

Pose and Semantic Map Based Probabilistic Forecast of Vulnerable Road Users' Trajectories

Viktor Kress, Fabian Jeske, Stefan Zernetsch, Konrad Doll, *Member, IEEE*, and Bernhard Sick, *Member, IEEE*

Abstract—In this article, an approach for probabilistic trajectory forecasting of vulnerable road users (VRUs) is presented, which considers past movements and the surrounding scene. Past movements are represented by 3D poses reflecting the posture and movements of individual body parts. The surrounding scene is modeled in the form of semantic maps showing, e.g., the course of streets, sidewalks, and the occurrence of obstacles. The forecasts are generated in grids discretizing the space and in the form of arbitrary discrete probability distributions. The distributions are evaluated in terms of their reliability, sharpness, and positional accuracy. We compare our method with an approach that provides forecasts in the form of Gaussian distributions and discuss the respective advantages and disadvantages. Thereby, we investigate the impact of using poses and semantic maps. With a technique called spatial label smoothing, our approach achieves reliable forecasts. Overall, the poses have a positive impact on the forecasts. The semantic maps offer the opportunity to adapt the probability distributions to the individual situation, although at the considered forecasted time horizon of 2.52 s they play a minor role compared to the past movements of the VRU. Our method is evaluated on a dataset recorded in inner-city traffic using a research vehicle. The dataset is made publicly available.

I. INTRODUCTION

A. Motivation

In the future, automated systems will operate in areas shared with humans and they must understand human behavior to make interactions safe, efficient, and comfortable. This is particularly important for automated vehicles and vulnerable road users (VRUs) in road traffic. A safe path planning of such vehicles requires a forecast of behavior and future trajectories of VRUs. However, future behavior is inherently fraught with uncertainty. Therefore, this work deals with probabilistic trajectory forecasting of VRUs, such as pedestrians and cyclists. The goal is the forecast of reliable probability distributions representing the movement capabilities of the VRUs fitted to the respective situation. The most important indicator for future behavior is the past movement. Besides the past trajectory, which is commonly used in the trajectory forecasting literature, we also include the body posture and movements of individual body parts. We address this by investigating the use of so-called 3D poses describing the three-dimensional positions of numerous joints along the body. Apart from past movements, the surrounding scenes of the VRUs is decisive for the future trajectory. It involves other road users, such as pedestrians, vehicles, static obstacles, lanes, and sidewalks.

V. Kress, F. Jeske, S. Zernetsch, and K. Doll are with the Faculty of Engineering, University of Applied Sciences Aschaffenburg, Aschaffenburg, Germany viktor.kress@th-ab.de, fabian.jeske@th-ab.de, stefan.zernetsch@th-ab.de, konrad.doll@th-ab.de

B. Sick is with the Intelligent Embedded Systems Lab, University of Kassel, Kassel, Germany bsick@uni-kassel.de

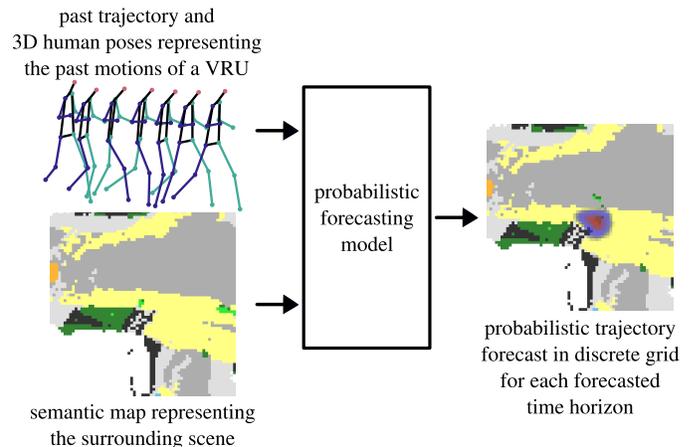


Fig. 1. Illustration of the approach for probabilistic trajectory forecasting. The model bases its forecasts on the VRU's past motions, represented by 3D poses, and semantic maps reflecting the surrounding scene. The colors correspond to the following semantic classes: *static obstacle* (black), *dynamic obstacle* (orange), *sidewalk* (yellow), *road* (dark grey), *walkable vegetation* (light blue), *person* (light green), *unknown obstacles* (dark green), and *unknown free space* (light grey). At the time of forecasting, the pedestrian is in the center of the map. In this scene, the pedestrian is walking on the sidewalk. For certain forecasted times, a discrete probability distribution (represented by the color gradient from red to blue on the right) is forecasted, expressing the likely future location of the VRU.

To encode such information, we use top-view semantic maps and consider them during forecasting. The overall approach is illustrated in Fig. 1. All the data used for trajectory forecasting is obtained solely using a stereo camera and LiDAR mounted on a research vehicle and is generated fully automatically. Such sensors could be incorporated into ordinary vehicles in the near future. In literature, approaches for forecasting continuous as well as discrete probability distributions, can be found. We compare methods of both types to explore their strengths and weaknesses. Given the importance of the forecasted probability distributions, e.g., for path planning of automated vehicles, we focus on evaluating the distributions in terms of their reliability and sharpness.

B. Related Work

This article mainly investigates probabilistic trajectory forecasting of VRUs while considering poses and surrounding scenes. By now, numerous works have been published regarding deterministic trajectory forecasting of VRUs in road traffic. For example, Keller and Gavrila [1] forecasted trajectories of pedestrians potentially crossing the street by using Gaussian process dynamical models (GPDMs). Goldhammer et al. [2] used polynomial approximations of past velocities

combined with a multi-layer perceptron (PolyMLP) to forecast the future trajectories of pedestrians and cyclists. However, forecasts are inevitably subject to uncertainty. An estimation of uncertainty is useful, for example, for safe path planning of automated vehicles. Without taking a negligible risk, an efficient traffic flow is not possible in areas shared by VRUs and automated vehicles. At the same time, path planning must consider the movement possibilities of VRUs and the risk taken must be quantifiable. Hence, methods for trajectory forecasting, including an estimation of the uncertainty, have been developed. For this purpose, some research forecasted several possible trajectories. Gupta et al. [3] used generative adversarial networks (GANs) and a pooling mechanism to incorporate dependencies among multiple people for forecasting multiple socially acceptable trajectories. By socially acceptable, the authors mean compliance with social norms, such as respecting a certain distance. A varying number of future trajectories were forecasted by memory augmented neural networks in [4]. However, such a set of forecasts only represents the uncertainty to a limited extent, and the quality of the uncertainty estimate is usually not evaluated. Instead, it is also feasible to forecast probability distributions that are either continuous or discretized regarding the spatial locations. Alahi et al. [5] focused on crowded spaces by forecasting bivariate Gaussian distributions for several pedestrians simultaneously via pooling based long short term memory (LSTM). The authors trained this model by minimizing the negative log-likelihood. Gaussian distributions describing the cyclists' future positions were forecasted in [6] using recurrent neural networks. In [7], forecasted trajectories were extended by an uncertainty estimation through an unconditional, constant model and used for path planning of autonomous vehicles. The authors of [8] proposed a method for forecasting all possible future positions of pedestrians in a set-based fashion using reachability analysis, contextual information, and traffic rules. The uncertainties of trajectory forecasts of cyclists were estimated in [9] in the form of unimodal Gaussian distributions with the help of a multi-layer perceptron (MLP). In addition, an approach to evaluate the reliability of the forecasted uncertainties was presented. Overall, this approach was able to make reliable forecasts for the motion types *start*, *stop*, *turn left*, and *turn right*, while the forecasts for the motion types *move* and *wait* were underconfident. Bieshaar et al. [10] extended single-output quantile regression to multivariate-targets for trajectory forecasting of cyclists. These so-called quantile surfaces represent star-shaped distributions using discrete quantile levels.

Besides forecasts of continuous distributions, distributions in discretized space are feasible as well. For example, Markov chains were used in [11] to forecast probabilistic distributions as occupancy grids considering the pedestrians' potential goals. Jain et al. [12] encoded the past in the form of multi-channel images and forecasts trajectories over long time horizons using a discrete residual flow network in the form of probability grids for each forecasted time step.

In addition to the past trajectory of the respective VRU, the posture and movements of individual body parts can be represented by 3D body poses and considered in trajectory forecasting. Quintero et al. [13] performed trajectory forecast-

ing of pedestrians based on 3D poses with several balanced GPDMs. They were trained on 3D poses for different motion types, while the most similar model for the individual pedestrian behavior was chosen for trajectory forecasting. In [14], 3D poses estimated from a moving vehicle were used for trajectory forecasting of pedestrians and cyclists. While both approaches achieved improvements by using poses, they did not include an estimate of the uncertainty.

The surrounding scenes, e.g., obstacles, lanes, sidewalks, or vegetation, influence the future trajectories of VRUs. Accordingly, numerous methods considered the surrounding for trajectory forecasting: Xue et al. [15] used an LSTM architecture for trajectory forecasting of pedestrians with an upstream convolutional neural network (CNN) extracting relevant features from top-view images of the scene. In [16], the surrounding was expressed through motion heat maps representing the prior movements of other VRUs, segmented maps describing the accessible areas, and aerial photography images. This representation is processed by a CNN and, together with the past trajectory and data about nearby road users, used to forecast several possible trajectories by means of a neural network. Ridel et al. [17] performed a semantic segmentation on top-view images using residual networks (ResNet). They encoded the past trajectory in binary 2D grids and forecasted discrete probability distributions in agent-centric grids for future time steps using a network architecture consisting of CNNs and convolutional LSTMs (ConvLSTM). However, these approaches require current images from infrastructure-based cameras or drones, which is usually not feasible in road traffic. In [12], images encoding semantic maps were used for trajectory forecasting of pedestrians. The images mask information of maps, e.g., crosswalks, drivable surfaces, lanes, and traffic light states, which were annotated semi-automatically, and the history of dynamic objects detected from LiDAR and camera. The maps were centered on the respective pedestrian and processed by a CNN. Marchetti et al. [4] projected semantic labels of static objects aggregated over time into a top-view map using a LiDAR point cloud and IMU data from a moving vehicle. Those maps were subsequently used to refine the forecasts to ensure compatibility with the surrounding.

C. Main Contributions and Outline of this Article

Our main contributions are the following: First, we combine the past trajectory with 3D poses of pedestrians and cyclists and maps of the surrounding scenes for discrete probabilistic trajectory forecasting by means of neural networks. The 3D poses represent the past motions of different body parts in great detail. To the best of our knowledge, this is the first time poses are used for probabilistic trajectory forecasting. We use semantic maps reflecting both the static (e.g., sidewalks or obstacles) and dynamic surrounding (e.g., vehicles or other VRUs) to align the forecasts with the surroundings. We individually examine the impact of the past trajectory, 3D poses, and maps on trajectory forecasting.

Second, we analyze the forecasted distributions in detail with respect to their reliability, sharpness, and positional

accuracy. We propose a method to calibrate the discrete probabilistic trajectory forecasting models regarding reliability and investigate the effects of the resolution of the discrete probability distributions on the quality of the forecasts, which has not been studied before in the literature.

Third, we evaluate the advantages and disadvantages of forecasting discrete versus continuous probability distributions. For this purpose, we use the approach from [9] as a continuous comparison method. Since this approach is originally based solely on the past trajectory, we extend it to include 3D poses as well.

Fourth, all used data were generated fully automatically with the help of a single vehicle’s sensors in real road traffic. In contrast to the literature, neither drones, infrastructure-based sensors, nor manual annotations were used. Accordingly, the approach is appropriate for an application in automated vehicles. The dataset is made publicly available [18].

The remainder of this article is organized as follows: Sec. II introduces the dataset, while in Sec. III the discrete (Sec. III-A) and continuous trajectory forecasting methods (Sec. III-B) are explained. This sections focuses in particular on the network architecture, the training methodology, and evaluation metrics. Next, the results are presented in Sec. IV. It covers the process of hyperparameter optimization, reliability calibration using spatial label smoothing, the impact of poses, semantic maps and the cell size, and a comparison of the discrete and continuous forecasting approaches. Finally, we conclude with a summary and an outlook on future work (Sec. V).

II. DATASET

The dataset was recorded with a moving vehicle in inner-city traffic in Aschaffenburg, Germany. The recordings were taken at different times of the day and year over four years and under different weather conditions to reflect road traffic variability. The driven routes cover different types of roads, such as multi-lane roads, roads with bike lanes, traffic-lighted and non-traffic-lighted intersections, traffic-calmed areas, crosswalks, bus stops, and more. Some of the recorded pedestrians and cyclists were instructed to walk or ride specific routes. The remaining VRUs were uninvolved and uninstructed individuals of all ages and agilities. The vehicle was equipped with a stereo camera behind the windshield, a LiDAR in the vehicle’s front, and a vehicle localization system. The trajectories of all VRUs up to a distance of 25 m were recorded using the stereo camera and a Kalman tracker. To obtain the 3D poses of the VRUs, first, the 2D poses, i.e., the two-dimensional coordinates of several joints in the images, were estimated using the CNN proposed in [19] followed by a reconstruction of plausible 3D poses using the approach from [20]. More details on this procedure and an evaluation of its accuracy can be found in [21]. Compared to 2D poses or other image-based methods, 3D poses have the advantage of being independent of the perspective of the recording camera and allow for a compensation of the vehicle’s own motion.

Further, semantic maps were generated describing the particular scene. For that purpose, semantic segmentation of the stereo camera images was performed using the approach

from [22]. To create a top-view semantic map, the semantic information about the surrounding was fused with an occupancy map based on the point cloud of the LiDAR. In this way, the high positional accuracy of the LiDAR is combined with the density of the camera images. The maps contain $n_s = 8$ semantic classes: *Static obstacles* include all types of non-movable barriers that are not negotiable by pedestrians and cyclists. Among these are, e.g., buildings, walls, fences, traffic signs, or vegetation. *Dynamic obstacles*, on the other hand, comprise all types of moving or parked vehicles. The other classes are *sidewalk*, *road*, *walkable vegetation*, such as meadows, and *person* including pedestrians and cyclists. The remaining two classes are *unknown obstacles* and *unknown free space*. They cover areas for which the semantic class could not be determined. By aggregating all data in the dataset, a map of the driven roadways was created, containing only the invariant semantic classes. This initial map is subsequently enhanced at each point in time by the measurements of the immediate past resulting in an up-to-date map including moving elements. In a real-world application, the initial map could be replaced by commercially available high precision maps, and the current maps could be shared between multiple vehicles. The maps have a spatial resolution of $0.35 \times 0.35 \text{ m}^2$ and a temporal frequency of 25 fps. Each pixel is classified as one of the eight classes.

The dataset contains 2351 trajectories of pedestrians and 1075 trajectories of cyclists with corresponding 3D poses and semantic maps. Each point in time of the trajectories was manually annotated with a motion type. The motion types comprise the states *wait*, *start*, *move*, and *stop* for pedestrians and cyclists. Additionally, a distinction is made between the motion types *turn left* and *turn right* for cyclists. These two motion types do not exist for pedestrians because they are difficult to annotate even for humans due to the agility of pedestrians. These annotated motion types are used solely for a differentiated evaluation in Sec. IV-C. About 60 % each of the pedestrian and cyclist data is used as the training set, 20 % as validation set, and the remaining 20 % as the test set. It was ensured that the same location does not occur in different sets and that the motion types are distributed as equally as possible. This is important to guarantee the generalization ability to arbitrary locations and for a fair comparison with methods without knowledge of the surrounding. To ensure rotational invariance, the trajectories, poses, and semantic maps are randomly rotated and augmented by a factor of 3 through repeated random rotation.

The dataset including trajectories, 3D poses, semantic maps, and motion types has been made publicly available [18]. All the data were recorded in accordance with the guidelines of the University of Applied Sciences Aschaffenburg and German privacy laws.

III. METHOD

A. Discrete Probabilistic Trajectory Forecast

This work aims at forecasting probability distributions describing the possible future locations of the VRUs discretized in time and space. Therefore, we forecast the distribution \hat{p}_{t_f} for the forecasted time horizons $t_f \in T_f$ for

which $\sum_{\vec{g} \in G} \hat{p}_{t_f}(\vec{g}) = 1$. Here, $G \subset \mathbb{R}^2$ is a grid centered on the current position of the respective VRU with cell size e^2 discretizing the space and $\hat{p}_{t_f}(\vec{g})$ denotes the forecasted probability for a cell with center point \vec{g} .

1) *Input Features*: We train and analyze three models based on different feature sets for trajectory forecasting: The model based on the past trajectory, referenced hereafter by the abbreviation **d_t** (**d**iscrete forecasting model based on past **t**rajectory), uses the VRU's two-dimensional head position of the ground plane for each time t_i in the observation period T_i . Accordingly, the feature set is given by $f_{d_t} = \{[x_{Head,t_i}, y_{Head,t_i}] \mid t_i \in T_i\}$. The origin of the corresponding coordinate system is the head position at the current time t_c .

The second model **d_tp** (**d**iscrete forecasting model based on past **t**rajectory and **p**oses) additionally uses 3D poses as input. Instead of the head position, the feature vector for time t_i consists of the three-dimensional position of 13 joints along the body. A sequence of 3D poses with associated joint positions is shown in Fig. 1. The coordinate system remains the same. All 3D poses are scaled in 3D space to have the same width of the hips. Again, the final feature set is obtained by concatenating the observation period: $f_{d_{tp}} = \{[x_{Head,t_i}, y_{Head,t_i}, z_{Head,t_i}, \dots, z_{LFoot,t_i}] \mid t_i \in T_i\}$. The feature sets of both models are normalized over all training samples using the statistical z-transformation.

Finally, the third model **d_tpm** additionally uses the semantic map m_{t_c} at current time t_c to adapt the trajectory forecasts to the particular scene. The origin of the semantic map is the current head position of the respective VRU. The semantic maps have the same dimensions and cell size as the forecasted grid G . They are converted into multichannel images. Each channel represents a binary image encoding the presence or absence of one of the eight semantic classes mentioned above at each particular location.

2) *Network Architecture*: The neural network has a two-stream architecture (Fig. 2). The trajectory net is responsible for processing the past trajectory and the 3D poses. Accordingly, depending on the model, the feature sets f_{d_t} or $f_{d_{tp}}$ serve as input for this stream. It consists of several fully connected layers, each followed by a Rectified Linear Unit (ReLU) activation function. The number of layers and the number of neurons per layer are hyperparameters chosen in an optimization step in Sec. IV-A. As a result, this network segment yields a feature map of size h^2 with $|T_f|$ channels. The dimensions correspond to those of the final network output \hat{p} . The second stream, called semantic map net, processing the semantic map consists of convolutional layers, each with a 3×3 filter and ReLU activation. This stream is omitted for the models **d_t** and **d_tp**, as they do not consider the semantic map. The two streams' feature maps are concatenated along the channel dimensions and processed in the so-called fusion net by further convolutional layers, again with 3×3 filter and ReLU activation. The number of convolutional layers and the number of filters of each convolutional layer are additional hyperparameters. Other activation functions, filter sizes, and the addition of pooling layers did not improve the forecasts. The last layer is a linear convolution with 1×1 filter producing

a grid for each forecasted time horizon. Finally, a softmax activation function is applied for each grid to obtain the final probability distributions. We have also tested other network architectures, such as UNets [23] or ConvLSTMs [24], achieving similar results.

3) *Training*: We train the models **d_t** and **d_tp** by minimizing the cross entropy (Eq. 1) between the forecasted distribution $\hat{p}_{t_f,i}$ and target distribution $p_{t_f,i}$ for all N training samples and forecasted time horizons T_f .

$$L = -\frac{1}{|T_f|} \sum_{t_f \in T_f} \frac{1}{N} \sum_{i=1}^N \sum_{\vec{g} \in G} p_{t_f,i}(\vec{g}) \log(\hat{p}_{t_f,i}(\vec{g})) \quad (1)$$

The cross entropy with one-hot encoded target tends to produce overconfident, i.e., to narrow, probability distributions that are not reliable. Typically, in order to avoid overconfidence, label smoothing with uniform distribution over the classes is applied. However, in our case, the grid cells are spatially related. Therefore, we use spatial label smoothing: Instead of a uniform distribution we use a Gaussian distribution $p_{t_f,i}(\vec{g}) = \mathcal{N}_{\vec{y}_{t_f,i}, \sigma_{t_f}}(\vec{g})$ with the actual grid position $\vec{y}_{t_f,i}$ as expected value and standard deviation σ_{t_f} as target distribution. The standard deviations are additional hyperparameters. Their effects are examined in Sec. IV-B.

For model **d_tpm** the loss function (Eq. 1) is used to consider obstacles in the surroundings. For this purpose we define obstacles $c_{t_f,i}(\vec{g})$ for sample i equaling 1 if there is an obstacle at the cell with center point \vec{g} and 0 otherwise. In this context, we declare obstacles as objects of the semantic classes *static obstacles* and *dynamic obstacles*. Thus, $c_{t_f,i}$ can be derived from the semantic map $m_{t_f,i}$ at the respective time. Instead of the Gaussian distribution a modified Gaussian distribution $p_{t_f,i}(\vec{g}) = \mathcal{N}_{\vec{y}_{t_f,i}, \sigma_{t_f}}^{c_{t_f,i}}(\vec{g})$ is used as target equaling 0 at locations with obstacles. The probabilities are scaled such that $\sum_{\vec{g} \in G} p_{t_f,i}(\vec{g}) = 1$.

4) *Evaluation Method*: Several scores describing different properties are used to evaluate the estimated probability distributions. They are introduced and explained in this section.

a) *Reliability*: The reliability measures whether the variances of the distributions are estimated correctly. In [9], reliability is calculated by comparing confidence intervals of the forecasts with the observed frequency of actual positions within the interval. In the following, we use a similar approach adapted to the discrete probability distributions. In favor of clarity, we omit the index for the respective forecasted time horizon t_f below. The procedure is the same for all forecasted time horizons. We define a set of cell center points Φ for which the associated cells have forecasted probabilities $\hat{p}(\vec{g})$ greater than the forecasted probability at the cell of the actual position $\hat{p}(\vec{y})$ for each sample (Eq. 2). The respective confidence level C given the actual grid position \vec{y} and the forecast \hat{p} is obtained by adding the forecasted probabilities of all cells with the center points within the set (Eq. 3).

$$\Phi(\vec{y}, \hat{p}) = \{\vec{g} \in G \mid \hat{p}(\vec{g}) \geq \hat{p}(\vec{y})\} \quad (2)$$

$$C(\vec{y}, \hat{p}) = \sum_{\vec{g} \in \Phi(\vec{y}, \hat{p})} \hat{p}(\vec{g}) \quad (3)$$

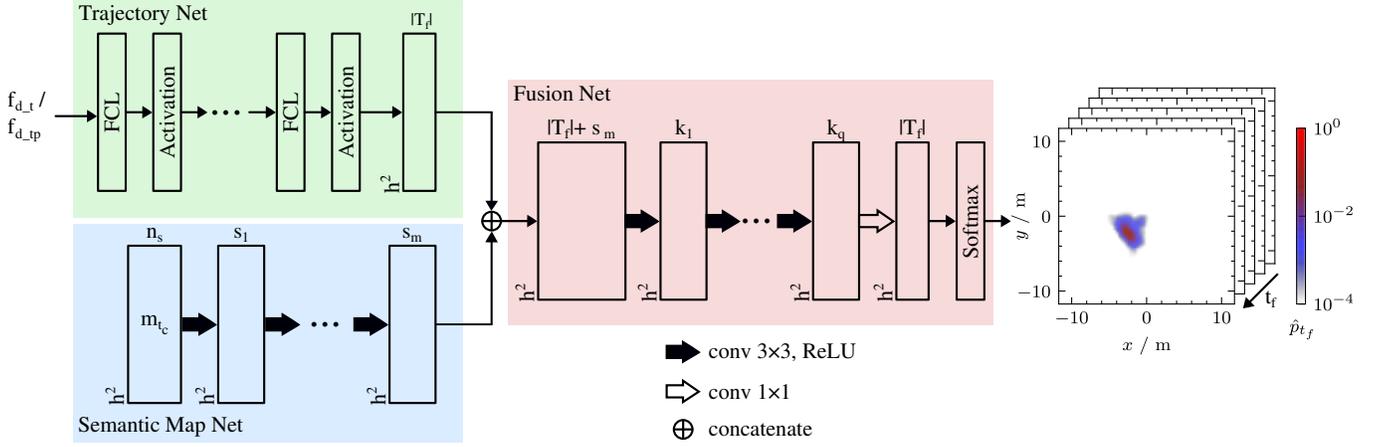


Fig. 2. Two stream network architecture with the so-called trajectory net processing the input trajectories and poses and the semantic map net receiving the semantic maps. Boxes represent the feature maps with their dimensions given in the lower-left corner and the number of channels on top of each box.

The observed frequency f_o in the dataset with M samples is calculated according to Eq. 4 for a given confidence level $1-\alpha$. A measure of reliability is obtained by comparing the observed frequencies to the given confidence levels. Ideally, they should be identical. A visualization is achieved by plotting the observed frequency given the confidence level, whereby the diagonal describes perfect reliability (see Fig. 3 for an example).

$$f_o(1-\alpha) = \frac{1}{M} \sum_{i=1}^M \begin{cases} 1, & \text{if } C(\vec{y}_i, \hat{p}_i) \leq 1-\alpha \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

For a numerical evaluation of the reliability, we use a score similar to the widely used Expected Calibration Error (ECE) [25] (Eq. 5). Here, the range of possible confidence levels is partitioned into B bins. The score is given by the averaged sum of the absolute difference between the observed frequency and given confidence level weighted by the number of forecasts j_b within the respective confidence level bin. It is averaged over all forecasted times T_f , to gain a score for the entire forecasts.

$$\text{ECE} = \frac{1}{|T_f| M} \sum_{t_f \in T_f} \sum_{b=1}^B j_{b,t_f} |(1-\alpha_b) - f_{o,t_f}(1-\alpha_b)| \quad (5)$$

This score is related to the definition of reliability obtained by a decomposition of the Brier score [26] and is used for model calibration in Sec. IV-B.

b) Sharpness: In order to evaluate the sharpness of the forecasted distributions, we determine the area covered by a specific confidence level. For this purpose we search for a threshold $\tau = \max(\xi)$ such that $\sum \{\hat{p}(\vec{g}) \mid \hat{p}(\vec{g}) \geq \xi\} \geq 1-\alpha$. The area of a given confidence level and sample is then given by:

$$A(1-\alpha) = |\{\vec{g} \mid \hat{p}(\vec{g}) \geq \tau\}|e^2 \quad (6)$$

The sharpness S (Eq. 7) is defined as the mean area \bar{A}_{t_f} of all samples of the forecasted time horizon t_f normalized

to the time horizon and averaged over all forecasted time horizons T_f .

$$S(1-\alpha) = \frac{1}{|T_f|} \sum_{t_f \in T_f} \frac{\bar{A}_{t_f}(1-\alpha)}{t_f} \quad (7)$$

c) Positional Accuracy: The weighted average Euclidean error (WAAEE) is used to evaluate the positional accuracy (Eq. 8). For a given forecasted time horizon, this metric determines the average sum of the M Euclidean distances between each cell center point \vec{g} in the grid and the actual grid position \vec{y}_{i,t_f} weighted according to the forecasted probability \hat{p}_{i,t_f} for the cell with that center position.

$$\text{WAAEE}_{t_f} = \frac{1}{M} \sum_{i=1}^M \sum_{\vec{g} \in G} \hat{p}_{i,t_f}(\vec{g}) \|\vec{g} - \vec{y}_{i,t_f}\|_2 \quad (8)$$

To gain a metric evaluating the forecasts over all forecasted time horizons, we calculate the average specific WAAEE (ASWAAEE). The metric normalizes and averages the WAAEE according to the forecasted time horizon (Eq. 9). It is the probabilistic equivalent of the average specific average Euclidean error introduced in [2].

$$\text{ASWAAEE} = \frac{1}{|T_f|} \sum_{t_f \in T_f} \frac{\text{WAAEE}_{t_f}}{t_f} \quad (9)$$

d) Forecasts of Obstacles Collisions: The impact of using the semantic maps for trajectory forecasting is difficult to quantify, as discussed in Sec. IV-C2. However, the proportion of the probability distribution located at areas occupied by static obstacles can be measured. This is obviously an incorrect forecast as long as the semantic maps are assumed to be correct. We define the so-called occupancy score O_{t_f} (Eq. 10) for a certain forecasted time horizon as the sum of forecasted probabilities at locations of static obstacles $o_{t_f,i}$. Here, $o_{t_f,i}$ is a binary map with value 1 at locations of static obstacles.

$$O_{t_f} = \frac{1}{M} \sum_{i=1}^M \sum_{\vec{g} \in G} \hat{p}_{t_f,i}(\vec{g}) o_{t_f,i}(\vec{g}) \quad (10)$$

The metric refers only to static obstacles and not to dynamic ones, since their movements are again a matter of probability and therefore collisions of the forecasts with dynamic obstacles are not necessarily wrong.

B. Continuous Probabilistic Trajectory Forecast

We compare our approach for probabilistic forecasting of trajectories in discrete space with the method from [9] for forecasting continuous probability distributions. This method estimates the forecasts' uncertainty in the form of Gaussian distributions whose parameters (expected values, variances, correlation coefficients) are forecasted by a neural network with feed-forward architecture and fully connected layers. As input, the approach originally uses the past head positions of cyclists in an ego coordinate system. The ego coordinate system has its origin in the current position of the respective VRU, and the x - and y -axes are defined depending on the movement direction. We train and optimize the neural network on our dataset separately for pedestrians and cyclists. This model is referred to as c_t (continuous forecasting model based on past trajectory). In addition, another model is created for pedestrians and cyclists, respectively, by extending the input feature space with poses (model c_tp). The poses are rotated according to the movement direction to obtain compatibility with the ego coordinate system. The feature set is defined by $f_{c_tp} = \{^{ego}[x_{Head,t_i}, y_{Head,t_i}, z_{Head,t_i}, \dots, z_{LFoot,t_i}] \mid t_i \in T_i\}$. We refer to these comparison methods hereafter as continuous models. The evaluation of the continuous forecasting models is done equivalent to the discrete models using reliability, sharpness, and ASWAE. However, the calculations are realized by sampling from the continuous probability distributions.

IV. EXPERIMENTAL RESULTS

A. Hyperparameter Optimization

We train the models d_t , d_tp , d_tpm , c_t , and c_tp separately for pedestrians and cyclists, resulting in a total of 10 models. For training the networks, the adaptive moment estimation (Adam) optimizer [27] is used as well as early stopping as regularization technique.

We forecast probability distributions for five time horizons $T_f = \{0.44\text{ s}, 0.96\text{ s}, 1.48\text{ s}, 2.00\text{ s}, 2.52\text{ s}\}$ using an observation period of 1 s $T_i = \{-0.96\text{ s}, \dots, -0.04\text{ s}, 0.00\text{ s}\}$ corresponding to 25 position measurements. The size of the forecasted grid is defined such that the positions of all pedestrians or cyclists within our dataset are located certainly within the grid for all time horizons. For pedestrians the size is $h^2 = 67 \times 67\text{ px}^2$ and for cyclists $147 \times 147\text{ px}^2$ corresponding to $23.45 \times 23.45\text{ m}^2$ and $51.45 \times 51.45\text{ m}^2$, respectively. The size of the grids could be chosen smaller for short time horizons. However, this was not done in favor of a simpler network architecture and implementation. The size of a cell equals $e^2 = 0.35 \times 0.35\text{ m}^2$, which corresponds to the resolution of the semantic maps and approximately to the area occupied by a human being. However, since the cell size is principally a hyperparameter, its effects are examined in Sec. IV-C3.

TABLE I

THE SELECTED HYPERPARAMETERS OF THE NETWORK ARCHITECTURE. NOTE THAT, IN ADDITION TO THE SPECIFICATIONS IN THE TABLE, THE FUSION NET ALWAYS CONCLUDES WITH A CONVOLUTIONAL LAYER WITH 1×1 FILTER.

parameter	d_t	d_tp	d_tpm
number of hidden layers in trajectory net	4	5	5
number of neurons in trajectory net	150	50	50
number of conv in semantic map net	-	-	1
number of filters in semantic map net	-	-	8
number of conv in fusion net	2	2	2
number of filters in fusion net	10	20	20

The hyperparameters are optimized by parameter sweeps using the validation dataset. Afterward, the networks obtaining the lowest ECE on the validation data are trained on a combined set of training and validation data and evaluated on the test dataset. All hyperparameters of the network architecture are optimized only for pedestrians and adopted for cyclists to keep the computational cost reasonable. The resulting parameters of the discrete models are summarized in Tab. I. The final networks of models c_t and c_tp each consist of 2 hidden layers with 100 neurons per layer.

B. Reliability Calibration

One focus of this work is the forecast of reliable probability distributions, as they have a safety-relevant meaning, e.g., for path planning of autonomous vehicles. However, the use of the cross entropy as loss function leads to overconfident distributions, which are critical for an application, e.g., path planning of automated vehicles. The left plot in Fig. 3 shows the reliability diagram for the model d_tp for pedestrians trained using the cross entropy evaluated on the validation dataset. The curves fall below the ideal diagonal for all forecasted time horizons, which characterizes an overconfident distribution. There are considerable differences in the reliability of the different forecasted time horizons: For increasing time horizons, the reliability approaches more and more the diagonal. Overall, an ECE of 8.6% is achieved. A popular approach to calibration is temperature scaling [28]. Therefore, we performed temperature scaling separately for each forecasted time horizon on the validation dataset. While the reliability indeed improves with an ECE of 6.4%, spatial label smoothing as described in Sec. III-A3 achieves better results on our dataset. Here, the ECE is reduced to 4.6%. The corresponding reliability diagram is shown on the right in Fig. 3. The standard deviation of the Gaussian distribution was optimized separately for each forecasted time horizon. In each case, the standard deviation achieving the smallest ECE on the validation dataset was chosen. For the forecasted time horizons T_f and cell size e^2 the standard deviations $\sigma = [0.48e, 0.48e, 0.53e, 0.55e, 0.55e]$ were found using model d_t for pedestrians and adopted for all discrete models. Due to the high computational costs required to optimize the standard deviations, the values were adopted for cyclists. An adjustment to the respective dataset could lead to further improvements. The procedure prevents overconfident distributions. However, the resulting reliability diagram shows

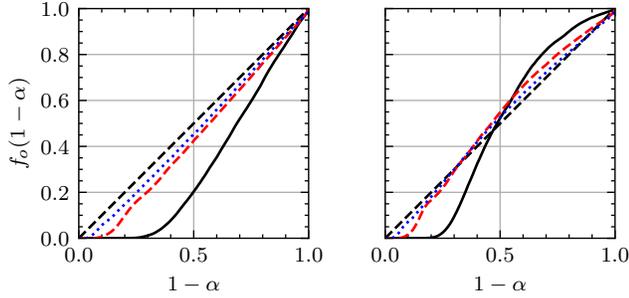


Fig. 3. Reliability diagrams for model d_tp for pedestrians using plain cross entropy as loss function (left) and spatial label smoothing (right) evaluated on the validation dataset. The reliability is shown for the forecasted time horizons 0.44 s (black solid line), 1.48 s (dashed red line) and 2.52 s (dotted blue line). The ideal diagonal is plotted as dashed black line.

TABLE II
ACHIEVED ECE IN % OF ALL MODELS FOR PEDESTRIANS STRUCTURED BY MOTION TYPES.

model	all	wait	start	move	stop
c_t	4.9	12.5	4.5	4.4	4.1
d_t	7.1	6.0	9.5	6.7	10.3
c_tp	3.4	8.2	3.3	3.2	4.2
d_tp	4.5	5.4	3.5	4.3	8.2
d_tpm	3.8	6.9	3.7	3.9	6.4

S-shaped curves, especially for short forecasted time horizons. This effect is caused by discretization and will be discussed in more detail in Sec. IV-C3. It should be noted here that reliable models can also be achieved using small network sizes, i.e., a small number of layers and neurons. But those networks are not able to make sharp forecasts at the same time because they are not able to learn the movement behavior precisely.

C. Comparison of the Various Methods

In the following, the results of the final network configuration of each of the five models for pedestrians and cyclists are presented. All results are measured on the separate test dataset. Quantitative results following the evaluation metrics defined in Sec. III-A4 are provided for all models in Tab. II to IV for pedestrians and for cyclists in Tab. V to VII. Here, the results are distinguished for the different motion types. We discuss the impact of using poses and semantic maps on the forecasts and investigate the influence of the cell size on the forecasting results. Last but not least, we compare the discrete and continuous trajectory forecasting models.

1) *Impact of Poses*: For pedestrians, with regard to positional accuracy, the use of poses results in a 9.7% reduction of the ASWAE for the continuous forecasting approach (motion type *all* for models c_t and c_tp in Tab. IV). For the discrete method, an improvement of 7.2% is achieved (motion type *all* for models d_t and d_tp in Tab. IV). Meanwhile, improvements are also observed in the reliability (Tab. II) and sharpness (Tab. III) for both the discrete (d_t and d_tp) and continuous (c_t and c_tp) approaches. Fig. 4 shows a comparison of discrete forecasts with and without the use of poses

TABLE III
SHARPNESS $S(0.95)$ IN $\frac{m^2}{s}$ OF ALL MODELS FOR PEDESTRIANS.

model	all	wait	start	move	stop
c_t	2.89	2.65	3.43	2.74	3.23
d_t	3.11	3.41	3.41	2.86	3.30
c_tp	2.35	2.28	2.67	2.24	2.52
d_tp	3.03	2.98	3.32	2.93	3.17
d_tpm	3.34	3.72	3.68	3.06	3.49

TABLE IV
ASWAE IN $\frac{m}{s}$ OF ALL MODELS FOR PEDESTRIANS.

model	all	wait	start	move	stop
c_t	0.57	0.51	0.65	0.56	0.61
d_t	0.60	0.54	0.67	0.60	0.64
c_tp	0.51	0.50	0.56	0.50	0.55
d_tp	0.56	0.53	0.59	0.55	0.60
d_tpm	0.57	0.55	0.62	0.55	0.61

TABLE V
ECE IN % OF ALL MODELS FOR CYCLISTS STRUCTURED BY MOTION TYPES.

model	all	wait	start	move	stop	turn left	turn right
c_t	6.83	15.34	12.71	6.87	4.53	12.26	20.68
d_t	5.90	13.06	10.63	5.58	7.04	17.20	21.78
c_tp	3.84	11.07	6.05	3.22	7.80	6.68	9.51
d_tp	7.56	11.75	9.41	7.59	8.19	19.15	21.45
d_tpm	8.99	11.46	12.27	9.09	9.30	18.99	19.49

TABLE VI
SHARPNESS $S(0.95)$ IN $\frac{m^2}{s}$ OF ALL MODELS FOR CYCLISTS.

model	all	wait	start	move	stop	turn left	turn right
c_t	4.72	1.76	6.41	4.84	5.51	9.34	7.34
d_t	4.61	2.48	6.61	4.66	5.57	8.19	7.52
c_tp	5.38	2.56	12.22	5.26	7.14	14.92	14.35
d_tp	3.47	2.01	5.36	3.48	4.13	5.87	5.87
d_tpm	3.33	1.70	5.03	3.37	3.81	6.14	5.85

TABLE VII
ASWAE IN $\frac{m}{s}$ OF ALL MODELS FOR CYCLISTS.

model	all	wait	start	move	stop	turn left	turn right
c_t	0.68	0.37	0.93	0.69	0.79	1.10	1.08
d_t	0.67	0.37	0.88	0.67	0.79	1.04	1.06
c_tp	0.64	0.40	1.09	0.63	0.79	1.18	1.15
d_tp	0.63	0.37	0.81	0.64	0.76	0.99	1.00
d_tpm	0.62	0.37	0.82	0.63	0.74	0.99	0.99

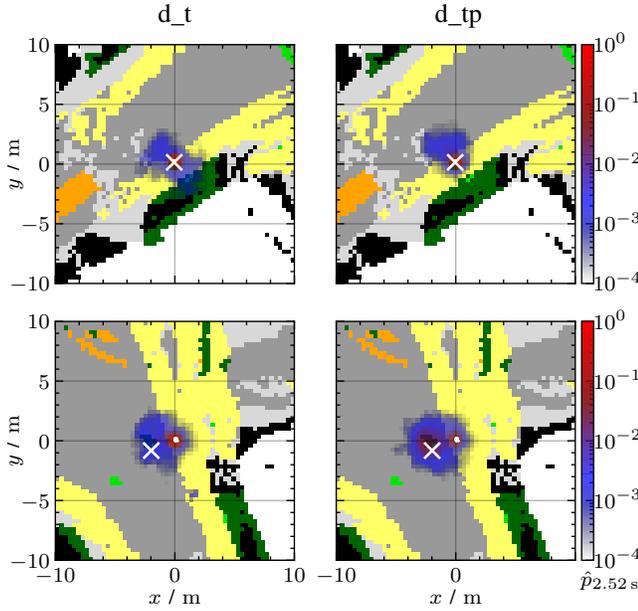


Fig. 4. Comparison of two exemplary forecasts of the models d_t (left) and d_{tp} (right) for pedestrians and a forecasted time horizon of 2.52 s. The forecasted distributions are illustrated using the logarithmic color scale on the right ranging from red (high probability) to blue (low probability). The past trajectory is represented by a white line/point and the actual position by a white cross. The semantic maps can be seen in the background. (Best viewed on screen).

to exemplify their effects. The example in the top row shows the forecasts 2.52 s into the future of a pedestrian waiting at the curbside. The forecast based on the past trajectory (left) indicates continued standing as the most likely event and a low probability of starting in different directions. In comparison, the model d_{tp} (right) can take the body orientation provided by the 3D poses into account. Accordingly, the forecast of a possible starting is made mainly towards the road. The second example (lower row) presents a pedestrian who initially stands at the curbside and then crosses the road. In contrast to the trajectory based model, the model d_{tp} using poses detects the pedestrian’s intention to start walking and forecasts a multimodal distribution covering a possible standing as well as crossing the street. This is an advantage of the discrete over the continuous methods used in this article, which cannot make multimodal forecasts. Overall, the forecasts for pedestrians are consistently improved by using poses for both the continuous and discrete approach. All motion types benefit from the additional information provided by poses. The improvement is indicated by better positional accuracy and reliability coupled with enhanced sharpness.

For cyclists, the ASWAEE is improved by 6.0% for the continuous method (c_t and c_{tp} in Tab. VII) and 4.9% for the discrete approach (d_t and d_{tp} in Tab. VII) by using poses. For the continuous model, the poses also lead to a considerable improvement in reliability (Tab. V) with slightly larger distributions with respect to the area (Tab. VI). The area covered by a confidence level of 95% is huge for turning cyclists since, in such situations, there is great variability and uncertainty. This is also evident from the high values

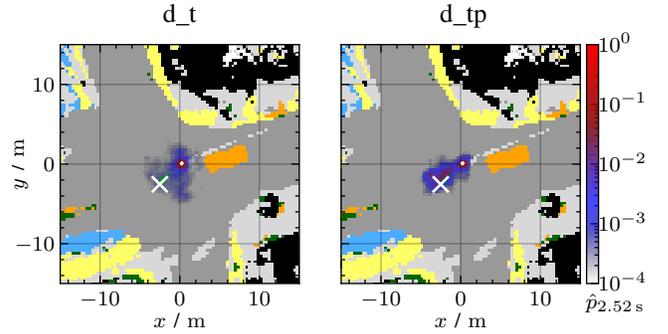


Fig. 5. Comparison of exemplary forecasts of the models d_t (left) and d_{tp} (right) for cyclists and a forecasted time horizon of 2.52 s.

obtained for the ASWAEE for these motion types (*turn left* and *turn right* in Tab. VII). While the use of poses improves the positional accuracy (Tab. VII) for the discrete method (model d_{tp} vs d_t), the reliability (Tab. V) is only enhanced for the motion types *wait*, *start*, and *turn right*. For the remaining motion types, the reliability deteriorates as the model tends to forecast overconfident distributions. This is evident in low values for the sharpness (d_{tp} in Tab. VI) and could be remedied by a separate optimization of the standard deviations for the spatial label smoothing. Fig. 5 presents an exemplary forecast for cyclists with (right) and without (left) poses. In this situation, a cyclist wants to turn left at an intersection. First, the cyclist stands in order to let oncoming traffic pass before initiating the turn. Similar to what we have seen with pedestrians, the model d_{tp} forecasts a multimodal distribution that covers continued standing as well as turning. The poses also provide the movement direction of the cyclist.

2) *Impact of Semantic Maps*: The semantic maps’ impact on the forecasts is difficult to evaluate because no measure is known that allows for a quantification of the benefits. The meaning of the surrounding scene can be quite different from case to case. Therefore, identifying potential relations and the definition of reasonable measures are difficult. For pedestrians, the additional use of the semantic maps leads to better reliability (d_{tpm} in Tab. II) while the ASWAEE (d_{tpm} in Tab. IV) remains roughly the same compared to model d_{tp} . For cyclists, the ASWAEE also remains similar (d_{tp} vs d_{tpm} in Tab. VII), while the reliability (Tab. V) deteriorates due to overconfident forecasts. However, these measures reflect the influence of semantic maps only to a limited extent. Therefore, in the following, we illustrate the impact of semantic maps qualitatively by some examples. In addition, we consider the occupancy score O_{t_f} as introduced in Sec. III-A4 since it allows the quantification of the influence of obstacles on the forecasts. Tab. VIII reports the occupancy score for pedestrians and cyclists, the different discrete models, and each forecasted time horizon. For each model, the values for cyclists are smaller than those for pedestrians (e.g. row 3 vs row 6). This is plausible since cyclists are usually on the road and thus have a greater distance to static obstacles. Moreover, the generally small values in Tab. VIII indicate a relatively small impact of the static obstacles on the forecasts in our dataset and the forecasted time horizons considered here. The score represents

TABLE VIII
OCCUPANCY SCORE O_{t_f} IN % OF ALL DISCRETE MODELS FOR DIFFERENT
FORECASTED TIME HORIZONS AND VRU TYPES.

	model	0.44 s	0.96 s	1.48 s	2.00 s	2.52 s
pedestrians	d_t	0.39	0.58	0.90	1.35	1.99
	d_tp	0.37	0.53	0.77	1.12	1.58
	d_tpm	0.26	0.37	0.48	0.58	0.68
cyclists	d_t	0.13	0.16	0.28	0.41	0.64
	d_tp	0.13	0.16	0.27	0.40	0.63
	d_tpm	0.09	0.12	0.17	0.17	0.17

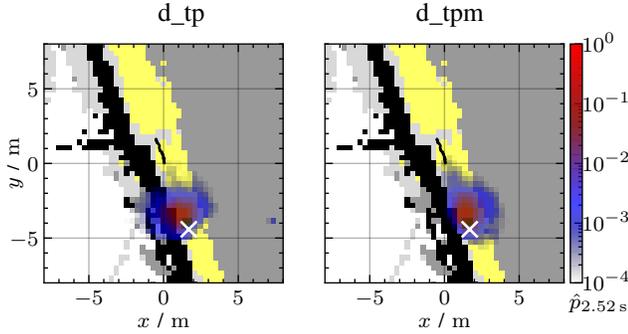


Fig. 6. Comparison of exemplary forecasts of the models d_tp (left) and d_tpm (right) for pedestrians and a forecasted time horizon of 2.52 s. For better visibility, the past trajectory is shown as black line here.

the average over the entire dataset. While for most pedestrians and cyclists, the future trajectory is sufficiently determined by the past trajectory and poses, obstacles are indeed important for individual examples. For the models d_t and d_tp without use of the semantic maps, the score in Tab. VIII rises as the forecasted time horizon grows, indicating an increase in the importance of the obstacles for the forecasts. The scores for all time horizons are reduced by considering the semantic maps for both pedestrians and cyclists (d_tpm in Tab. VIII) which means that the forecasts of this model collide less with static obstacles. It should be noted that the optimal score is not necessarily zero, since the semantic maps contain errors that can be considered by the forecasting models. Fig. 6 illustrates the forecasts of model d_tp (left) and d_tpm (right) of a pedestrian walking on the sidewalk alongside a house wall. While the forecasted distribution of model d_tp is approximately symmetric with respect to the movement direction, model d_tpm considers the house wall resulting in a skewed distribution. For cyclists, model d_tpm in particular takes into account the course of the road and the curbside. An example of this can be seen in Fig. 7. It shows a cyclist who is just entering an intersection. In addition to driving straight and turning left, model d_tpm also forecasts a low probability for turning right following the curb.

3) *Impact of Cell Size:* To investigate the effects of the grid cell size, the model d_tp is trained for pedestrians with a cell size of $e^2 = 0.175 \times 0.175 \text{ m}^2$ and compared with the previously reported results using a cell size of $0.35 \times 0.35 \text{ m}^2$. Increasing the resolution offers two potential advantages: First, the model can provide a finer resolution of the forecasted distribution, and second, the actual position used to train the

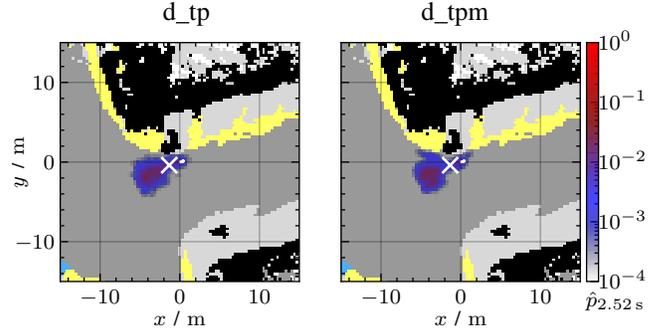


Fig. 7. Comparison of exemplary forecasts of the models d_tp (left) and d_tpm (right) for cyclists and a forecasted time horizon of 2.52 s.

TABLE IX
THE ACHIEVED RESULTS OF THE MODEL d_tp FOR PEDESTRIANS USING A
CELL SIZE OF $0.175 \times 0.175 \text{ m}^2$ FOR EACH MOTION TYPE. FOR
COMPARISON, THE RESULTS OF THE SAME MODEL FOR CELL SIZE
 $0.35 \times 0.35 \text{ m}^2$ ARE PROVIDED IN PARENTHESES.

motion type	ECE in %	$S(0.95)$ in $\frac{\text{m}^2}{\text{s}}$	ASWAAE in $\frac{\text{m}}{\text{s}}$
all	2.99 (4.48)	2.81 (3.03)	0.53 (0.56)
wait	2.90 (5.42)	3.05 (2.98)	0.49 (0.53)
start	4.94 (3.48)	3.04 (3.32)	0.58 (0.59)
move	2.44 (4.34)	2.63 (2.93)	0.51 (0.55)
stop	7.71 (8.16)	2.91 (3.17)	0.58 (0.60)

model is also resolved more precisely. Both can have a beneficial effect on the forecasting results. However, this comes at the cost of increased demand for computational and memory resources. The training time for a single batch of size 40 increases from 7.1 ms to 18.6 ms using an Nvidia RTX 2080 Ti and the required graphics memory from 1.8 GB to 5.0 GB. The increase is even greater for cyclists due to the larger grid: Here, the training time for a single batch increases from 21.9 ms to 83.3 ms and the needed memory from 3.1 GB to 10.3 GB. As a result, a detailed examination of the effects of the grid cell size is not possible for cyclists and we limit ourselves to pedestrians. However, it should be noted here that the cell size has little effect on the inference time: For pedestrians, the inference time remains at 1.1 ms for both resolutions, while for cyclists it increases from 1.1 ms to 1.6 ms. This shows that the high resource demand for a fine resolution is required mainly during training and not during the use of the models.

Tab. IX shows the results for pedestrians obtained using the smaller cell size. Compared to the larger cell size (given in parentheses in Tab. IX), there is a considerable improvement in all scores. This model even outperforms the continuous model c_tp in terms of reliability (c_tp in Tab. II vs Tab. IX) while achieving slightly worse results for the ASWAAE (c_tp in Tab. IV vs Tab. IX). A look at the corresponding reliability diagram on the right in Fig. 8 reveals the main reason for the improvement in reliability. Compared to the reliability diagram of model d_tp with the larger cell size (on the left in Fig. 8), the S-shape of the curves is reduced. This particularly affects short forecasted time horizons, since for them, the forecasted distributions are spread over relatively small areas.

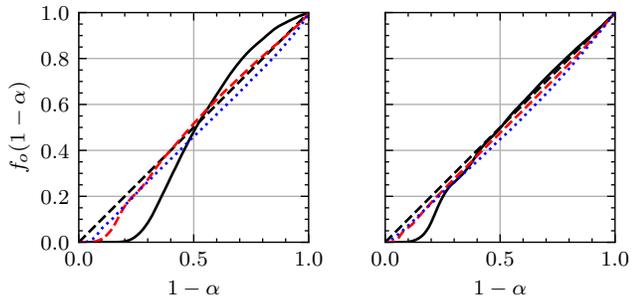


Fig. 8. Reliability diagrams for model d_{tp} and pedestrians using spatial label smoothing and a cell size of $0.35 \times 0.35 \text{ m}^2$ (left) and $0.175 \times 0.175 \text{ m}^2$ (right) evaluated on the test dataset. The reliability is shown for the time horizons 0.44 s (black solid line), 1.48 s (dashed red line) and 2.52 s (dotted blue line). The ideal diagonal is plotted as dashed black line.

4) *Comparison of Discrete and Continuous Forecasting Models:* Overall, the continuous models c_{t} and c_{tp} achieve better results for pedestrians in terms of reliability (Tab. II), sharpness (Tab. III), and ASWAE (Tab. IV) compared to the discrete models d_{t} and d_{tp} with cell size $0.35 \times 0.35 \text{ m}^2$. In contrast, for cyclists, using solely the past trajectory as input, the discrete model achieves better scores (c_{t} vs d_{t} in Tab. V to VII). The use of poses and semantic maps (d_{tp} and d_{tpm}), on the other hand, leads to overconfident forecasts so that the reliability of the continuous model c_{tp} is in front. Both approaches have specific advantages and disadvantages: The used continuous method forecasts a single Gaussian distribution, so the type of the distribution is fixed and multimodal distributions cannot be represented. This disadvantage can be observed, for example, in the reliability of waiting pedestrians (c_{tp} vs d_{tp} for *wait* in Tab. II). Here, the continuous method c_{tp} cannot express the unlikely but possible chance of a starting motion. The discrete method d_{tp} has no such limitations. As the examples have demonstrated, it can express unsymmetric and multimodal distributions. However, this comes at the cost of discretizing the space leading to less accurate forecasts. The positions of the VRUs can only be forecasted within the chosen grid, and the required computational resources for the training limit the cell size. The use of the cross entropy as loss function leads to overconfident distributions. It can be compensated by spatial label smoothing, but the determination of the necessary parameters is a tedious process. On the other hand, the discrete method allows the consideration of the surroundings in the form of semantic maps. Their impact depends on the particular situation and the forecasted time horizon. There is no continuous approach in the literature allowing the consideration of semantic maps. It is difficult for continuous methods to learn and predict arbitrary distributions and to adjust them to the semantic maps.

V. CONCLUSIONS AND FUTURE WORK

In this article, we presented an approach for probabilistic trajectory forecasting of pedestrians and cyclists considering past movements represented by 3D poses and the surroundings in the form of semantic maps. The forecasts are generated in discrete grids allowing the model to learn arbitrary distribu-

tions. The impact of poses and semantic maps on the forecasts was examined, and the forecasted distributions were evaluated by their reliability, sharpness, and positional accuracy. We compared our approach with a method for forecasting continuous Gaussian distributions discussing their respective advantages and disadvantages. Using the plain cross entropy as loss function leads to overconfident distributions. This can be prevented by applying spatial label smoothing. The resolution of the grids has a non-negligible influence on the quality of the results. Here, a trade-off must be made, taking into account the required computational and memory resources. Overall, the use of 3D poses improves the forecasts, especially through better detection of motion type changes and body orientation. To adequately forecast the behavior of VRUs, multimodal and skewed distributions are advantageous which are possible with the proposed method. The semantic maps allow a precise adaptation of the forecasts to the individual situation. While this leads to major improvements of the forecasts in individual cases, the semantic maps have an overall subordinate impact on the entire dataset and the forecasted time horizon of 2.52 s. However, their impact increases for longer time horizons. Therefore, the surroundings must be taken into account in long-term forecasting.

Our future work will focus on using our approach to model specific motion types to improve the forecasts. A preceding motion type detection should be used to weigh the individual models. Furthermore, we want to improve the quality of the semantic maps contributing to better trajectory forecasts. The impact of considering semantic maps in motion type detection of VRUs will also be investigated.

ACKNOWLEDGMENT

This work was supported by “Zentrum Digitalisierung.Bayern”. In addition, the work is backed by the project DeCoInt², supported by the German Research Foundation (DFG) within the priority program SPP 1835: “Kooperativ interagierende Automobile”, grant numbers DO 1186/1-2 and SI 674/11-2.

REFERENCES

- [1] C. Keller and D. Gavrilu, “Will the pedestrian cross? a study on pedestrian path prediction,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 15, no. 2, pp. 494–506, 2014.
- [2] M. Goldhammer, S. Köhler, S. Zernetsch, K. Doll, B. Sick, and K. Dietmayer, “Intentions of vulnerable road users—detection and forecasting by means of machine learning,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 7, pp. 3035–3045, 2020.
- [3] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, “Social GAN: Socially acceptable trajectories with generative adversarial networks,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2255–2264.
- [4] F. Marchetti, F. Becattini, L. Seidenari, and A. D. Bimbo, “MANTRA: Memory augmented networks for multiple trajectory prediction,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 7141–7150.
- [5] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, “Social lstm: Human trajectory prediction in crowded spaces,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 961–971.
- [6] E. A. I. Pool, J. F. P. Kooij, and D. M. Gavrilu, “Context-based cyclist path prediction using Recurrent Neural Networks,” in *IEEE Intelligent Vehicles Symposium (IV)*, 2019, pp. 824–830.

- [7] J. Eilbrecht, M. Bieshaar, S. Zernetsch, K. Doll, B. Sick, and O. Stursberg, "Model-predictive planning for autonomous vehicles anticipating intentions of vulnerable road users by artificial neural networks," in *IEEE Symposium Series on Computational Intelligence (SSCI)*, 2017, pp. 1–8.
- [8] M. Koschi, C. Pek, M. Beikirch, and M. Althoff, "Set-based prediction of pedestrians in urban environments considering formalized traffic rules," in *IEEE Intelligent Transportation Systems Conference (ITSC)*, 2018, pp. 2704–2711.
- [9] S. Zernetsch, H. Reichert, V. Kress, K. Doll, and B. Sick, "Trajectory forecasts with uncertainties of vulnerable road users by means of neural networks," in *IEEE Intelligent Vehicles Symposium (IV)*, 2019, pp. 810–815.
- [10] M. Bieshaar, J. Schreiber, S. Vogt, A. Gensler, and B. Sick, "Quantile surfaces – generalizing quantile regression to multivariate targets," arXiv: 2010.05898, <https://arxiv.org/abs/2010.05898>, 2020.
- [11] J. Wu, J. Ruenz, and M. Althoff, "Probabilistic Map-based Pedestrian Motion Prediction Taking Traffic Participants into Consideration," in *IEEE Intelligent Vehicles Symposium (IV)*, 2018, pp. 1285–1292.
- [12] A. Jain, S. Casas, R. Liao, Y. Xiong, S. Feng, S. Segal, and R. Urtasun, "Discrete Residual Flow for Probabilistic Pedestrian Behavior Prediction," in *Conference on Robot Learning (CoRL)*, ser. Proceedings of Machine Learning Research, vol. 100, 2019, pp. 407–419.
- [13] R. Quintero, I. Parra, D. Fernández-Llorca, and M. A. Sotelo, "Pedestrian path, pose, and intention prediction through gaussian process dynamical models and pedestrian activity recognition," *IEEE Transactions on Intelligent Transportation Systems*, vol. 20, no. 5, pp. 1803–1814, 2019.
- [14] V. Kress, S. Zernetsch, K. Doll, and B. Sick, "Pose based trajectory forecast of vulnerable road users," in *IEEE Symposium Series on Computational Intelligence (SSCI)*, 2019, pp. 1200–1207.
- [15] H. Xue, D. Q. Huynh, and M. Reynolds, "SS-LSTM: A Hierarchical LSTM Model for Pedestrian Trajectory Prediction," in *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2018, pp. 1186–1194.
- [16] H. Cheng, W. Liao, M. Y. Yang, M. Sester, and B. Rosenhahn, "MCNET: Multi-Context Encoder Network for Homogeneous Agent Trajectory Prediction in Mixed Traffic," in *IEEE Intelligent Transportation Systems Conference (ITSC)*, 2020, pp. 1–8.
- [17] D. Ridel, N. Deo, D. Wolf, and M. Trivedi, "Scene Compliant Trajectory Forecast With Agent-Centric Spatio-Temporal Grids," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 2816–2823, 2020.
- [18] V. Kress, S. Zernetsch, M. Bieshaar, G. Reitberger, E. Fuchs, K. Doll, and B. Sick, "Pedestrians and Cyclists in Road Traffic: Trajectories, 3D Poses and Semantic Maps," 2021. [Online]. Available: <https://doi.org/10.5281/zenodo.4898838>
- [19] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1302–1310.
- [20] D. Tome, C. Russell, and L. Agapito, "Lifting from the deep: Convolutional 3d pose estimation from a single image," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5689–5698.
- [21] V. Kress, J. Jung, S. Zernetsch, K. Doll, and B. Sick, "Human pose estimation in real traffic scenes," in *IEEE Symposium Series on Computational Intelligence (SSCI)*, 2018, pp. 518–523.
- [22] Z. Wu, C. Shen, and A. van den Hengel, "Wider or Deeper: Revisiting the ResNet Model for Visual Recognition," *Pattern Recognition*, vol. 90, pp. 119–133, 2019.
- [23] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention – MICCAI*, ser. Lecture Notes in Computer Science, 2015, pp. 234–241.
- [24] X. Shi, Z. Chen, H. Wang, D. Yeung, W. Wong, and W. Woo, "Convolutional LSTM network: A machine learning approach for precipitation nowcasting," in *International Conference on Neural Information Processing Systems (NIPS)*, 2015, pp. 802–810.
- [25] M. P. Naeni, G. F. Cooper, and M. Hauskrecht, "Obtaining Well Calibrated Probabilities Using Bayesian Binning," in *AAAI Conference on Artificial Intelligence*, 2015, pp. 2901–2907.
- [26] A. H. Murphy, "A New Vector Partition of the Probability Score," *Journal of Applied Meteorology (1962-1982)*, vol. 12, no. 4, pp. 595–600, 1973.
- [27] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, 2015.
- [28] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, "On Calibration of Modern Neural Networks," in *International Conference on Machine Learning*. PMLR, 2017, pp. 1321–1330.



Viktor Kress received the B.Eng. degree in Mechatronics and the M.Eng. degree in Electrical Engineering and Information Technology from the University of Applied Sciences Aschaffenburg, Germany, in 2015 and 2016, respectively. Currently, he is working on his PhD thesis in cooperation with the Faculty of Electrical Engineering and Computer Science of the University of Kassel, Germany. His research interests include sensor data fusion, pattern recognition, machine learning, behavior recognition and trajectory forecasting of traffic participants.



Fabian Jeske received a B.Eng. degree Industrial Engineering from the University of Applied Sciences Aschaffenburg, Germany in 2019. At the moment he is working on his M.Sc. in Industrial Engineering at the University of Applied Sciences Aschaffenburg. His research interests include trajectory forecasting of traffic participants, sensor data fusion and pattern recognition.



Stefan Zernetsch received the B.Eng. and the M.Eng. degree in Electrical Engineering and Information Technology from the University of Applied Sciences Aschaffenburg, Germany, in 2012 and 2014, respectively. Currently, he is working on his PhD thesis in cooperation with the Faculty of Electrical Engineering and Computer Science of the University of Kassel, Germany. His research interests include cooperative sensor networks, data fusion, multiple view geometry, pattern recognition and behavior recognition of traffic participants.



Konrad Doll received the Diploma (Dipl.-Ing.) degree and the Dr.-Ing. degree in Electrical Engineering and Information Technology from the Technical University of Munich, Germany, in 1989 and 1994, respectively. In 1994 he joined the Semiconductor Products Sector of Motorola, Inc. (now Freescale Semiconductor, Inc.). In 1997 he was appointed to professor at the University of Applied Sciences Aschaffenburg in the field of computer science and digital systems design. His research interests include intelligent systems, their real-time implementations,

and their applications in advanced driver assistance systems and automated driving. He received several thesis and best paper awards. Konrad Doll is member of the IEEE.



Bernhard Sick received the diploma, the Ph.D. degree, and the "Habilitation" degree, all in computer science, from the University of Passau, Germany, in 1992, 1999, and 2004, respectively. Currently, he is full Professor for Intelligent Embedded Systems at the Faculty for Electrical Engineering and Computer Science of the University of Kassel, Germany. There, he is conducting research in the areas autonomic and organic computing and technical data analytics with applications in, e.g., energy systems, automotive engineering, physics and materials science. He

authored more than 200 peer-reviewed publications in these areas. Dr. Sick is associate editor of the IEEE TRANSACTIONS ON CYBERNETICS. He holds one patent and received several thesis, best paper, teaching, and inventor awards. He is a member of IEEE (Systems, Man, and Cybernetics Society, Computer Society, and Computational Intelligence Society) and GI (Gesellschaft fuer Informatik).