

Pedestrian Behavior Prediction for Automated Driving: Requirements, Metrics, and Relevant Features

Michael Herman, Jörg Wagner, Vishnu Prabhakaran, Nicolas Möser,
Hanna Ziesche, Waleed Ahmed, Lutz Bürkle, Ernst Kloppenburg, and Claudius Gläser

Abstract—Automated vehicles require a comprehensive understanding of traffic situations to ensure safe and anticipatory driving. In this context, the prediction of pedestrians is particularly challenging as pedestrian behavior can be influenced by multiple factors. In this paper, we thoroughly analyze the requirements on pedestrian behavior prediction for automated driving via a system-level approach. To this end we investigate real-world pedestrian-vehicle interactions with human drivers. Based on human driving behavior we then derive appropriate reaction patterns of an automated vehicle and determine requirements for the prediction of pedestrians. This includes a novel metric tailored to measure prediction performance from a system-level perspective. The proposed metric is evaluated on a large-scale dataset comprising thousands of real-world pedestrian-vehicle interactions. We furthermore conduct an ablation study to evaluate the importance of different contextual cues and compare these results to ones obtained using established performance metrics for pedestrian prediction. Our results highlight the importance of a system-level approach to pedestrian behavior prediction.

Index Terms—Autonomous vehicles, Automated driving, Prediction methods, Machine learning.

I. INTRODUCTION

ROAD safety is a key driver for the development of Driver Assistance (DA) and Automated Driving (AD) systems. According to a report of the World Health Organization [1] traffic accidents cause more than 1.3 million fatalities annually, almost half of them being Vulnerable Road Users (VRUs). Therefore, the protection of VRUs, in particular pedestrians, constitutes a major goal of intelligent vehicles. The Automatic Emergency Braking system for Pedestrians (AEB-P) is a good example on how driver assistance systems already protect pedestrians today. AEB-P detects pedestrians in the predicted vehicle's path and, if a collision cannot be avoided by the driver, automatically initiates emergency braking. By either avoiding the collision or, if avoidance is not possible, reducing the velocity of an impact, pedestrian AEB systems mitigate pedestrian fatality and injury [2]. The detection of pedestrians and the prediction of their behavior are essential components of an AEB-P system. Prediction of an AEB-P is generally

M. Herman, J. Wagner, V. Prabhakaran, H. Ziesche, and E. Kloppenburg are with the Bosch Center for Artificial Intelligence, Germany.

N. Möser, L. Bürkle, and C. Gläser are with the Robert Bosch GmbH, Corporate Research, Germany.

W. Ahmed is with the Robert Bosch GmbH, Cross-Domain Computing Solutions – Automated Driving, Germany.

e-mail: Michael.Herman@de.bosch.com



Fig. 1. Exemplary scenario in which the behavior of a pedestrian depends on multiple contextual cues, e.g. the present road infrastructure or interactions with other traffic participants.

restricted to short prediction horizons in the order of 1 to 2s and is typically based on kinematic models.

On the other hand, automated driving not only addresses near-collision situations, but broadens the scope to everyday driving scenarios. Thus, besides collision mitigation, comfortable driving that imitates human driving behavior shifts into focus. In order to react appropriately at an early stage automated driving requires an extended pedestrian behavior prediction to correctly reason about a situation on a longer time-scale.

Recent years have seen an increased interest in using deep learning based methods for the purpose of long-term pedestrian prediction. Most of the works on this topic however used generic metrics both for development and evaluation of the prediction models. While these generic metrics are suitable for measuring the overall accuracy of a predicted behavior, they do not take into account the actual requirements of downstream functions, like e.g. an automated driving system. We argue that due to this, important task-specific requirements are not considered in model development and evaluation, since generic metrics do not or only partially cover those requirements. As a result the proposed models are often too complex or suboptimal for the downstream task. Therefore, we propose a new function-specific metric, which is based on system-level requirements. This metric allows for a more task-informed assessment of models for pedestrian prediction.

Especially on a longer time-scale the behavior of pedestrians cannot be investigated in isolation, but rather has to be considered within the context of the overall traffic scene. The exemplary scenario of a girl running along a sidewalk depicted in Fig. 1 illustrates this relationship: The future behavior of the

girl is likely to depend on a variety of factors. This includes the static driving infrastructure (e.g. the road layout, the zebra crossing), interactions of the pedestrian with other traffic participants (e.g. an approaching vehicle), and appearance or communication cues (e.g. gestures). Yet, the importance and influence of individual contextual cues on prediction accuracy and on the downstream AD function is not obvious a priori.

In this paper, we derive requirements and a performance metric for pedestrian behavior prediction in the context of automated driving. For evaluation purposes, we present a prediction model based on a Conditional Variational Autoencoder (CVAE) that specifically addresses long-term prediction by including contextual cues of the traffic scene. Finally, we investigate the importance of contextual cues in terms of prediction accuracy. Our results highlight the importance of a system-level approach to pedestrian behavior prediction. In detail, the contributions of the paper are:

- An appropriate system reaction pattern for interactions of an automated vehicle with pedestrians is derived from an analysis of human driving behavior.
- Requirements for pedestrian behavior prediction in automated driving are specified and a corresponding metric to assess prediction performance is derived.
- On the basis of the proposed metric, a CVAE prediction model is evaluated on a large-scale dataset comprising thousands of real-world pedestrian-vehicle interactions.
- The relevance of different contextual cues is assessed based on an ablation study and the results are compared to the ones obtained using established performance metrics for pedestrian prediction.

The remainder of this paper is organized as follows: We first give an overview on related work in Sec. II. We then introduce our large-scale dataset of vehicle-pedestrian interactions in Sec. III. In Sec. IV, we analyze pedestrian-vehicle interactions to determine human driving behavior and to derive requirements on pedestrian prediction for automated driving. Our prediction model is introduced in Sec. V and thoroughly evaluated in Sec. VI. Finally, we conclude the paper with a discussion of our results in Sec. VII.

II. RELATED WORK

In this section, we give an overview of state of the art approaches as well as typical input features applied to pedestrian prediction. Furthermore, evaluation metrics and public datasets are briefly summarized.

A. Pedestrian Prediction Approaches

Pedestrian prediction has already been studied for a long time, resulting in numerous approaches that address this problem in the automotive domain [3], [4] and beyond [5]. It is important to note, however, that in the context of Advanced Driver Assistance Systems (ADAS) and AD the employed output representation of prediction models and consequently the system integration may significantly differ. Approaches range from the prediction of pedestrian crossing intentions [6], [7], over walking destinations [8], [9], to the prediction of paths [10] or trajectories [11], where the latter can further

be distinguished into trajectory prediction in image-view and bird's-eye view.

We consider a prediction of bird's-eye view pedestrian trajectories best suited for an AD system as the downstream planning of appropriate system reactions usually relies on this kind of representation. That is why we will focus on such approaches in the following.

To the best of our knowledge, requirements on pedestrian behavior prediction for AD have not been thoroughly analyzed so far. We strongly believe that it is essential to mirror pedestrian prediction to requirements of the downstream task in order to obtain an optimal overall system performance. By taking a system-level approach to pedestrian behavior prediction, this paper provides an important contribution for this.

B. Trajectory Prediction Models

Many traditional approaches for predicting the future motion of traffic participants depend on a set of explicitly defined dynamics equations that are generally derived from physics-based motion models [5]. Often, these approaches are used in combination with Probabilistic Graphical Models [12], [13], [14].

Recently, pattern-based methods that learn behavioral patterns from data have outperformed traditional approaches. Especially, deep learning based solutions became state of the art for most of the problems related to public datasets. Often, Recurrent Neural Networks (RNN) are used for encoding trajectories of interacting agents and decoding future behavior [15], [16]. One of the major problems associated with these approaches is to accurately capture the probabilistic, multi-modal distribution over trajectories. In order to address this issue recent deep learning based methods predict parametric distributions [15], learn mixtures of Gaussian trajectory distributions [17], use adversarial training approaches [18], or introduce discrete [19], [20] or continuous [21], [22], [23], [24] latent variables.

For the investigations in this paper, we use models that build on Conditional Variational Autoencoders (CVAE) [25]. Their use of continuous latent variables renders them suitable for capturing complex, multi-modal probability distributions.

C. Contextual Cues

Human behavior is influenced by contextual cues of internal and external stimuli. The survey [5] groups them into three categories: cues of the target agent, the dynamic environment, and the static environment. Potential target agent cues are the motion state (e.g. position, velocities) [14], [15], [16], appearance-based cues (e.g. head or body pose) [26], or semantic attributes (e.g. age or gender) [27]. While traditional models often do not take into account influences of the dynamic environment [28], [29], [30], other approaches exist that model interactions with other agents [31], [32], [15] or even social groups [33]. Regarding cues of the static environment, there are several approaches that neglect this influence [34], [35], some others only model influence of individual static objects [36], while still others model more complex

TABLE I
COMPARISON OF EXISTING DATASETS

Dataset	Representation		Features			recording time [min]	Size # pedestrian tracks (with attributes / total)
	image plane	top view	map	ego	pedestrian attributes		
STIP [40]	x					923	- / 25,000
JAAD [41]	x				x	43	686 / 2,786
PIE [42]	x			x	x	360	1,842 / 1,842
TITAN [43]	x			x	x	175	8,592 / 8,592
SDD [33]		x	x			620	- / 11,216
INTERACTION [44]		x	x			991	- / 1,700
inD [45]		x	x			600	- / 3,107
Argoverse [46]	x	x	x	x		30	- / 1,322
nuScenes [47]	x	x	x	x		333	- / 11,512
PePScenes [48]	x	x	x	x	x	333	719 / 11,512
ours	x	x	x	x	x	4,434	9,438 / 93,162

influences from environment geometry and topology [19], [37]. In addition, recent work [38] studies the contribution of different contextual cues on an action classification task.

D. Evaluation Metrics

Performance evaluation is an integral part of the process of creating a prediction model. The survey [5] extensively discusses different metrics for models that predict trajectories. These metrics fall into two classes, geometric and probabilistic. A widely used geometric metric is the Average Displacement Error (ADE), which applies to models providing point predictions. ADE measures the euclidean distance of a predicted trajectory from ground truth positions at a specific prediction time interval, averaged over the trajectory, or over multiple trajectories. Quite commonly the ADE metric is also applied to probabilistic predictions by averaging over the predictive distribution as well.

Probabilistic metrics are used for models that provide predictions in the form of probability densities. This kind of metrics measures how well the predictions capture the uncertainties inherent to the prediction process as well as the true process. A typical metric here is average Negative Log Likelihood (NLL) of the ground truth positions. Unfortunately, these kind of metrics tend to lack intuitive interpretability. One subtlety with the different probabilistic metrics is whether they encourage multimodal predictive distributions or not.

This raises the question whether the metrics discussed so far measure properties relevant for possible applications of the model under consideration. The authors of [39] discuss the problem of evaluating generative (probabilistic) models and come to the conclusion that application specific metrics are generally required.

E. Existing Datasets

In contrast to the large number of datasets devoted to pedestrian detection there is only a limited number of datasets available that address pedestrian prediction in an automated driving context. Table I summarizes and compares the most important ones. The comparison takes into account the employed representation (image plane vs. top view), the availability of features which are relevant for pedestrian prediction

(a semantic map, ego-vehicle data, and detailed pedestrian attributes at a perceptual and/or behavioral level), and the size of the datasets (in terms of the number of unique pedestrian tracks and recording time).

Datasets from the first group provide video recordings including annotated objects. Some of them additionally include rich behavioral annotations for pedestrians. The most notable contributions to this group include the Joint Attention in Autonomous Driving (JAAD) dataset [41], the Pedestrian Intention Estimation (PIE) dataset [42], and the TITAN dataset [43]. However, even though these datasets are suitable for benchmarking vision-based algorithms for pedestrian detection or intention recognition, they are limited in use for investigating trajectory prediction where an object representation in bird's-eye view is required.

Datasets from the second group like the Stanford Drone Dataset (SDD) [33], the INTERACTION dataset [44], or the Intersection Drone dataset (inD) [45] comprise trajectories from a bird's-eye view and are thus suitable for pedestrian prediction. However, since drones were used to record these datasets they are limited to a small number of locations, and more important, do not provide detailed pedestrian attributes.

Recently, various large scale automotive datasets were published e.g. Argoverse [49] and nuScenes [50], many of which also comprise tracking or motion forecasting challenges. However, none of these datasets provide detailed pedestrian attributes and mainly focus on the prediction of vehicles trajectories. Furthermore, they lack a significant number of scenes with pedestrian-vehicle interactions. The authors of [48] conducted a post-labeling of nuScenes to provide additional pedestrian attributes. However, for their PePScenes dataset only 6% of the pedestrians in nuScenes were taken into account, since all others were not of interest for the driving task.

III. DATASET

To overcome the limitations of existing datasets, we created a large-scale dataset that specifically addresses pedestrian prediction in automated driving context. In detail, the requirements on the datasets were as follows:

- 2D pedestrian positions in the ground plane.

- Annotations for additional features that are relevant for the prediction task. Most notably, a semantic map of the road infrastructure, ego-vehicle data for modeling pedestrian-vehicle interactions, and additional appearance-based attributes for pedestrians.
- A large number of pedestrians that are relevant for the driving task. The respective pedestrian-vehicle interactions shall cover various scenarios, e.g. interactions at different traffic control elements.

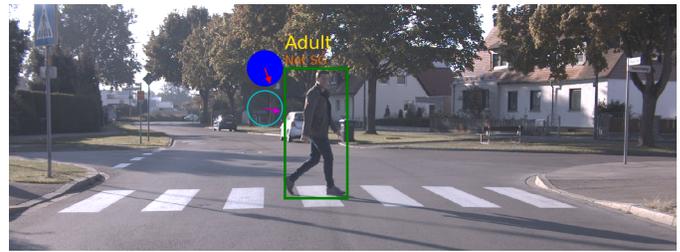


Fig. 3. Example of labeled pedestrian attributes: The compass plots depict head and body orientation, where the filled upper compass indicates that the pedestrian is looking at the ego-vehicle.

A. Data Collection

The dataset we refer to in this paper comprises vehicle-pedestrian interactions from inner-city traffic in southern Germany. The data was recorded on three different round-courses with lengths between 2 and 4 km. The routes were chosen to maximize variability of traffic scenarios including both downtown and suburban areas, different road sizes, traffic densities, and number of pedestrians. Furthermore, the routes contain various traffic control elements which are relevant for vehicle-pedestrian interactions, most notably zebra crossings, pedestrian refuge islands, or a combination of both.

To further increase scenario coverage a number of actors were positioned at different locations along the courses. The instruction of actors followed a semi-scripted approach where actors were told not to perform specific interactions with the recording vehicle but rather to arbitrarily vary their walking routes, interactions and behavior in a realistic manner. Since actors had to adapt their behavior to the respective traffic situation, including the recording vehicle and other traffic participants, the setup proved to result in a great variety of realistic vehicle-pedestrian interactions.

The recordings were carried out during three weeks in fall of 2018. Overall, we recorded 74 hours of data using seven different drivers to reduce driving style biases. Data acquisition was performed using a Bosch test vehicle equipped with various surround sensors. The recorded data comprises 3D point clouds of a 360 degree LiDAR sensor and images of two front-facing cameras with horizontal opening angles of 45 and 90 degrees, respectively. Furthermore, an IMU with differential GPS provides precise information on the ego-vehicle’s position and motion.

B. Data Labeling

Data post processing included an automatic extraction of pedestrian trajectories. As a first step, we applied a Mask-RCNN object detector to the recorded image data. Pedestrian

detections were subsequently lifted to 3D by matching them with respective LiDAR point clusters. To this end, LiDAR point clouds were ego-motion compensated, projected onto the image plane, and finally matched based on their overlap with pedestrian masks. The resulting 3D detections were tracked by Kalman filtering using a constant velocity model and a greedy association scheme. Finally, the 3D annotations were projected into the ground plane to obtain 2D pedestrian positions.

Overall, 93,162 unique pedestrian tracks with an average track length of 6.6s were extracted from the recorded data, including approximately 5,500 actor trajectories. 7 % of all pedestrian trajectories are crossing the road, out of which approximately 40 % originate from actors. This demonstrates that even though actor trajectories constitute only a small portion of the overall dataset, it was possible to significantly increase scenario coverage by employing actors during the recordings. Fig. 2 illustrates the distribution of pedestrian tracks with respect to the different round courses and crossing locations. For our further investigations we excluded pedestrian tracks at traffic lights because behavior there is primarily determined by the traffic light states.

To enable an investigation of potentially behavior-relevant features, a subset of 9,438 pedestrian tracks was manually labeled. For the labeling, the tracks were randomly selected while ensuring that tracks are balanced between crossing and

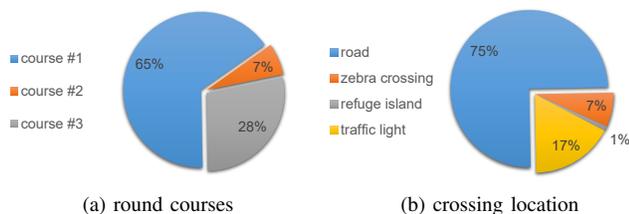


Fig. 2. Dataset statistics: (a) distribution of pedestrian tracks with respect to the three courses; (b) locations at which pedestrians crossed the street.

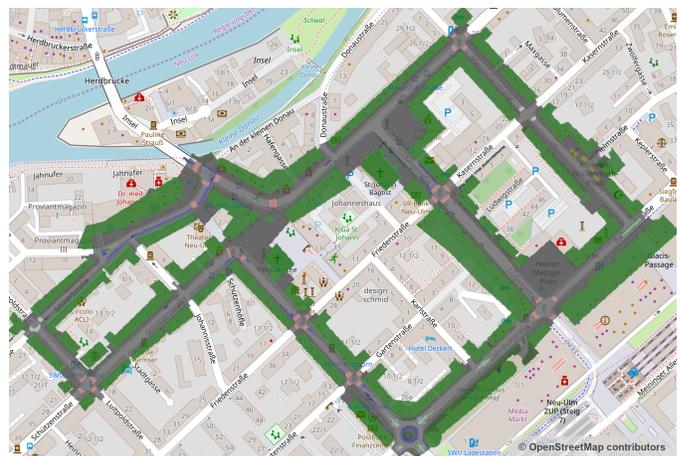


Fig. 4. Semantic map shown on top of a road map. Colors denote different semantic classes. © OpenStreetMap contributors (openstreetmap.org, open-datacommons.org)

non-crossing pedestrians as well as location (i.e. presence of different traffic control elements). As depicted in Fig. 3, the frame-wise labels include body and head orientation (in degrees).

Finally, we created semantic maps of the round courses via hand-labeling of decimeter-level accurate aerial orthoimages (see Fig. 4). These maps comprise the locations of roads, sidewalks, cycle tracks, bus lanes, barred areas, zebra crossings, refuge islands, traffic lights, lawn, and buildings. The trajectories of the ego-vehicle as well as those of the detected pedestrians hence can be projected onto the map, given the precise global position recordings of the ego-vehicle.

The large variety of scenarios and rich feature description renders the dataset suitable for pedestrian behavior prediction. However, the recordings only comprise German traffic. Thus, the generalization of trained prediction models and our results may be limited to regions with similar cultural factors or traffic rules.

IV. REQUIREMENTS ON PEDESTRIAN BEHAVIOR PREDICTION

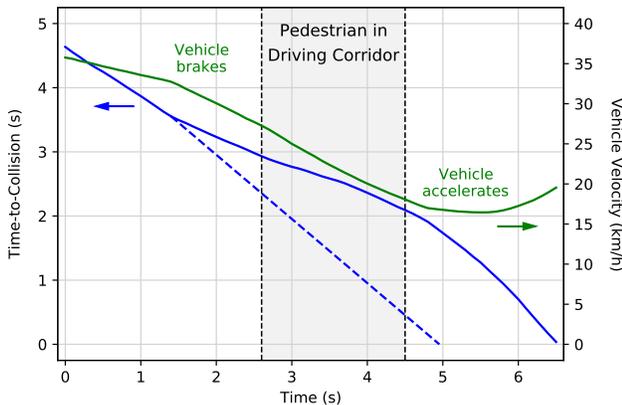
A. Analysis of Human Driving Behavior

The interaction of vehicles and pedestrians has been studied extensively in recent years [51]. While early studies on the crossing behavior of pedestrians date back to the 1950s, the complex interaction process between drivers and pedestrians that e.g. occurs at crosswalks is a relatively new field of research [52]. At roadways without areas distinctly labeled for

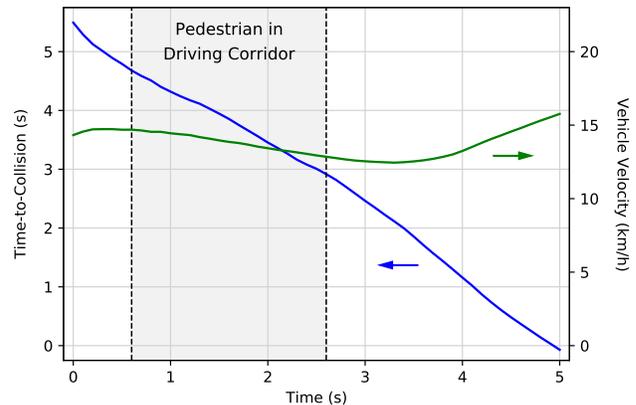
pedestrian crossing, it is the pedestrian’s responsibility to find a safe gap in traffic. Generally, the minimum gap which is still accepted by pedestrians when they decide to cross the road in front of an approaching vehicle is denoted as acceptance gap. A common measure for how safe a specific gap is, is the Time-to-Collision $TTC = d/v$. It is calculated from the distance d of an approaching vehicle, in relation to its driving velocity v and thus indicates the time required for the vehicle to arrive at the pedestrians location assuming constant velocity. Although common sense might suggest that TTC is the basis for pedestrians’ gap selection it has been shown that other factors such as vehicle speed [53] or crossing distance [54] have an influence on the size of chosen gaps.

When interacting with pedestrians it is crucial that the system behavior of an automated vehicle is perceived neither uncomfortable nor even critical by both the pedestrian and the passengers of the vehicle. In that respect, an automated vehicle should ideally imitate the behavior of a defensive human driver [55]. In order to define a suitable system behavior of an automated vehicle, we investigated different scenarios where a human driver interacts with pedestrians crossing the roadway in front of the vehicle. These scenarios were taken from the dataset described in Section III. Two typical examples are depicted in Fig. 5.

The graph in Fig. 5a shows the velocity (green) and the TTC (blue) of a vehicle approaching a pedestrian who crosses the road from left to right. The two vertical dashed lines at $t_1 = 2.6$ s and $t_2 = 4.5$ s mark the points in time at which the



(a)



(b)

Fig. 5. Two traffic scenarios of a vehicle approaching a pedestrian crossing the roadway. In scenario (a) the driver brakes in order to maintain a Time-to-Collision (TTC) above approximately 2 s while the pedestrian is traversing the driving corridor of the vehicle. In scenario (b) the driver maintains a sufficiently large time gap between his vehicle and the pedestrian without having to brake.

pedestrian enters and leaves the driving corridor of the vehicle, respectively. For the analysis in this paper we assume a driving corridor width of 3 m. Overall it takes the pedestrian 1.9 s to traverse the danger zone of the driving corridor. At time t_1 , when they enter the corridor, the TTC of the approaching vehicle is 2.9 s and at time t_2 , when the pedestrian finally leaves the corridor, TTC is 2.1 s.

As soon as the driver recognizes that the pedestrian intends to cross the roadway, they slightly start to brake their vehicle at $t = 1.5$ s. This can be clearly seen by the decrease in slope of the vehicle velocity and by the increase in slope of the TTC curves, respectively. By slightly braking, the driver maintains a TTC above approximately 2 s while the pedestrian is traversing the driving corridor. As soon as the pedestrian has left the corridor and is thus out of the danger zone, the driver releases the brake at $t = 4.8$ s and shortly afterwards accelerates the vehicle.

In contrast, the dashed blue line in the graph indicates how the TTC would evolve over time if the driver did not brake and maintained a constant velocity. In this case, the vehicle would approach the pedestrian faster leading to smaller TTC values while the pedestrian is traversing the driving corridor. Here, the TTC would drop to a value of 0.5 s at t_2 when the pedestrian leaves the driving corridor.

Fig. 5b depicts another traffic scenario with a vehicle approaching a pedestrian crossing the roadway from right to left. Again, the dashed vertical lines indicate when the pedestrian enters and leaves the driving corridor, respectively. In contrast to the previous case, the driver neither brakes nor accelerates while their vehicle is approaching the pedestrian. This is due to the fact that the time gap already is at a safe level > 2.9 s while the pedestrian is traversing the driving corridor and, therefore, the driver does not have to take any action but rather maintains a constant velocity of approximately 14 km/h.

These two examples suggest that human drivers try to establish a time gap between the pedestrian and their vehicle that does not fall below a certain threshold value for the time span while a pedestrian traverses the driving corridor. This hypothesis is also confirmed by a statistical analysis of the minimum time gaps which occur in crossing scenarios. To this end, 2,238 scenarios of pedestrians crossing the roadway from right to left and vice versa in front of an approaching vehicle were identified in the dataset. For these scenarios, we determined the minimum time gaps occurring while pedestrians were traversing the vehicle corridor.

Fig. 6 shows the cumulative distribution function (red) and the number of occurrences (blue) of these scenarios as a function of the minimum time gap. The onset of individual scenarios occurs at a minimum time gap of 0.6 s and the distribution reaches the maximum number of occurrences in the right-open interval between 2.3 s and 2.4 s. The median of the distribution is at a minimum time gap of 2.84 s, the first and third quartiles are at 2.09 s and 3.92 s, respectively.

The distribution function of minimum time gaps thus shows that drivers try not to fall below a certain threshold value for the time gap between their vehicle and the pedestrian crossing, which is perceived as safe and comfortable by both parties. If the initial situation of a crossing scenario leads to a time gap

lower than the threshold value, the driver reacts by adjusting the speed of their vehicle accordingly, e.g. by applying the brakes. Furthermore, the result shows that a minimum time gap of at least 2.8 s (i.e. the median of the distribution) is perceived as safe and comfortable by the majority of traffic participants.

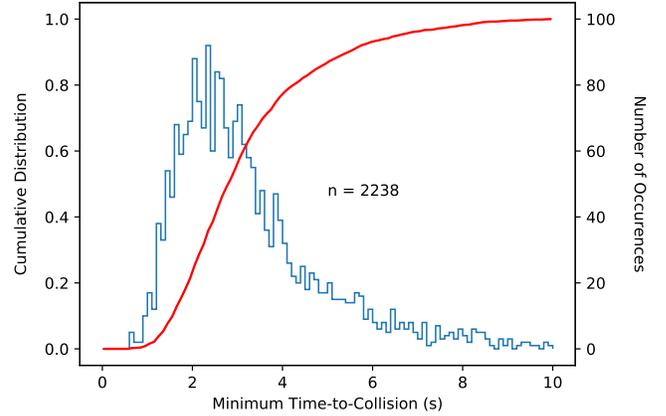


Fig. 6. Cumulative distribution function (red) and number of occurrences (blue) of the minimum time gap while pedestrians traverse the driving corridor of an approaching vehicle. The median of the distribution is at a minimum time gap of 2.84 s.

As stated earlier, the acceptance gap is defined as the minimum gap still accepted by pedestrians when they decide to cross the road in front of an approaching vehicle. While it has been widely used to characterize the crossing behavior of pedestrians it is not an appropriate measure to describe the interaction between drivers and pedestrians as they cross the roadway. This is due to the fact that the acceptance gap solely depends on the pedestrians' decision to cross in a specific traffic situation and does not account for the drivers' reaction. In contrast, our results show that the minimum time gap of pedestrians while they traverse the driving corridor can be used to analyze the interaction of both parties as it combines the crossing intent of the pedestrian given a specific traffic situation and the reaction of the driver to the crossing pedestrian.

Still both metrics are related to each other: While the acceptance time gap corresponds to the TTC when a pedestrian enters the roadway, the minimum time gap usually refers to the TTC when a pedestrian is about to leave the driving corridor of the approaching vehicle. Consequently, the minimum time gap takes on smaller values than the acceptance time gap. The difference between both metrics thereby depends on the time span it takes the pedestrian to traverse the driving corridor and the reaction of the driver of the approaching vehicle.

B. Derived AD System Reaction Pattern

Based on our analysis of human driving behavior we next derive system reactions of an automated vehicle that should imitate those of human drivers. For the sake of simplicity we focus on longitudinal vehicle control, i.e. the adaption of the vehicle's velocity along a predefined or planned path.

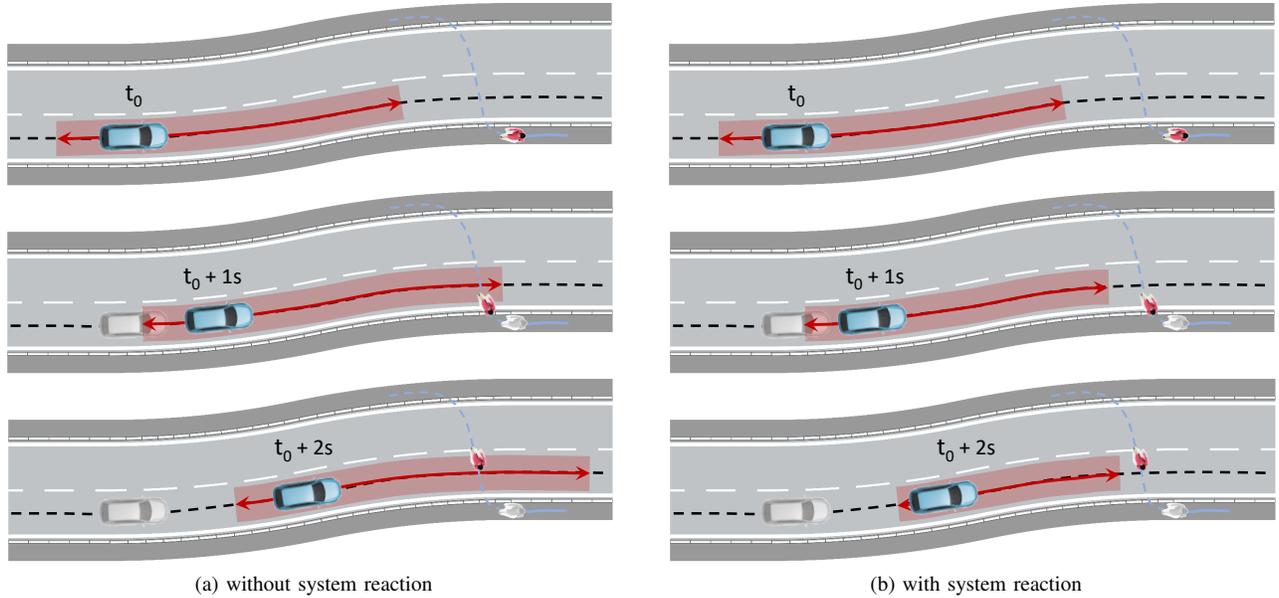


Fig. 7. Defining an appropriate system reaction based on the concept of a driver's comfort zone (red regions). In (a) the system first evaluates whether a pedestrian will violate the comfort zone in the future. In (b) a braking maneuver finally slows down the vehicle such that future comfort zone violations are circumvented. Dashed lines depict the paths of the vehicle and the pedestrian.

As shown in the previous section, the observed minimum time gap between a vehicle and a pedestrian traversing the driving corridor is of particular importance in this respect. Our analysis suggests, that there is a threshold value below which a situation is perceived uncomfortable or even critical. Acceptance of minimum time gaps is, of course, subjective and may depend on driving style, street layout, or traffic density. However, the histogram in Fig. 6 suggests that time gaps below 2s (i.e. approximately the first quartile of the distribution) are seldom observed and hence also should be avoided by an automated vehicle.

To realize an adequate system behavior we suggest that an automated vehicle monitors its future driving corridor, in particular with respect to violations of a defined minimum time gap. For this purpose, we define a comfort zone that extends along the vehicle's future path, as illustrated in Fig. 7. The extent of the comfort zone thereby reflects a time gap which is considered comfortable by drivers and pedestrians. Consequently, its length depends on the vehicle's speed and can thus be adapted via braking (zone shrinks) or acceleration (zone is enlarged). In our implementation we choose a time gap of 3s, which is slightly above the median value of the distribution in Fig. 6 and therefore corresponds to an average driving style.

Furthermore, we propose that an automated vehicle evaluates whether its comfort zone is violated by a pedestrian – currently or in the future. This is done by predicting the probability distribution of future pedestrian locations and intersecting them with the vehicle's future comfort zones (see Fig. 7a). The latter are derived by shifting the current comfort zone along the planned path assuming constant vehicle velocity (i.e. no change in system behavior). In case of violations, the automated vehicle issues a system reaction, e.g. braking,

such that the pedestrian's future trajectory does not violate the adapted comfort zones anymore (see Fig. 7b).

From the above considerations it becomes evident that comfortable and anticipatory driving requires large prediction horizons. Pedestrians at least have to be predicted for a time horizon that corresponds to the comfort zone's time gap, i.e. 3s in our case. To realize natural driving behavior even larger horizons are required such that future comfort zone violations can be anticipated. However, with increasing prediction horizons the requirement on prediction accuracy can be gradually relaxed, since then an increasing time span is available to correct for an erroneous absence of a system reaction. On the other hand, to minimize false system reactions it is required that behavior planning takes prediction uncertainties into account. Pedestrian prediction models hence should yield uncertainty estimates (e.g. in terms of probability distributions over future pedestrian locations).

It should be noted that the comfort zone is not necessarily restricted to regions in front of the vehicle. Rather, it may also cover a region behind the vehicle. This allows to evaluate whether the vehicle can pass a pedestrian with a sufficient time gap before the pedestrian enters the driving corridor. Similarly, the comfort zone may be extended to include infrastructure elements (e.g. zebra crossings) such that country-specific traffic rules are taken into account.

C. Prediction Performance Metric: In-ROI Sensitivity (IRS)

Building on the AD system reaction, we now derive an application specific performance metric for our prediction model. To this end we define pedestrian behavior prediction as a binary classification task.

Fig. 8a shows a typical traffic scene with the trajectory of a crossing pedestrian as well as the ego vehicle's comfort zone.

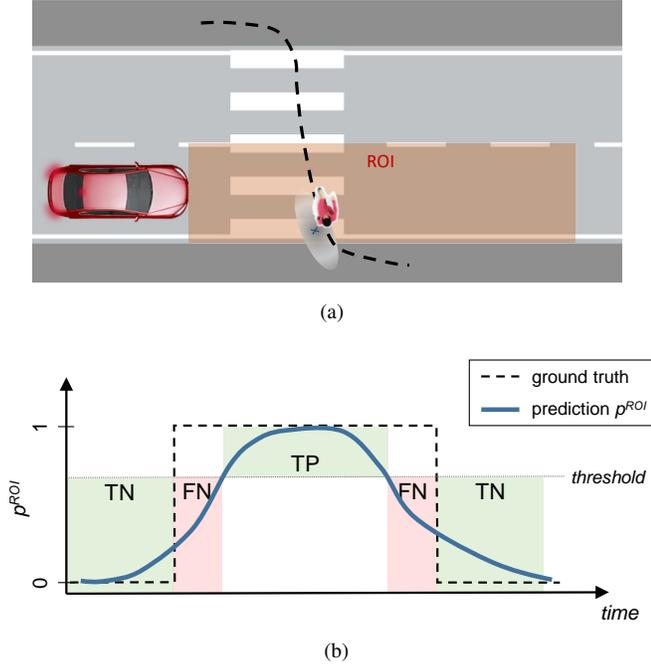


Fig. 8. Derivation of the proposed performance metric. (a) Ground truth pedestrian trajectory crossing the region of interest. (b) Prediction of in-ROI probability (solid) and ground truth of in-ROI state (dashed). For one recorded traffic scene, 3 s predictions of pedestrian and of ego vehicle ROI have been started from all possible points in time (x-axis).

The latter is referred to as Region of Interest (ROI) in the following. The task of the AD system is to anticipate predicted violations of the ROI by pedestrians, i.e. a classification whether a pedestrian will be located inside the ROI in the future. We define the in-ROI probability P_{t+T}^{ROI} at time $t + T$ as

$$P_{t+T}^{\text{ROI}} = \int_{x_{t+T} \in \text{ROI}_{t+T}} p(x_{t+T} | x_{0:t}, C) dx_{t+T}, \quad (1)$$

where ROI_{t+T} denotes the predicted ROI and $p(x_{t+T} | x_{0:t}, C)$ the predictive distribution of pedestrian locations given the pedestrian's past trajectory $x_{0:t}$ and contextual cues C . Computing the in-ROI probability requires integration of the predictive distribution over the ROI. This can be approximated e.g. with a Monte-Carlo approach using samples from the distribution.

Fig. 8b illustrates the evolution of the scene and the prediction over time. Specifically, the dashed line represents the true state of the pedestrian with regard to being inside the ROI for a given prediction horizon, e.g. $T = 3$ s, whereas the solid blue line shows the predicted in-ROI probability for this prediction horizon. Thresholding P_{t+T}^{ROI} finally yields an in-ROI classification that is compared to the true state for each sample t .

Thus, we have defined a classification problem (with predicted class probabilities) for which textbook metrics can be applied. In particular, we choose the True Positive Rate (TPR) and the False Positive Rate (FPR) which are defined as

$$\text{TPR} = \text{TP} / (\text{TP} + \text{FN}) \quad (2)$$

$$\text{FPR} = \text{FP} / (\text{FP} + \text{TN}), \quad (3)$$

where TP, TN, FP, and FN denote the number of true positive, true negative, false positive, and false negative samples, respectively. For metric calculation, we accumulate classification results from the set of test samples that are relevant for the AD system reaction. In our evaluation, we consider samples with $\text{TTC} < 5$ s as relevant, which excludes samples that will not result in a system reaction as defined in section IV-B.

In general, there is a trade off between TPR and FPR that can be represented as a ROC curve, where points on the curve are generated by varying the classification threshold. Fig. 11 shows ROC curves for different prediction horizons. The models and the confidence bands shown in the figure will be discussed in Sec. VI. These curves allow for studying said trade off which would be difficult to assess from requirements given a priori. In particular, we use ROC curves to determine a FPR for each prediction horizon according to the following reasoning.

For an application of prediction models in an AD system, the performance (TPR) at low FPR is of particular interest, since high FPR would result in an unacceptable number of false system reactions. However, acceptable FPRs differ with respect to prediction horizons. Small prediction horizons potentially correspond to critical situations that require a rather strong reaction by the automated vehicle. In such situations an erroneous system reaction is highly unacceptable and potentially dangerous. On the other hand, large prediction horizons correspond to situations that are still uncritical but which are relevant for anticipatory driving. These situations require much weaker reactions by the automated vehicle and therefore erroneous system reactions are more acceptable. Hence, we allow for larger FPRs for long-term prediction compared to very small FPRs in the short-term prediction case. Specifically, we consider FPRs of 2.5%, 5%, 10%, and 15% as suitable working points for the 1 s, 2 s, 3 s, and 4 s predictions, respectively. Thus, we propose In-ROI Sensitivity (IRS) as a prediction performance metric, which measures the In-ROI TPR at the respective FPRs for the different prediction horizons.

D. Metric Assessment

In the previous section, we derive an application-specific IRS-metric with an intuitive function-level interpretation. We postulate that it is preferable to traditional metrics for pedestrian future prediction, as it focuses on aspects of the prediction, which affect the corresponding AD system reaction. To highlight differences between the proposed metric and traditional ones, we perform a metric assessment based on toy data. The toy example is illustrated in Fig. 9, where background colors specify different ground types. It corresponds to a crossing scene, where a pedestrian (black line) is walking on a sidewalk (dark gray) close to a street (bright gray) towards a zebra crossing (white). We visualize distributions with violet contour plots, where Fig. 9a corresponds to the three-modal ground truth distribution, Fig. 9b illustrates a model with inaccurate on-sidewalk prediction, Fig. 9c shows a model with inaccurate crossing prediction, while Fig. 9d does not capture multiple modes at all. The captions contain corresponding

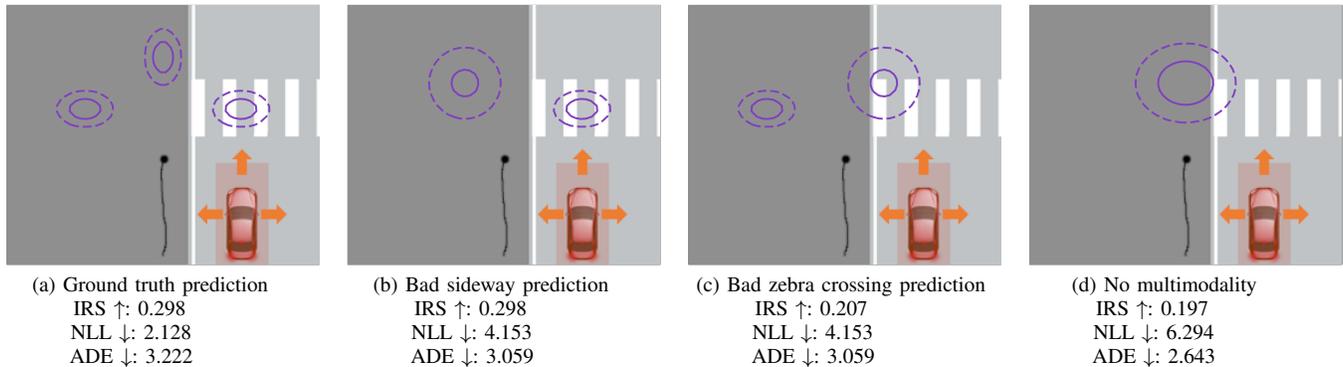


Fig. 9. Toy example for metric assessment: The toy example highlights differences between our metrics in a crossing scene. Background colors specify different ground types, where a pedestrian (black line) is walking on a sidewalk (dark gray) close to a street (bright gray) towards a zebra crossing (violet). (a) corresponds to a three-modal ground truth distribution, (b) represents a model with inaccurate on-sidewalk prediction, (c) shows a model with inaccurate crossing prediction, while (d) does not capture multiple modes at all. The captions contain corresponding metrics when comparing the predictions (a) - (d) to the ground truth distribution (a). For computing IRS, we have varied the positions of car ROIs on the street.

metrics when comparing the prediction to the ground truth distribution in Fig. 9a. For computing IRS, we have varied the positions of ego vehicle ROIs on the street.

We observe that matching the ground truth distribution yields highest NLL, while wrong predictions always lead to worse scores. Interestingly, ADE results in contradictory findings, as it favors uni-modal distributions due to measuring expected Euclidean error. Finally, according to our IRS metric, two predictions lead to similarly good results: Fig. 9a, which matches the ground truth distribution correctly, and Fig. 9b, which only matches the relevant road / zebra mode correctly. Consequently, the IRS metric focuses on relevant aspects of the prediction, while not evaluating irrelevant parts (e.g. accurate prediction of pedestrians on sidewalks). This helps in finding the right balance between model complexity and predictive performance.

V. PEDESTRIAN PREDICTION MODEL

In this section, we outline a simple yet powerful enough pedestrian prediction model for assessing our proposed metric and the relevance of different feature combinations. Section IV analyzes human driving behavior and derives an expected AD system reaction pattern. Such a reaction pattern requires judging whether the probability of a pedestrian being inside a future comfort zone of the ego vehicle exceeds a certain threshold. Based on this reaction pattern, we can derive desired properties of prediction models:

- *Prediction of a continuous distribution over future pedestrian locations* for computing the likelihood of a pedestrian being inside the future comfort zone of the ego vehicle.
- *Predictive distributions over different prediction horizons* to evaluate comfort zone violations for different prediction horizons.
- *Ability to model complex multi-modal distributions* as future behavior can have several non-trivial modes (e.g. going straight or crossing the street).
- *Learning influences of contextual cues* that affect pedestrian behavior. Concretely, we focus on pedestrian motion

and poses, the static map, and interactions with the ego vehicle, as we expect that these features have the strongest influence on the AD system reaction.

Due to the limited prediction horizon and only a subset of pedestrians being relevant for the IRS metric, we believe that single agent prediction models that condition on static and dynamic contextual cues are sufficient for our purpose.

A. Model and Feature Integration

In Section II, we refer to several recent approaches for pedestrian behavior prediction, which differ in terms of model types, architectures, learning method, or availability of features. However, most recent approaches use deep learning based models for learning complex, non-linear functional dependencies from input features to the predicted behavior. Especially Conditional Variational Autoencoders (CVAEs) [25] and Conditional Generative Adversarial Networks (CGANs) [56] have proven to be suitable for such use cases as they can learn complex, continuous distributions by introducing continuous latent variables. Furthermore, they allow to incorporate feature observations by conditioning on them. We specifically focus on CVAEs due to two reasons: their explicit density modeling, which allows to evaluate the NLL, as well as not suffering from mode collapse, which ensures to capture non-trivial modes in real human behavior.

In the following, we give an overview on the proposed CVAE architecture, with a focus on simplicity and flexibility for enabling feature relevance assessment with different sets of features under the proposed metric. As indicated in Fig. 10 the model conditions on both static map information at a particular time step s_t as well as dynamic features of the pedestrian $x_{t-H+1:t}$ over the last H time steps. We encode dynamic features of the agent $x_{t-H+1:t}$ (e.g. relative motion between two time steps, head pose, body pose, distance to the ego vehicle) via recurrent encoders into an embedding space. In addition, we represent the static environment around the pedestrian via a grid in bird's-eye view, where different colors indicate different semantics (e.g. sidewalk, road, zebra

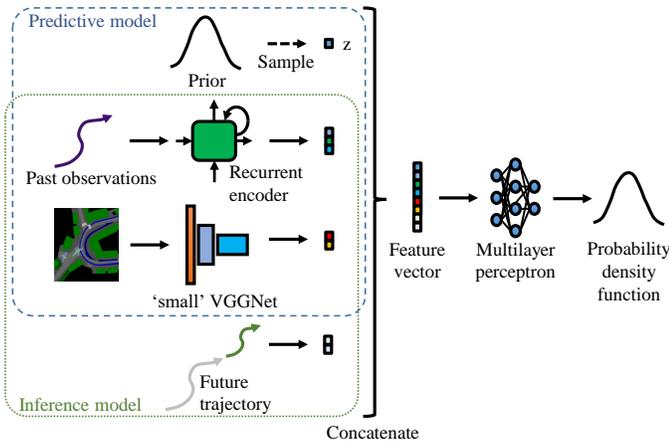


Fig. 10. DL-based architecture of the pedestrian prediction CVAE. For inferring the latent variable, the inference model uses encoded past observations of a pedestrian, encoded static context of the environment, as well as the observed future trajectory. In contrast, the predictive model does not have access to the future trajectory, but predicts future locations given past observation, static context, and a sample of the latent variable z . The predictive and the inference model share the same encoders of past observations and static context.

crossing, building, isle, bicycle lane, unknown). These bird’s-eye view grids are encoded via a small VGGNet architecture [57] to provide a static environment embedding vector.

A CVAE comprises two models: The predictive model $p_\theta(x_{t+1:T}|x_{t-H+1:t}, s_t, z)$ (dashed blue rectangle) predicts an 8-dimensional Gaussian distribution over the future pedestrian locations in x and y at four prediction time steps (1s, 2s, 3s, 4s) conditioned on past observations and a latent variable sample z from the Z -dimensional multivariate standard Gaussian prior $p(z)$. The inference model $q_\phi(z|x_{t+1:T}, x_{t-H+1:t}, s_t)$ (dotted green rectangle) infers the latent variable z from all available observations (future and past) and is only used during training. Additional implementation details can be found in the Appendix in Sec. A.

We train the model in the usual way by minimizing the evidence lower bound (ELBO):

$$\mathbb{E}_{q_\phi(z|x_{t-H:T}, s_t)} [\log p_\theta(x_{t+1:T}|x_{t-H:t}, s_t, z)] - D_{KL}(q_\phi||p(z)) \quad (4)$$

The ELBO comprises a reconstruction term that aims to maximize the expected likelihood of future observations under the latent posterior distribution, as well as a Kullback-Leibler (KL) divergence term that tries to enhance agreement between prior and posterior latent distribution. So while the reconstruction term serves the purpose of obtaining an accurate prediction of the future trajectory given the observed past, the KL term acts as a regularizer.

B. Feature combinations:

The proposed architecture allows for flexible changes of input features and enables an ablation study regarding the influence of different features on the prediction performance. Table II denotes the features that we evaluate in this study. The most basic model only conditions on the past motion of the pedestrian and does not use any additional information.

Hence, it can only learn future predictions based on cues in the pedestrian trajectory itself.

TABLE II
DESCRIPTION OF FEATURES USED IN THE CVAE.

Feature	Description
Motion	Trajectory encoded as relative motion $(\Delta x, \Delta y)$ between time steps.
Egodist	Distance (x, y) to ego vehicle at every time step.
Head	Head pose of the pedestrian at every time step.
Body	Body pose of the pedestrian at every time step.
Map	Semantic map around the last position of the pedestrian.

VI. EXPERIMENTAL RESULTS

In the following, we evaluate the importance of different sets of contextual cues in terms of our proposed IRS metric. Furthermore, we highlight differences in the respective conclusions when comparing to traditional metrics.

A. Experimental Setup

Training, validation, and test subsets are drawn stratified from the three round courses in our dataset. The dataset is split by first assembling the test set. In order to be able to evaluate how pedestrian features like head pose affect the prediction performance, the whole test set needs to consist of labeled tracks. With the fraction of labeled tracks in the dataset being comparatively small, the test set has to be limited in size to ensure enough labeled tracks are available for training. We decided to constrain the test set size to 500 tracks while at the same time guaranteeing a large variety of scenes. This is achieved by drawing from the labeled subset via stratified random sampling, with each category, e.g. crossing/not crossing, distance from ego vehicle, crossing in front/behind the ego vehicle, being represented according to its proportion in the complete data set.

For training and validation, a minimum track length of 5s is required to enable training of prediction horizons of up to 4s. 5% of the tracks fulfilling this criterion are randomly assigned to the validation set, the remaining tracks form the training set. Training and validation set contain 42,551 (7,175 labeled) and 2,278 (389 labeled), respectively.

In order to assess the variance of the test results, evaluations are repeated using the bootstrap method [58] with $B = 10000$ replications. Each replication uses an artificial test set created from the original test set by resampling 500 tracks with replacement. The bootstrap method produces an estimate of the variability in the results that would be seen with completely new test data taken from the same ground truth distribution. From the bootstrap replications we compute 50% confidence intervals for the metrics IRS, NLL, and ADE, using the ‘‘BCa’’ method [59], which corrects for bias and skewness of the sampling distribution.

B. Quantitative Evaluation

In order to introduce the procedure and to gain some intuition we first show evaluation results for two models in comparison. One is a CVAE model that solely uses a

pedestrian’s past trajectory (*CVAE motion*), while the other one is a CVAE model that employs a full feature set (*CVAE full*) consisting of the pedestrian’s past trajectory, their head and body pose, the distance to the ego vehicle’s position, and the semantic map. In this way the effect of contextual cues on the prediction performance can be assessed. We run both models on the test dataset and evaluate the predictions with the system-specific ROI-metric proposed in Sec. IV-C, using a range of FPR values.

The ROC curves for different prediction horizons are depicted in Fig. 11 (solid line: TPR per FPR, shaded area: 50% confidence interval per FPR). It becomes evident that the TPRs

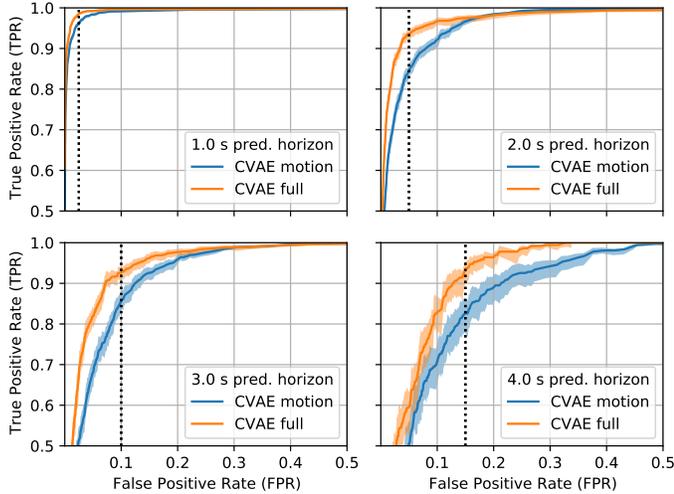


Fig. 11. ROI-based metric results of the baseline CVAE motion model (pedestrian motion feature only) and of the full CVAE model (all features) for different prediction horizons. The shaded bands denote the TPR confidence interval at each FPR. Dashed lines represent the target FPR values chosen to define the IRS metric, see Sec. IV-C.

of the models decrease with increasing prediction horizons. For a prediction horizon of 1 s both models achieve very good TPRs with the full CVAE model slightly outperforming the simpler one. This result confirms that for short-term pedestrian prediction the simpler model is already well suited in terms of our IRS metric, and that contextual cues only have a minor benefit. In contrast, for prediction horizons of 2 to 4 s, the full CVAE model strongly outperforms the CVAE motion model. Our results thus confirm the importance of using contextual cues for long-term pedestrian prediction. A more detailed analysis of feature relevance follows in Sec. VI-C.

With longer prediction horizons the confidence bands become wider. This results from the fact that there is less data available for longer prediction horizons in our data set. We visualize the chosen FPRs from Sec. IV as vertical dotted lines in Figure 11. Our prediction performance metric IRS is equal to the TPR at the chosen FPRs. The IRS as well as the corresponding 50% confidence intervals are reported in Table III for eight CVAE models with different input feature combinations.

C. Feature Relevance Assessment

To analyze the relevance of different contextual cues, we perform an ablation study and report the results in Table III.

We list our proposed In-ROI Sensitivity (IRS), the Negative Log-Likelihood (NLL), and the Average Displacement Error (ADE) for four prediction horizons and eight variants of the CVAE model, each using different input feature combinations. In order to keep the complexity of this analysis manageable, we treat the head and body orientation as one combined feature named head&body by concatenating the orientations. Along with the metrics itself Table III reports 50% confidence intervals for all values in order to show the uncertainty of the metric estimation. This allows to better interpret the amount to which one model improves about another, but these intervals are not suited for concluding whether the models perform significantly different. For that, one actually has to look at, e.g., 90% confidence intervals for pairwise differences of metrics. We do not report such intervals due to space considerations but they have been computed and used to check the statistical significance of the statements made in the following.

Based on our pairwise significance analysis, we additionally report in Table IV for each metric and prediction horizon the optimal CVAE model (input feature set). A model is deemed optimal if there is no other model with a significantly better performance. In cases of ties, we report the model with the least amount of input features, i.e. the least complex model.

Using our proposed IRS metric, we first analyze the relevance of different contextual cues. In comparison to other ablation studies, this analysis thus focuses on the importance of different input feature combinations with respect to the overall system performance. By taking the specific requirements of the system function into account, one can ensure that the prediction model is perfectly tailored towards a downstream task (e.g. AEB-P). As shown in Sec. IV-D, the results of general prediction metrics (e.g. NLL) and the IRS metrics may vary. An evaluation based on a system-level metric can consequently avoid the use of unnecessary complex models.

It is obvious that the map is crucial for good prediction performance, when comparing models 1-4 (without map) to models 5-8 (with map). The positive effect of using a map is particularly pronounced for prediction horizons of 2 s to 4 s. However, our further significance analysis indicates three deviations from this general trend. Enhancing model CVAE motion+egodist with map does not yield a significant performance improvement for 3 s and 4 s predictions. Additionally, no significant improvement for a prediction horizon of 1 s can be observed when comparing model 3 to 7 and 4 to 8.

Overall the IRS metric can be significantly increased from 96.1% to 98.5% for short prediction horizons (1 s) and from 82.6% to 93.3% for long prediction horizons (4 s) by extending the conditioning of the CVAE motion model to all available contextual cues. The IRS of model CVAE motion+map+head&body is with 93.4% even a little better, but this improvement is not significant.

As can be seen in Table IV-D, an additional feature does not always yield a significant performance improvement when using our system-level IRS metric based on our selection criterion, dataset, and model family. Short-term predictions (1 s, 2 s) benefit from considering egodist and head&body features. For long-term predictions (2 s, 3 s, and 4 s), it is advantageous to use map features. Comparing to CVAE motion+map, we

TABLE III
ABLATION STUDY OF CONTEXTUAL CUES

CVAE Model		IRS (%) \uparrow				NLL \downarrow				ADE (m) \downarrow			
		1 s	2 s	3 s	4 s	1 s	2 s	3 s	4 s	1 s	2 s	3 s	4 s
1	motion	96.5	86.1	87.8	86.3	0.18	1.59	2.47	3.21	0.54	1.12	1.77	2.54
		96.1	84.6	85.4	82.6	0.13	1.54	2.39	3.10	0.53	1.10	1.73	2.48
		95.5	82.8	82.6	77.4	0.08	1.46	2.28	2.92	0.52	1.08	1.70	2.41
2	motion+egodist	97.8	89.6	91.1	85.3	-0.10	1.33	2.22	2.99	0.50	1.05	1.68	2.45
		97.4	88.0	89.0	82.2	-0.14	1.27	2.14	2.88	0.49	1.03	1.64	2.39
		96.9	86.4	86.9	76.8	-0.20	1.19	2.02	2.69	0.48	1.01	1.61	2.32
3	motion+head&body	98.2	91.3	91.0	88.2	0.09	1.50	2.39	3.13	0.52	1.11	1.77	2.55
		97.7	89.7	89.1	84.6	0.05	1.45	2.32	3.04	0.51	1.09	1.74	2.49
		97.3	88.2	87.1	80.1	0.00	1.39	2.22	2.88	0.50	1.07	1.70	2.43
4	motion+egodist+head&body	98.4	92.2	92.1	89.7	-0.13	1.29	2.17	2.91	0.48	1.02	1.62	2.35
		98.1	91.0	90.8	86.8	-0.18	1.23	2.09	2.81	0.47	0.99	1.58	2.29
		97.7	89.4	88.5	82.2	-0.23	1.15	1.99	2.62	0.46	0.97	1.55	2.23
5	motion+map	98.3	92.1	92.6	93.7	-0.02	1.30	2.11	2.81	0.52	1.04	1.59	2.25
		97.9	90.7	90.6	91.2	-0.06	1.25	2.04	2.71	0.50	1.01	1.55	2.19
		97.4	89.1	88.8	88.1	-0.11	1.18	1.94	2.57	0.49	0.99	1.52	2.13
6	motion+map+egodist	98.4	93.3	92.5	89.8	-0.19	1.16	2.00	2.72	0.48	0.98	1.53	2.19
		98.0	92.2	90.4	87.2	-0.24	1.10	1.91	2.61	0.47	0.96	1.50	2.13
		97.6	90.5	88.5	82.3	-0.29	1.02	1.79	2.42	0.46	0.94	1.46	2.06
7	motion+map+head&body	98.3	93.4	94.1	95.8	-0.07	1.25	2.05	2.77	0.49	1.00	1.54	2.22
		98.0	92.2	92.6	93.4	-0.11	1.19	1.98	2.69	0.48	0.97	1.50	2.16
		97.6	90.6	91.1	90.7	-0.15	1.13	1.89	2.56	0.47	0.95	1.47	2.10
8	motion+map+egodist+head&body	98.9	94.5	94.0	95.3	-0.21	1.11	1.94	2.65	0.48	0.97	1.51	2.18
		98.5	93.5	92.7	93.3	-0.26	1.05	1.86	2.55	0.47	0.95	1.47	2.12
		98.2	92.1	90.8	90.0	-0.31	0.98	1.76	2.37	0.46	0.93	1.44	2.05

TABLE IV
OPTIMAL INPUT FEATURE SET PER METRIC AND TIME HORIZON

	IRS	NLL	ADE
1 s	motion+egodist+head&body	motion+map+egodist	motion+map+egodist
2 s	motion+map+egodist+head&body	motion+map+egodist+head&body	motion+map+egodist+head&body
3 s	motion+map	motion+map+egodist+head&body	motion+map+egodist+head&body
4 s	motion+map	motion+map+egodist+head&body	motion+map+egodist

did not observe significant improvements for 3s, and 4s by adding egodist or head&body. However, this could be due to limited head and body orientation labels in the dataset.

In a second step, we compare the IRS score of the models with the corresponding ADE and NLL scores. Based on a toy example, we have shown in Sec. IV-D that improvements in NLL do not necessarily pay off in the proposed application-specific IRS-metric. Indeed, we also make similar observations in the ablation study results of Table III. For example, model 8 significantly improves over model 7 for the 4s prediction horizon in terms of NLL, while this is not the case in terms of IRS. The difference are also visible in Table IV and especially for a prediction horizon of 3s and 4s. According to our selection criterion, NLL and ADE suggest to use all available input features, our IRS metric, on the other hand, indicates that features motion and map are sufficient.

D. Qualitative Evaluation

The system-level feature relevance assessment in Section VI-C indicates that CVAE-based models with additional features can significantly outperform a simple CVAE motion

baseline model. In the following, we provide qualitative examples, which highlight the influence of features on pedestrian prediction. Fig. 12 illustrates two scenes that often occur in urban scenarios. In the first scene, Fig. 12a and 12b, a pedestrian is walking on the sidewalk (dark gray) parallel to the street (light gray) and approaching a zebra crossing (violet). The CVAE motion model in Fig. 12a does not have information about the static environment and thus mainly predicts straight walking with some uncertainty. Instead, the CVAE with map feature in Fig. 12b was able to learn that zebra crossings increase likelihoods of pedestrians to cross and it correctly skews the distribution towards the zebra crossing, while still keeping a mode for straight walking. The second scene, Fig. 12c and 12d, shows a pedestrian that stands at the side of a street and waits for crossing, while a car is decelerating to let the pedestrian pass. The CVAE motion model in Fig. 12c with only pedestrian features is not able to predict that the likelihood increases for the pedestrian to start walking, while the CVAE model with map and ego vehicle features in Fig. 12d picks up this information very fast and predicts the pedestrian to cross the street.

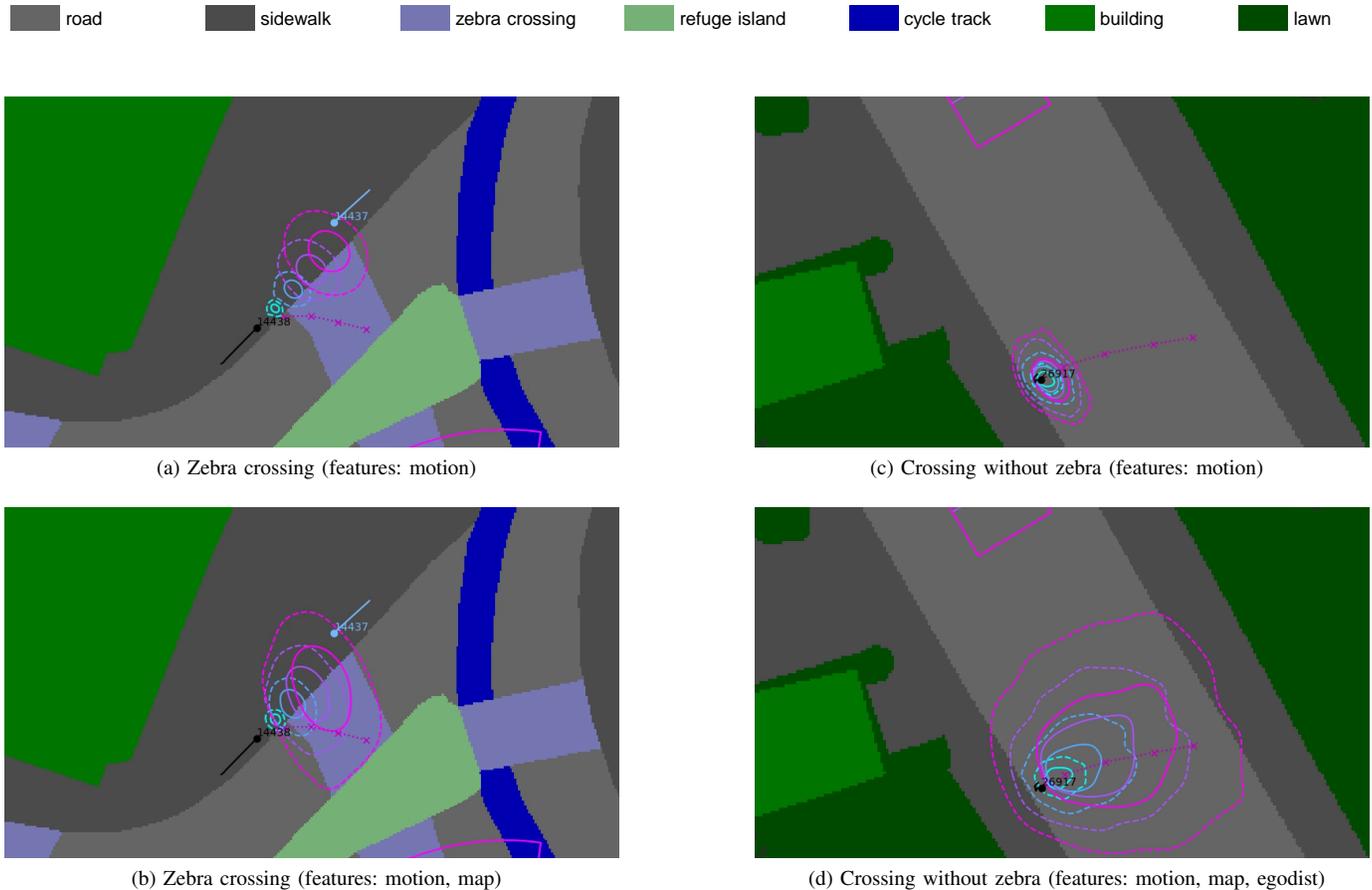


Fig. 12. Pedestrian behavior prediction of the proposed model in two different scenarios. The circles indicate predictions for 1s (turquoise), 2s (blue), 3s (purple), and 4s (magenta), while the dashed line with crosses indicate the ground truth future trajectory. In the first scenario, depicted in (a) and (b), a pedestrian is walking towards a zebra crossing after walking straight on the sidewalk, but parallel to the street. In the second scenario, depicted in (c) and (d), a pedestrian is starting to cross the street without a zebra crossing, after a car has reduced its speed (prediction of the car is illustrated by a magenta box at the top).

VII. CONCLUSION

With the shift from advanced driver assistance systems towards fully automated driving, novel requirements on pedestrian behavior prediction arise. In this paper, we argued that these requirements are not fully taken into account by established evaluation procedures – particularly in terms of metrics and datasets that are usually used to quantitatively assess prediction performance. We proposed a system-level approach to bridge this gap: based on a large dataset comprising thousands of pedestrian-vehicle interactions, we analyzed human driving behavior, derived appropriate reaction patterns of an AD system, and finally specified corresponding requirements on a pedestrian behavior prediction component. Moreover, we proposed a novel evaluation metric that measures the fulfillment of these requirements. It eases interpretation of prediction performance from a system-level perspective and thus allows for balancing model complexity vs. system-level performance. Our contribution thus shall stimulate future research on system-level evaluation and optimization of prediction models.

The proposed metric was evaluated on a large-scale dataset comprising thousands of real-world pedestrian-vehicle interactions using a CVAE-based model. A thorough ablation study shed light on the relative importance of different features. We

demonstrated that considering additional contextual cues does not always yield a significant performance improvement when using a system-level metric (i.e. our IRS metric). We also showed that results of general prediction metrics (e.g. NLL) and system-level metrics differ. Consequently, an evaluation based on a system-level metric can avoid the use of unnecessary complex models, highlighting the importance of a system-level approach to pedestrian behavior prediction.

Future work could extend the ablation study towards additional datasets and contextual cues such as considering the influence of other traffic participants besides the ego vehicle (e.g. other pedestrians or vehicles). Furthermore, the evaluation of additional model types could yield interesting insights, such as analyzing models that make joint predictions for multiple traffic participants at once. This would allow comparing the IRS metric to multi-agent prediction metrics. Finally, our future work will focus on the integration and optimization of the developed prediction component in an AD system.

VIII. ACKNOWLEDGMENTS

This work is a result of the research project @CITY-AF — Automated Cars and Intelligent Traffic in the City: Automated

Driving Functions. The project is supported by the Federal Ministry for Economic Affairs and Energy (BMWi), based on a decision taken by the German Bundestag. The authors are solely responsible for the content of this publication.

APPENDIX A IMPLEMENTATION DETAILS

A. Model Architecture

Recurrent encoder: We compute an embedding of the observed pedestrian trajectory $x_{t-H:t}$ using a recurrent encoder. This encoder consists of two stacked Long Short-Term Memory (LSTM) cells producing a 128 dimensional embedding vector. Both cells use the same state size which is determined for each input feature combination by means of a hyperparameter search. The observation horizon H is set to 10 time steps, which corresponds to one second.

Map encoder: The static environment, represented as a bird's-eye view semantic grid, is encoded via a small Convolutional Neural Network (CNN). The CNN processes grids of size 256×256 pixels, containing the agent's local environment of size $25.6\text{m} \times 25.6\text{m}$ centered around the position of the pedestrian at the last conditioning time step. The architecture of the CNN is defined by the shortcut notation:

$$C'_4-P-C'_8-P-C'_{16}-C'_{16}-P-C'_{32}-C'_{32}-P-C'_{64}-C'_{64}-P-F'_{512}-F_{100},$$

where C_i is a convolutional layer with i filters, a stride of one and a filter size of 3×3 , P a max-pooling layer with non-overlapping 2×2 regions, and F_i a fully connected layer with i output features. A prime marks layers which apply a Rectified Linear Unit (ReLU) nonlinearity.

Feature transformers: The predictive model and the inference model each use a Multilayer Perceptron (MLP) to transform feature vectors to parameters of an n -dimensional Gaussian distribution. The dimensionality n is set to eight for the predictive model and to ten for the inference model. Both MLPs consist of three fully connected layers and utilize ReLU nonlinearities. The number of output features is identical in each layer and derived using a hyperparameter search. Missing features are replaced by a constant value of zero.

B. Model training

We perform a grid search to determine the best model for each input feature combination. The parameters of our hyperparameter search space are listed in Table V. In total, we train 60 models for each input feature combination and pick the best model based on the validation IRS. Models are trained for 3000 epochs using the Adam optimizer [60] with a constant learning rate of 0.001. To augment the training data, we randomly rotate the trajectories and environment maps. We use a batch size of 1024 for the models using the semantic map as an input feature and a batch size of 2048 for models without map. All weights of the model are reparametrized using weight normalization [61].

TABLE V
HYPERPARAMETER SEARCH SPACE

	Models with map	Models without map
State size of LSTM cells	256, 384	256, 384
Features per MLP layer	256, 384, 512	384, 512, 640
Random seeds	1, ..., 10	1, ..., 10

REFERENCES

- [1] *Global status report on road safety 2018*. Geneva: World Health Organization, 2018.
- [2] S. H. Haus, R. Sherony, and H. C. Gabler, "Estimated benefit of automated emergency braking systems for vehicle-pedestrian crashes in the United States," *Traffic Injury Prevention*, vol. 20, no. Suppl. 1, pp. S171–S176, 2019.
- [3] D. Ridel, E. Rehder, M. Lauer, C. Stiller, and D. Wolf, "A literature review on the prediction of pedestrian behavior in urban scenarios," in *21st International Conference on Intelligent Transportation Systems (ITSC)*, 2018, pp. 3105–3112.
- [4] F. Camara, N. Bellotto, S. Cosar, F. Weber, D. Nathanael, M. Althoff, J. Wu, J. Ruenz, A. Dietrich, G. Markkula, A. Schieben, F. Tango, N. Merat, and C. Fox, "Pedestrian models for autonomous driving part II: High-level models of human behavior," *IEEE Transactions on Intelligent Transportation Systems (T-ITS)*, pp. 1–20, 2020.
- [5] A. Rudenko, L. Palmieri, M. Herman, K. M. Kitani, D. M. Gavrila, and K. O. Arras, "Human motion trajectory prediction: A survey," *The International Journal of Robotics Research (IJRR)*, vol. 39, no. 8, pp. 895–935, 2020.
- [6] F. Schneemann and P. Heinemann, "Context-based detection of pedestrian crossing intention for autonomous driving in urban environments," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2016.
- [7] B. Völz, K. Behrendt, H. Mielenz, I. Gilitschenski, R. Siegwart, and J. Nieto, "A data-driven approach for pedestrian intention estimation," in *IEEE 19th International Conference on Intelligent Transportation Systems (ITSC)*, 2016.
- [8] E. Rehder and H. Kloeden, "Goal-directed pedestrian prediction," in *Proc. of the IEEE International Conference on Computer Vision Workshop (ICCVW)*, 2015.
- [9] E. Rehder, F. Wirth, M. Lauer, and C. Stiller, "Pedestrian prediction by planning using deep neural networks," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2018.
- [10] N. Deo and M. M. Trivedi, "Trajectory forecasts in unknown environments conditioned on grid-based plans," *arXiv:2001.00735*, 2021.
- [11] N. Radwan, A. Valada, and W. Burgard, "Multimodal interaction-aware motion prediction for autonomous street crossing," *arXiv:1808.06887*, 2018.
- [12] E. Mazor, A. Averbuch, Y. Bar-Shalom, and J. Dayan, "Interacting multiple model methods in target tracking: A survey," *IEEE Transactions on aerospace and electronic systems*, vol. 34, no. 1, pp. 103–123, 1998.
- [13] T. Gindele, S. Brechtel, and R. Dillmann, "A probabilistic model for estimating driver behaviors and vehicle trajectories in traffic environments," in *Proc. of the IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2010.
- [14] G. Agamennoni, J. I. Nieto, and E. M. Nebot, "Estimation of multivehicle dynamics by considering contextual information," *IEEE Transaction on Robotics (TRO)*, vol. 28, no. 4, pp. 855–870, 2012.
- [15] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, "Social LSTM: Human trajectory prediction in crowded spaces," in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [16] N. Deo and M. M. Trivedi, "Multi-modal trajectory prediction of surrounding vehicles with maneuver based LSTMs," in *Proc. of the IEEE Intelligent Vehicles Symposium (IV)*, 2018, pp. 1179–1184.
- [17] Y. Chai, B. Sapp, M. Bansal, and D. Anguelov, "Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction," in *Conference on Robot Learning (CoRL)*, 2020.
- [18] A. Kuefler, J. Morton, T. Wheeler, and M. Kochenderfer, "Imitating driver behavior with generative adversarial networks," in *Proc. of the IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2017.
- [19] C. Tang and R. R. Salakhutdinov, "Multiple futures prediction," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

- [20] T. Salzmann, B. Ivanovic, P. Chakravarty, and M. Pavone, “Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data,” in *Proc. of the European Conference on Computer Vision (ECCV)*, 2020.
- [21] N. Lee, W. Choi, P. Vernaza, C. B. Choy, P. H. Torr, and M. Chandraker, “DESIRE: Distant future prediction in dynamic scenes with interacting agents,” in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [22] P. Felsen, P. Lucey, and S. Ganguly, “Where will they go? Predicting fine-grained adversarial multi-agent motion using conditional variational autoencoders,” in *Proc. of the European Conference on Computer Vision (ECCV)*, 2018.
- [23] A. Bhattacharyya, M. Hanselmann, M. Fritz, B. Schiele, and C.-N. Strachle, “Conditional flow variational autoencoders for structured sequence prediction,” in *Bayesian Deep Learning Workshop (NeurIPS)*, 2019.
- [24] S. Casas, C. Gulino, S. Suo, K. Luo, R. Liao, and R. Urtasun, “Implicit latent variable model for scene-consistent motion forecasting,” in *Proc. of the European Conference on Computer Vision (ECCV)*, 2020.
- [25] K. Sohn, H. Lee, and X. Yan, “Learning structured output representation using deep conditional generative models,” in *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- [26] J. F. P. Kooij, N. Schneider, F. Flohr, and D. M. Gavrila, “Context-based pedestrian path prediction,” in *Proc. of the European Conference on Computer Vision (ECCV)*, Cham, 2014.
- [27] W.-C. Ma, D.-A. Huang, N. Lee, and K. M. Kitani, “Forecasting interactive dynamics of pedestrians with fictitious play,” in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [28] A. Elnagar and K. Gupta, “Motion prediction of moving objects based on autoregressive model,” *IEEE Transactions on Systems, Man, and Cybernetics (SMC) - Part A: Systems and Humans*, vol. 28, no. 6, pp. 803–810, 1998.
- [29] M. Bennewitz, W. Burgard, G. Cielniak, and S. Thrun, “Learning motion patterns of people for compliant robot motion,” *The International Journal of Robotics Research (IJRR)*, vol. 24, no. 1, pp. 31–48, 2005.
- [30] C. Schöller, V. Aravantinos, F. Lay, and A. Knoll, “What the constant velocity model can teach us about pedestrian motion prediction,” *IEEE Robotics and Automation Letters (RA-L)*, vol. 5, no. 2, pp. 1696–1703, 2020.
- [31] P. Trautman and A. Krause, “Unfreezing the robot: Navigation in dense, interacting crowds,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2010.
- [32] M. Kuderer, H. Kretschmar, C. Sprunk, and W. Burgard, “Feature-based prediction of trajectories for socially compliant navigation,” in *Proc. of Robotics: Science and Systems (RSS)*, Sydney, Australia, 2012.
- [33] A. Robicquet, A. Sadeghian, A. Alahi, and S. Savarese, “Learning social etiquette: Human trajectory prediction in crowded scenes,” in *Proc. of the European Conference on Computer Vision (ECCV)*, 2016.
- [34] N. Schneider and D. M. Gavrila, “Pedestrian path prediction with recursive bayesian filters: A comparative study,” in *German Conference on Pattern Recognition (GCPR)*. Springer, 2013.
- [35] M. Lubner, L. Spinello, J. Silva, and K. O. Arras, “Socially-aware robot navigation: A learning approach,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2012.
- [36] F. Bartoli, G. Lisanti, L. Ballan, and A. Del Bimbo, “Context-aware trajectory prediction,” in *24th International Conference on Pattern Recognition (ICPR)*. IEEE, 2018.
- [37] J. Gao, C. Sun, H. Zhao, Y. Shen, D. Anguelov, C. Li, and C. Schmid, “Vectornet: Encoding HD maps and agent dynamics from vectorized representation,” in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [38] A. Rasouli, I. Kotseruba, and J. K. Tsotsos, “Pedestrian action anticipation using contextual feature fusion in stacked RNNs,” in *Proc. of the British Machine Vision Conference (BMVC)*, 2019.
- [39] L. Theis, A. van den Oord, and M. Bethge, “A note on the evaluation of generative models,” in *Proc. of the International Conference on Learning Representations (ICLR)*, 2016.
- [40] B. Liu, E. Adeli, Z. Cao, K.-H. Lee, A. Shenoi, A. Gaidon, and J. C. Nibbles, “Spatiotemporal relationship reasoning for pedestrian intent prediction,” in *IEEE Robotics and Automation Letters (IEEE RA-L) and International Conference on Robotics and Automation (ICRA)*, 2020.
- [41] A. Rasouli, I. Kotseruba, and J. K. Tsotsos, “Are they going to cross? A benchmark dataset and baseline for pedestrian crosswalk behavior,” in *Proc. of the IEEE International Conference on Computer Vision Workshop (ICCVW)*, 2017.
- [42] A. Rasouli, I. Kotseruba, T. Kunic, and J. Tsotsos, “PIE: A large-scale dataset and models for pedestrian intention estimation and trajectory prediction,” in *Proc. of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [43] S. Malla, B. Dariush, and C. Choi, “TITAN: Future forecast using action priors,” in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [44] W. Zhan, L. Sun, D. Wang, H. Shi, A. Clause, M. Naumann, J. Kümmerle, H. Königshof, C. Stiller, A. de La Fortelle, and M. Tomizuka, “INTERACTION Dataset: An INTERnational, Adversarial and Cooperative moTION dataset in interactive driving scenarios with semantic maps,” *arXiv:1910.03088 [cs, eess]*, 2019.
- [45] J. Bock, R. Krajewski, T. Moers, S. Runde, L. Vater, and L. Eckstein, “The inD dataset: A drone dataset of naturalistic road user trajectories at german intersections,” *arXiv preprint arXiv:1911.07602*, 2019.
- [46] M.-F. Chang, J. W. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan, and J. Hays, “Argoverse: 3D tracking and forecasting with rich maps,” in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [47] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, “nuScenes: A multimodal dataset for autonomous driving,” *arXiv preprint arXiv:1903.11027*, 2019.
- [48] A. Rasouli, T. Yau, P. Lakner, S. Malekmohammadi, M. Rohani, and J. Luo, “PePScenes: A novel dataset and baseline for pedestrian action prediction in 3D,” in *Machine Learning for Automated Driving Workshop (NeurIPS)*, 2020.
- [49] M.-F. Chang, J. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan, and J. Hays, “Argoverse: 3D tracking and forecasting with rich maps,” in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [50] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, “nuScenes: A multimodal dataset for autonomous driving,” in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [51] A. Rasouli and J. K. Tsotsos, “Autonomous vehicles that interact with pedestrians: A survey of theory and practice,” *IEEE Transactions on Intelligent Transportation Systems (T-ITS)*, vol. 21, no. 3, pp. 900–918, 2020.
- [52] F. Schneemann and I. Gohl, “Analyzing driver-pedestrian interaction at crosswalks: A contribution to autonomous driving in urban environments,” in *Proc. of the IEEE Intelligent Vehicles Symposium (IV)*, 2016.
- [53] T. Petzoldt, “On the relationship between pedestrian gap acceptance and time to arrival estimates,” *Accident Analysis & Prevention*, vol. 72, pp. 127–133, 2014.
- [54] J. Zhao, J. O. Malenje, Y. Tang, and Y. Han, “Gap acceptance probability model for pedestrians at unsignalized mid-block crosswalks based on logistic regression,” *Accident Analysis & Prevention*, vol. 129, pp. 76–83, 2019.
- [55] M. Houtenbos, H. M. Jagtman, M. Hagenzieker, P. Wieringa, and A. Hale, “Understanding road users’ expectations: An essential step for ADAS development,” *European Journal of Transport and Infrastructure Research*, vol. 5, no. 4, pp. 253–266, 2005.
- [56] M. Mirza and S. Osindero, “Conditional generative adversarial nets,” *arXiv preprint arXiv:1411.1784*, 2014.
- [57] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *Proc. of the International Conference on Learning Representations (ICLR)*, Y. Bengio and Y. LeCun, Eds., 2015.
- [58] B. Efron and R. J. Tibshirani, *An Introduction to the Bootstrap*, ser. Mono. Stat. Appl. Probab. London: Chapman and Hall, 1993.
- [59] P. Hall, “Theoretical Comparison of Bootstrap Confidence Intervals,” *The Annals of Statistics*, vol. 16, no. 3, pp. 927 – 953, 1988.
- [60] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *Proc. of the International Conference on Learning Representations (ICLR)*, Y. Bengio and Y. LeCun, Eds., 2015.
- [61] T. Salimans and D. P. Kingma, “Weight normalization: A simple reparameterization to accelerate training of deep neural networks,” in *Advances in Neural Information Processing Systems (NIPS)*, 2016.



Michael Herman is a research scientist and sub-project lead at the Bosch Center for Artificial Intelligence working on machine learning-based motion prediction for automated driving. He received his Ph.D. degree in computer science from the University of Freiburg in an industrial Ph.D. program, where his research focused on learning generalizable representations of an experts motivations from observed behavior in complex, unknown environments. His current research is focused on learning prediction models in multi-agent systems with a focus on

automated driving use cases, while his research interests include Probabilistic Inference, Deep Learning, and Imitation Learning.



Waleed Ahmed received his M.Sc degree in Automation and Robotics in 2018 from the Technical University of Dortmund, Germany, with a focus on robotics and machine learning. Since 2018, he is a development engineer at the Bosch Cognitive Systems Group with a focus on Automated Driving. He has been involved in various machine learning activities for automated driving and identification of driving strategies for driver assistance systems.



Jörg Wagner received his M.Sc degree in electrical engineering and information technology from the Karlsruhe Institute of Technology, Germany, in 2014. He is a research engineer at Bosch Center for Artificial Intelligence (BCAI) working on machine learning-based motion prediction for automated driving. While working at BCAI, he is currently pursuing the Ph.D. degree with the University of Bonn, Germany. His research interests include deep learning, interpretable and explainable AI, computer vision and time series modeling.



Lutz Bürkle received his Dr. rer. nat. (Ph.D.) in physics from the University of Freiburg, Germany, in 2001. From 1997 he was a Research Scientist at the Fraunhofer Institute of Applied Solid State Physics in Freiburg, Germany. In 2002 he joined Robert Bosch GmbH, where he has been involved in the development of various driver assistant and automated driving systems. He is currently a project manager at Bosch Corporate Research in Renningen, Germany.



Vishnu Prabhakaran received his M.Sc degree in Information Technology in 2018 from the University of Stuttgart, Germany, with a focus on machine learning and robotics. He is a research engineer at Bosch Center for Artificial Intelligence, working on multi-agent behavior prediction models for automated driving. His research interests include probabilistic inference, deep learning and time series modeling.



Ernst Kloppenburg received his Dr.-Ing. (Ph.D.) degree from University of Stuttgart in 1999 in Technical Cybernetics / Control Engineering. He is a senior expert at Bosch Center for Artificial Intelligence, working in the fields of probabilistic inference, and of verification of machine learning systems.



Nicolas Möser received his Dr. rer. nat. (Ph.D.) in 2011 from the University of Bonn, Germany, in experimental elementary particle physics at the ATLAS experiment at CERN, Geneva. In 2012 he joined Corporate Research of Robert Bosch GmbH, where he has been involved in various data mining and machine learning activities related to driver assistance and automated driving.



Claudius Gläser received his Dr.-Ing. (Ph.D.) degree in computer science from Bielefeld University, Germany, in 2012. From 2006 he was a Research Scientist with the Honda Research Institute Europe GmbH, Offenbach/Main, Germany, working in the fields of speech processing and language understanding for humanoid robots. In 2011, he joined the Corporate Research of Robert Bosch GmbH in Renningen, Germany, where he developed perception algorithms for driver assistance and highly automated driving functions. He is currently senior

expert for sensor data fusion in autonomous systems. His research interests include environment perception, multimodal sensor data fusion, multi object tracking, and machine learning for highly automated driving.



Hanna Ziesche is a research scientist at the Bosch Center for Artificial Intelligence working on robotics and information theoretic deep reinforcement learning. She received her Dr. rer. nat. (Ph.D) in 2016 from the University of Karlsruhe, Germany, in theoretical particle physics. In 2017 she joined Robert Bosch GmbH, where she has been involved in various projects on machine learning and reinforcement learning topics.