization process aimed to standardize the distribution of both predictors and targets to have a mean of zero and a unit variance. The mean value μ and standard deviation σ for each state variable were computed across all tracks, employing the subsequent equations:

$$\mu = \frac{1}{m} \sum_{i=1}^m S_i \quad \text{and} \quad \sigma = \sqrt{\frac{1}{m-1} \sum_{i=1}^m |S_i - \mu|^2} \quad (3)$$

where S_i the total number of all time stamps of all tracks and m is the number of states per track.

In our [approach](#), we incorporate pre-padding as delineated by Reddy et al. in [24], where the authors elucidate the effect of padding strategies on the performance of neural networks in sequence-oriented tasks. Their investigation reveals that

while both pre-padding and post-padding are viable options, the selection of padding technique plays a pivotal role in the model's efficiency. This is particularly salient for LSTM networks, where the contextual integrity of the sequence is paramount for optimal performance. To handle sequences of varying lengths in our LSTM network, we utilized the padding technique. In this process, shorter sequences are padded with special padding tokens so that all sequences within a batch have the same length. We used zeros as padding tokens, which were appended to the end of a sequence as needed. This allows for efficient batch processing and ensures consistent training of the network without the model performance being affected by the varying sequence lengths.

As highlighted by [24], while padding introduces noise into the network, it is essential for aligning sequences within each mini-batch to facilitate the training of LSTM networks. To mitigate this noise, we sorted the sequences in the training dataset by their length prior to applying mini-batch padding. As demonstrated in Figure 3, this approach significantly reduces the padding (depicted in turquoise) required for each mini-batch by utilizing a pre-sorted dataset. This adjustment was critical for achieving convergence during the training process.

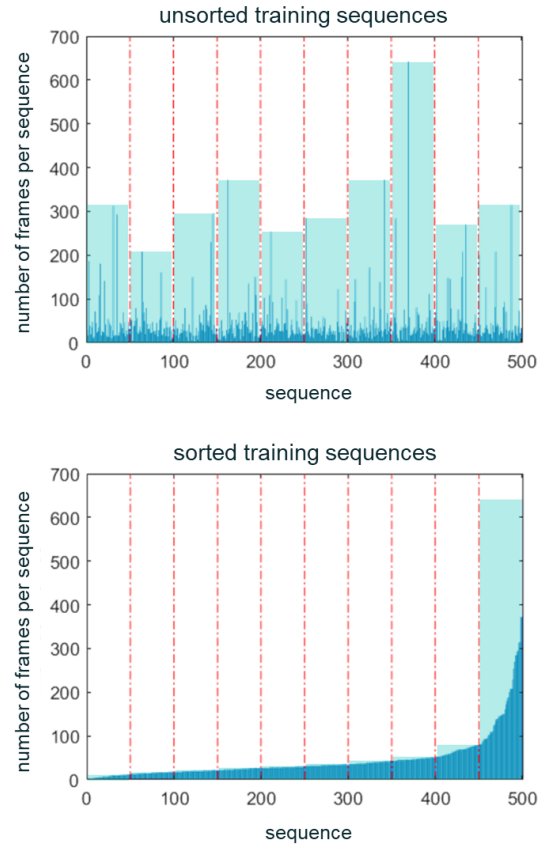


Fig. 3. Analysis of Sequence Padding: Unsorted vs. Sorted Data. Delineates the differential impact of padding on LSTM network training, contingent upon whether the input data sequences are sorted by length. In the upper panel, unsorted data necessitates extensive padding to equalize the sequence lengths within a batch, thereby augmenting the computational overhead. Conversely, the lower panel illustrates that sorting data by sequence length prior to batching significantly curtails the requisite padding.

B. Single Association Network (SANT)

The association network facilitates a data-driven methodology to address the combinatorial optimization challenge of data association, which is known to be NP-hard. Unlike traditional association methods referenced in [4], [5], our SANT model innovates by not relying on a predefined distance matrix as input. Instead, it processes the extant tracks alongside each newly measured sensor object directly.

This approach obviates the necessity for a predefined distance metric, thereby endowing the network with the autonomy to devise its own association logic, which it learns from the training data. Consequently, the association network supplants the traditional computation of a distance metric and the implementation of an association algorithm, such as the Hungarian Algorithm, with a data-driven, learning-based methodology.

a) *Data preprocessing*: In the formulation of SANT, we conceptualized data association as a temporally structured challenge, adopting a sequence-to-vector paradigm. This framework posits the association of a singular sensor object, denoted as $Z_{(t,1)}$, with a collection of tracks, represented as $X_{(t,1:n)}$. These tracks were meticulously curated from the KITTI dataset, which comprises annotated camera recordings. It is imperative to note, however, that genuine ground truth data for the specific association problem at hand are not available. To guarantee the unambiguity of assigning a sensor object to a single track within a specified set of tracks, each sensor object was synthetically generated at a given timestep from the pre-existing track set of that timestep. Furthermore, to simulate realistic sensor data from the GT data, artificial noise was introduced. This process was meticulously designed to reflect the inherent inaccuracies and uncertainties present in real-world sensor measurements.

To achieve this, a maximum of 3% of the value of the current state vector was randomly subtracted or added to each value. The data set was created in 7 iterations, so that the noise intensity was increased by 0.5% per iteration.

As shown in figure 4 the data format was created accordingly to enable index-based track assignment for SANT. The size of the input matrix therefore corresponds to $m \times n + 1$. With $m = 5$ as the number of state values for our work and $n = 16$ as the maximum number of tracks per time step.

The actual number of tracks can vary between 0 and a maximum of 16 objects in relation to the KITTI data set and the selected objects (cars and vans).

b) *Network architecture*: The network as depicted in fig. 4 is designed as a sequence-to-classification network. At each time step, a matrix is passed as input holding both, the information of to-be-assigned sensor object and the currently tracked objects.

Architectures with RNN, GNU, LSTM and Bidirectional Long Short-Term Memory Network (BiLSTM) layers were tested. As part of these investigations, the best association performance was achieved with the BiLSTM layer. The network shown in Figure 4 achieved the best performance in comparison with a Training Accuracy of 95% and Validation Accuracy of 95%. By combining the outputs of two LSTM layers that pass the information in opposite directions, [25] demonstrates the ability to capture the context from both ends

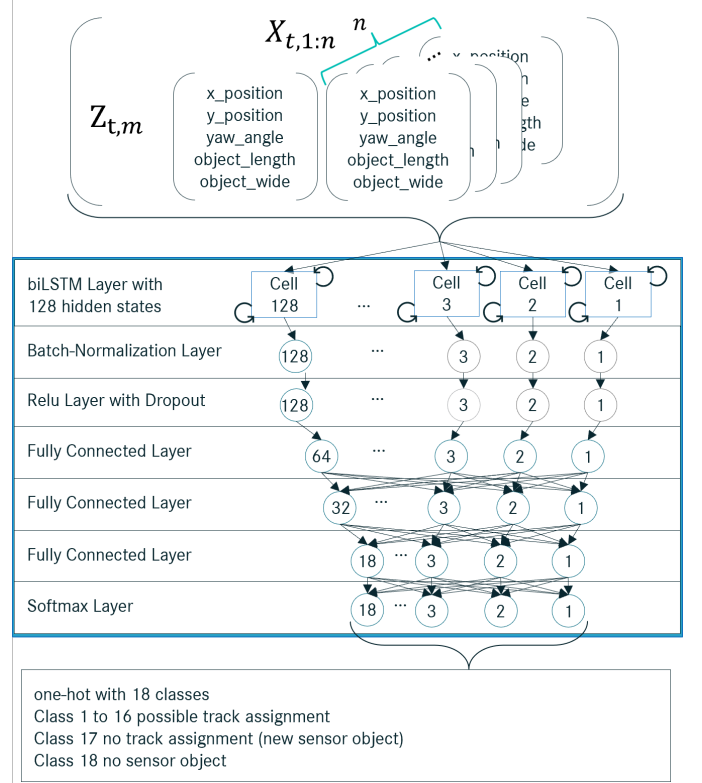


Fig. 4. Schematic representation of the generic network structure of the Single Association Network.

of the sequence. The resulting architecture is called BiLSTM. The output mode has been configured in the BiLSTM layer, so that the layer is able to receive a sequence as input and output value vector. This form of dimension reduction is necessary in order to carry out a corresponding classification. The last FC layer specifies the number of classes via the number of output values. The classes are calculated in the softmax layer by applying the softmax function resulting in a probability distribution. The softmax function converts a number of values z_i into a probability vector with i values. The cross-entropy cost function is utilized to quantify the discrepancy between the network's probabilistic predictions and the ground truth values, a method particularly suited for tasks involving categorically exclusive classes. This approach employs one-hot encoding to transform class representations into binary vector formats, thus enabling the delineation of each class within a 1-to- n coding scheme. The cross-entropy loss for each prediction, relative to its corresponding target value, is computed as follows:

$$loss = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C T_{ij} \log(Y_{ij}), \quad (4)$$

where N denotes the total number of samples, C represents the number of class categories, T_{ij} is the GT indicator for whether class j is the correct classification for sample i , and Y_{ij} is the predicted probability that sample i belongs to class j , as derived from the softmax function output.

C. Association network for multiple sensor objects (MANTa)

The developed single association network shows that a data-based association logic can be learnt from a deep learning model. The aim was also to develop a network to recognise a number of m sensor objects SO to an existing number of n tracks in one operation step. A multi-association network was developed, which is able to solve the following association problems:

- **1 to n** - one SO to n tracks
- **m to 1** - m SO to one track
- **m to n** - m SO to n tracks
- **m to 0** - m SO to no tracks
- **0 to n** - no SO to n tracks
- **0 to 0** - no SO to no tracks

With the integration of MANTa into a Multi-Object Tracking framework, the question can be asked whether a 0 to n and 0 to 0 assignment is a task to be solved. If no new SO is detected,

the prediction of the last operation step can be continued until the track is deleted on the basis of the decreasing probability of existence within the track management module. The association algorithm or the MANTa does not need to be called if no tracks and sensor objects are detected. Although these assignment options can therefore be resolved via the programme structure in the MOT framework, these options are also taken into account. This is intended to ensure that the network also learns to deal with SOs and tracks that are no longer available.

a) **Data preprocessing:** The data set for the training and validation of MANTa was created according to the described objective. Figure 5 shows the input data structure with corresponding association tasks in the displayed time step (85) of sequence 20 of the KITTI data set. The respective tracks per time step were extracted across all sequences. The extracted tracks were each modified with noise as in the SANT development and then normalised.

| | | | | | | | | | | | | | | | | | | | | |
|-------------------------------------|------|---------|---------|---------|---------|---------|---------|---------|---------|---|---|---|---|---|---|---|---|---|-----|------------|
| tracks | posX | -0.7321 | -1.4924 | -1.1651 | -1.0206 | -0.5031 | -0.0960 | 0.5941 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | T_{max} |
| | posY | 0.2551 | -0.2959 | -0.6997 | -0.2835 | -1.1377 | -1.1094 | 0.5258 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | T_{max} |
| | Yaw | -0.7592 | -0.7885 | -0.7877 | -0.7889 | -0.7901 | -0.7879 | -0.6944 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | T_{max} |
| | dimX | -0.2770 | -0.8776 | -0.1167 | 0.5767 | -0.0208 | -0.0884 | -0.9578 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | T_{max} |
| | dimY | 0.7674 | -0.5407 | -0.0992 | -0.1049 | -0.4269 | 1.2493 | -0.6291 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | T_{max} |
| sensor objects | posX | -1.4924 | 0.5941 | -1.0206 | 1.9263 | -0.0960 | -0.5031 | -0.7321 | -1.1651 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | T_{max} |
| | posY | -0.2959 | 0.5258 | -0.2835 | 1.2180 | -1.1094 | -1.1377 | 0.2551 | -0.6997 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | T_{max} |
| | Yaw | -0.7885 | -0.6944 | -0.7889 | -0.5229 | -0.7879 | -0.7901 | -0.7592 | -0.7877 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | T_{max} |
| | dimX | -0.8776 | -0.9578 | 0.5767 | -0.4118 | -0.0884 | -0.0208 | -0.2770 | -0.1167 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | T_{max} |
| | dimY | -0.5407 | -0.6291 | -0.1049 | 0.0500 | 1.2493 | -0.4269 | 0.7674 | -0.0992 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | T_{max} |
| one-hot vector 1:18 (from 1:288) | | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | OH_{max} |

Fig. 5. MANTa, data structure, shows the non-noisy sensor objects to enable a visual assignment and increase understanding of the association procedure.

Figure 5 shows the non-noisy sensor objects to enable a visual assignment and increase understanding of the association procedure. The values obtained in this way were assigned as shown in Fig. 5 as measurements in pseudo-random order per time step assigned to an $F * T_{max}$ matrix. Where F represents the number of features. The number of features results from the status values per track and SO set. (here, $F = 2 * m$, with $m = 5$ (number of state values in this investigation)). T_{max} stands for the maximum number of existing tracks per time step. For the defined test case with cars and vans, the KITTI data set results in $T_{max} = 16$. There are 7 tracks in the time step of the sequence shown. Each track receives a new measurement in this time step, and an additional object was detected (new sensor object).

The GT assignment is displayed at the bottom of the section. The size of $OH_{max} = 288$ results from the maximum number of tracks $T_{max} = 16$ and the number of possible assignment classes = 18. The assignment classes result from the described index class 1 to 16 and additional degrees of freedom. One degree of freedom of the assignment represents the case that no measurement exists, another that the measurement should not be assigned. The section of the one-hot vector shown (1:18) thus shows the GT assignment of the first sensor object to the track at position two.

b) **Network development:** The schematic representation 6 shows the developed network architectures for the simultaneous association of a large number of sensor objects to a large number of tracks. This is what we call a Multi-Association Network.

The BILSTM layer processes the input data as already explained for SANT. The task of associating a sensor object list with a track list requires a separate network part for each track. This extension is labelled accordingly in the graphic 6. For each track (from 1 to T_{max}), the MANTa has been developed with the familiar Fully Connected, Softmax stack. Each softmax output consists of a vector with classes = 18 elements, which represents the most probable assignment. This means that a single assignment can be realised for each track. The vectors 1 to T_{max} are linked together in the Concatenation Layer. This creates a vector with 288 elements, whereby 18 elements each represent the most probable assignment of a measurement to a track.

The cost function was implemented according to the format of the GT and output data of the network. The cross-entropy loss function already introduced was used and a clear assignment was realised by means of additional iteration per time step through the respective track blocks. The following formula is used to calculate the cross-entropy loss values for each input value Y_i and associated target value T_i element by

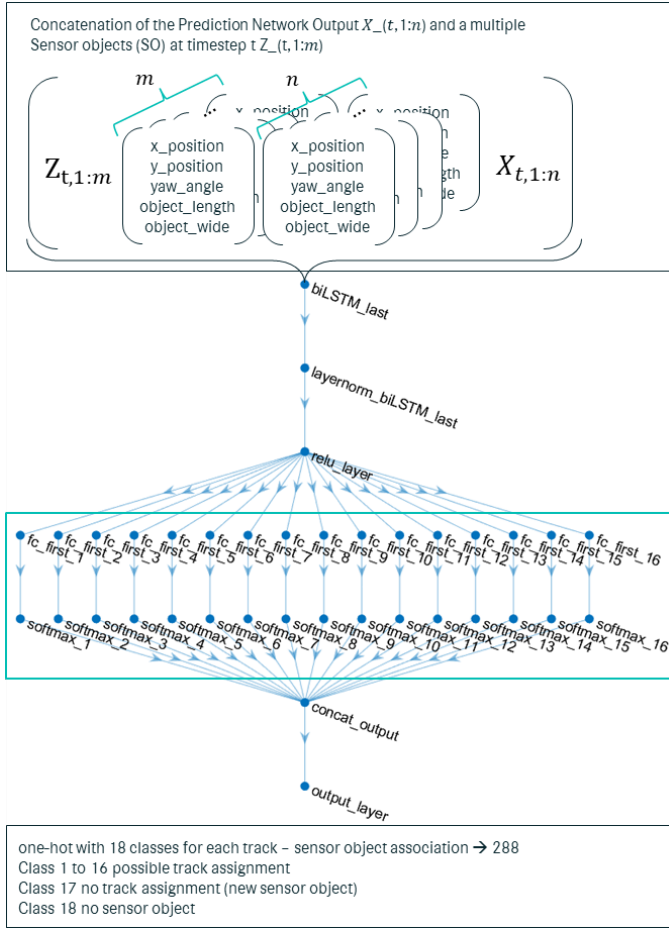


Fig. 6. Schematic representation of the generic network structure of MANTa.

element. To obtain a scalar per time step $loss_k$, all loss values are summed up and divided by the number of *classes*. This results in the average loss per time step for all tracks.

$$loss_k = \frac{1}{classes} \sum_{i=1}^{classes} -(T_i \ln Y_i + (1 - T_i) \ln(1 - Y_i)) \quad (5)$$

Then, all scalars obtained per time step are summarised and divided by the number of samples N of a minibatch:

$$loss = \frac{1}{N} \sum loss_k \quad (6)$$

With the described procedure, the Multi-Association Network could be trained, for assigning a list of sensor objects SO to a list of tracks in one operation step.

V. EXPERIMENTAL EVALUATION

In relation to the entire KITTI data set (Cars and Vans objects), MANTa achieves an average allocation accuracy of 70%. The main reason for this limitation was identified by analysing the extracted data. Figure 7 shows a distribution of the number of existing tracks per time step. For example, it can be seen that time steps containing a track account for almost a third of the entire available data set with 29.9% and an absolute number of 2315 samples.

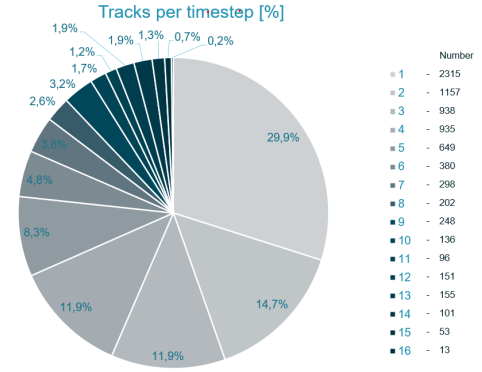


Fig. 7. KITTI Cars and Vans diagram, distribution of available data, number of tracks per sample

Time steps containing one to six tracks together make up 81.5% of the samples. Accordingly, tests were carried out with a reduced data set in order to demonstrate the multi-association capability of the network. MANTa correctly assigns 95% of the data set with time steps containing one to six tracks. This confirms MANTa's ability to assign data with the appropriate data set, because SANT also achieves a validation accuracy of approx. 95% in relation to the entire KITTI data set (cars and vans).

VI. CONCLUSION

The developed networks can be modularly integrated into the TbD framework and thus replace classic heuristic algorithms within the MOT process.

SPENT replaces the state predictions of the KF. The trained network estimates the predictions per time step without the need for a dynamics model. The implementation allows the recurrent network to update the internal hidden states per time step, thus achieving an accurate state prediction without an external correction.

The model is suitable for use in real time applications and represents an alternative approach to classical prediction methods. The network verification shows an RMSE of 0.026 averaged over the training, validation and test data set. In relation to the predictions of the positions of an object in the X coordinate, this corresponds to an average deviation of 0.42m. In relation to the position in the Y direction relative to the ego vehicle, the average deviation of the verification performed is 0.23m.

For another main task of the TbD MOT method, data association, SANT was developed as a replacement for the classic GNN method. This means that SANT can replace the algorithms for calculating a distance metric and assignment procedures such as HA with the learned, training data-based assignment logic. Based on a defined validation data set with approx. 2700 samples, SANT achieves an accuracy of 95%.

MANTa is a further development of SANT and addresses the limitation of individual assignment. A network extension could be implemented that assigns a set of sensor objects m to a set of tracks n . The data situation was analysed in more detail and identified as a limitation for the network performance. The verification carried out shows that MANTa

achieves an assignment accuracy of 95% in relation to the six most frequently occurring association sets.

ACKNOWLEDGMENTS

We would like to express our sincere thanks to the organisations that provided financial support for this research project. In particular, we would like to thank Daimler Truck AG and the DHBW Stuttgart. Without this support, our work would not have been possible.

Our special thanks also go to our academic supervisors and mentors. I would like to thank Prof Dr Matthias Drüppel for his tireless support, valuable advice and continuous guidance throughout the duration of this project. His expertise and commitment have contributed significantly to the success of our research.

Further thanks go to my colleague. I would like to thank Christian Bader for his essential help in conducting the experiments and providing technical resources. His support and fruitful discussions have greatly enriched this work.

Finally, we would like to thank our families and friends for their unconditional support and encouragement. Their patience, understanding and motivation have seen us through the challenges of this project and have been an invaluable help.

APPENDIX

PROOF OF THE ZONKLAR EQUATIONS

PROOF OF THE FIRST ZONKLAR EQUATION

Appendix goes here.

PROOF OF THE SECOND ZONKLAR EQUATION

And here.

REFERENCES

- [1] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, "Sort: Simple online and realtime tracking," in *2016 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2016, pp. 3464–3468.
- [2] Jenny Seidenschwarz, Guillem Brasó, Victor Serrano, Ismail Elezi, Laura Leal-Taixé, "Simple cues lead to a strong multi-object tracker," *Paper*, 2022.
- [3] J. Krejčí, O. Kost, O. Straka, and J. Duník, "Bounding box dynamics in visual tracking: Modeling and noise covariance estimation," *2023 26th International Conference on Information Fusion (FUSION)*, pp. 1–6, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:261126559>
- [4] Anton Milan, Seyed Rezatofighi, Anthony Dick, Ian Reid, Konrad Schindler, "Online multi-target tracking using recurrent neural networks," *Paper*, 2016.
- [5] H. L. H. Z. C. Mertz, "Deepda: Lstm-based deep data association network for multi-targets tracking in clutter," *Paper*, 2019.
- [6] Q. C. W. O. H. L. X. W. B. L. N. Yu, "Online multi-object tracking using cnn-based single object tracker with spatial-temporal attention mechanism," *ICCV*, 2017.
- [7] T. M. A. K. L. L.-T. C. Feichtenhofer, "Trackformer: Multi-object tracking with transformers," *Paper*, 2022.
- [8] W.-C. H. H. K. T.-Y. L. Y. C. R. Yu, "Soda: Multi-object tracking with soft data association," *Paper*, 2020.
- [9] H. W. Kuhn, "The hungarian method for the assignment problem," *Naval Research Logistics Quarterly*, 1955.
- [10] Ba-Tuong Vo, "code set for research use: Multi-sensor multi-target tracking," *Code*, 2013. [Online]. Available: <https://ba-tuong.vo-au.com/codes.html>
- [11] B. Ristic, S. Arulampalam, and N. Gordon, *Particle filters for tracking applications*. Artech House, 2004.
- [12] S. J. Julier and J. K. Uhlmann, "The unscented kalman filter for nonlinear estimation," *Proceedings of the IEEE*, vol. 92, no. 3, pp. 401–422, 2004.
- [13] E. A. Wan and R. Van Der Merwe, "The unscented kalman filter for nonlinear estimation," in *Proceedings of the IEEE Adaptive Systems for Signal Processing, Communications, and Control Symposium (AS-SPCC)*. IEEE, 2000, pp. 153–158.
- [14] Q. C. W. O. H. L. X. W. B. L. N. Yu, "Online multi-object tracking using cnn-based single object tracker with spatial-temporal attention mechanism," *Paper*, 2017.
- [15] N. Wojke, A. Bewley, and D. Paulus, "Deepsort: Simple online and realtime tracking with a deep association metric," in *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2017, pp. 3645–3649.
- [16] Roberto Henschel, Laura Leal-Taixé, Daniel Cremers, Bodo Rosenhahn, "Fusion of head and full-body detectors for multi-object tracking," *Paper*, 2017.
- [17] Johannes Fitz, "Datenassoziation für multi-objekt-verfolgung mittels deep learning," *Paper*, 2020.
- [18] Qi Chu, Wanli Ouyang, Hongsheng Li, Xiaogang Wang, Bin Liu, Nenghai Yu, "Online multi-object tracking using cnn-based single object tracker with spatial-temporal attention mechanism," *Paper*, 2017.
- [19] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *NIPS*, 2017.
- [20] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *ICML*, 2015.
- [21] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, 2014.
- [22] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," *AISTATS*, 2010.
- [23] I. G. Y. B. A. Courville, *Deep Learning*. MIT Press, 2019.
- [24] D. M. R. N. V. S. Reddy, "Effect of padding on lstms and cnns," *Paper*, 2019.
- [25] S. H. J. Schmidhuber, "Long short-term memory," *Neural Computation*, 1997.