



# Driver and Pedestrian Mutual Awareness for Path Prediction and Collision Risk Estimation

Markus Roth, Jork Stapel, Riender Happee and Dariu M. Gavrila

**Abstract**—We present a novel method for vehicle-pedestrian path prediction that takes into account the awareness of the driver and the pedestrian towards each other. The method jointly models the paths of vehicle and pedestrian within a single Dynamic Bayesian Network (DBN). In this DBN, sub-graphs model the environment and entity-specific context cues of the vehicle and pedestrian (incl. awareness), which affect their future motion and allow to increase the prediction horizon. These sub-graphs share a latent state which models whether vehicle and pedestrian are on collision course; this accounts for a certain degree of motion coupling.

The method was validated with real-world data obtained by on-board vehicle sensing (stereo vision, GNSS and proprioceptive). Data consist of 93 vehicle and pedestrian encounters, spanning various awareness conditions and dynamic characteristics of the participants. In ablation studies, we quantify the benefits of various components of our proposed DBN model for path prediction and collision risk estimation.

Results show that at a prediction horizon of 1.5s, context-aware models outperform context-agnostic models in path prediction for scenarios with a dynamics change, while performing similarly otherwise. Results further indicate that driver attention-aware models improve collision risk estimation compared to driver-agnostic models.

**Index Terms**—Driver Awareness, Pedestrian Awareness, Path Prediction, Collision Risk Estimation

## I. INTRODUCTION

**M**ORE than 1.35 million people are killed yearly in traffic worldwide, according to a much cited report of the World Health Organization [1]. Pedestrians make up 23% of this number. More than half of serious crashes between vehicles and pedestrians occur outside dedicated crossing locations (e.g. zebras, traffic lights) with marked right-of-way [2].

Despite the recent interest and effort spent on higher levels of automated driving (SAE level 3+), for the foreseeable future, the reality on the road (and the accident numbers) will largely be determined by assistance systems where the driver is still required to keep the eyes on the road. This especially holds for active pedestrian safety in urban traffic.

Pedestrians are highly manoeuvrable; they can stop walking or change direction in an instant. This makes it challenging to predict their paths. Current active pedestrian safety systems on the market provide driver assistance (SAE level 0-2). They are conservatively designed in their warning and control strategy, emphasizing the current pedestrian state (i.e. position) rather than prediction, in order to avoid false system activations (i.e. automatic braking and evasive steering [3]).

M. Roth, J. Stapel, R. Happee and D. M. Gavrila are with the Intelligent Vehicles Group at TU Delft. M. Roth is also with the Environment Perception Department at Mercedes-Benz AG. markus.r.roth@daimler.com

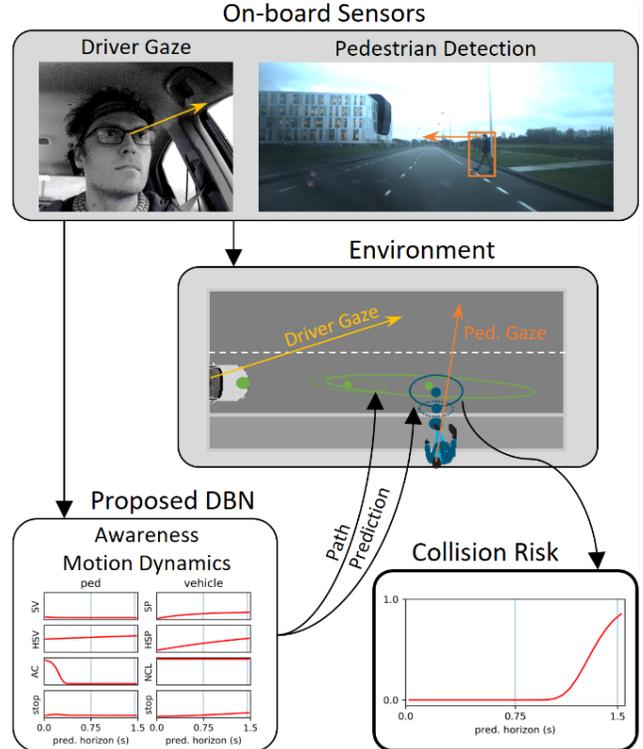


Fig. 1: The system assesses mutual awareness of pedestrian and driver in a scenario of a potentially crossing pedestrian. Cues about the driver, pedestrian and spatial environment are collected from on-board sensors. A probabilistic framework based on a Dynamic Bayesian Network (DBN) estimates latent states of awareness of the driver and pedestrian to predict their future motion. Based on the predicted paths, collision risk is estimated.

Camera-based driver monitoring systems can detect fatigue, drowsiness, distraction, gestures, signs of being drunk and readiness to take over from automated driving. On-market systems for collision warning have been employed as early as 2007 (Toyota/Lexus) by monitoring head pose and eye opening. Recent releases allow for extended SAE level 2 capabilities while driving on specially mapped highways (Cadillac Super Cruise, 2018), in traffic jams with restricted velocity (BMW Extended Traffic Jam Assistant, 2018), or in single-lane cruising (Nissan ProPilot, 2019). Mercedes-Benz's latest S-Class features a driver camera that monitors driver's readiness to take over from automated driving mode on highways in an SAE level 3 system. This legally allows the driver to perform non-

61 driving related tasks for up to 10s under specific conditions.

62 Active safety systems on the market stand to gain from  
63 improved path prediction capability of both ego-vehicle and  
64 other road users. Furthermore, they can benefit from more  
65 information regarding which specific parts of the scene have  
66 been perceived by the driver, to ascertain whether this includes  
67 the potential hazard. Ideally, a prediction horizon of 2.5s is  
68 achieved, at which point the driver “feels no danger” [4]. For  
69 the pedestrian case, we will be hard pressed to achieve accurate  
70 predictions for a 1.5s time horizon, as will become apparent.

71 In this paper, we consider the setting of a potentially crossing  
72 pedestrian and an approaching vehicle which has the right-of-  
73 way (i.e. no dedicated crossing location). We present a method  
74 which uses context cues about the spatial environment, driver-  
75 pedestrian mutual awareness and potential motion coupling to  
76 estimate the future paths of both participants and associated  
77 collision risk. See Figure 1 for an illustration of the overall  
78 system.

79 Specifically, we extend the Dynamic Bayesian Network  
80 (DBN) method from Kooij *et al.* [5], [6], which performs path  
81 prediction for an individual pedestrian, to the mutual vehicle-  
82 pedestrian case. As in [5], [6], we capture that pedestrian  
83 awareness of the on-coming vehicle will likely affect his/her  
84 future path. In our method we also model that driver awareness  
85 of the pedestrian will likely affect the future ego-vehicle path.  
86 We use head pose (pedestrian, driver) and eye gaze (driver)  
87 as proxies for awareness, as the latter cannot be determined  
88 directly.

89 There are several reasons for choosing a physics-based DBN  
90 approach for path prediction, as opposed to the popular neural  
91 networks. First, a DBN allows more easily to incorporate expert  
92 domain knowledge by means of its graphical model structure.  
93 Second, a DBN is interpretable, one can inspect the values of  
94 its latent variables and follow how it reaches its output. This is  
95 especially important for safety-critical applications. Third, one  
96 can expect a DBN to deal well with smaller datasets, as it has a  
97 comparatively small set of parameters, which will minimize the  
98 effects of over-training. Finally, recent work by Pool *et al.* [7]  
99 suggests that a DBN can deliver competitive path prediction  
100 results compared to a RNN, when its parameters are optimized  
101 by backpropagation as well.

102 The paper outline is as follows. Section II presents the  
103 related work. Section III describes the proposed context-based  
104 path prediction model for vehicle and pedestrian. Sections IV  
105 and V describe the collected dataset and the procedures for  
106 parameter estimation. Section VI describes the experimental  
107 results. Section VII provides a discussion and Section VIII  
108 lists the conclusions.

## 109 II. RELATED WORK

110 Road user path prediction has attracted a lot of attention  
111 in recent years, see surveys regarding the ego-vehicle [8] and  
112 Vulnerable Road Users [9], [10]. Path prediction methods  
113 require positions as input. Ground plane positions relative  
114 to a vehicle coordinate system can be obtained from detections  
115 in various sensors (e.g. camera [11], radar [12], LiDAR [13],  
116 or a combination thereof [13], [14]). If ground plane positions

relative to a global coordinate system are needed (e.g. this  
paper), then vehicle ego-motion compensation is necessary  
as an additional pre-processing step. For this, a combination  
of GNNS, INS and vehicle proprioceptive sensing can be  
used. Following sub-sections focus on context cues and motion  
models used for path prediction.

### 123 A. Context cues for path prediction

124 In the most rudimentary form, cues for path prediction  
125 consist of point kinematics, i.e. positions and velocities of  
126 the relevant object. It has however been well established that  
127 the use of additional “context” cues can improve path prediction  
128 performance [10]. These can be categorized into object cues,  
129 and static and dynamic environment cues.

130 Object context cues refer to cues pertaining to the object of  
131 interest itself. For example, Keller and Gavrila [15] improve  
132 pedestrian path prediction by using dense optical flow features  
133 extracted from a pedestrian bounding box. Kooij *et al.* [5]  
134 use relative head orientation as a “proxy” for the pedestrian’s  
135 awareness of the oncoming ego-vehicle while crossing. Kooij *et*  
136 *al.* [6] and Pool *et al.* [7] incorporate the arm gesture of a  
137 cyclist to predict its turn at an intersection. Quintero *et al.* [16]  
138 recover a full 3D articulated pose of a pedestrian to better  
139 predict crossing action.

140 Object context cues can also refer to properties derived from  
141 the driver of the ego-vehicle, when interested in predicting  
142 the future ego-vehicle path. Typical such cues are driver head  
143 orientation or gaze, or performed driver actions, as inferred  
144 from accelerator pedal position, braking force and steering  
145 wheel angle. For example, Roth *et al.* [17] employ driver head  
146 pose to capture the driver’s awareness of a crossing pedestrian.

147 Static environment context cues refer to elements of the  
148 static traffic infrastructure which will likely influence road user  
149 motion, such as road topology [7], [18], road markings and  
150 traffic lights.

151 Dynamic environment context cues capture the presence and  
152 motion properties of other road users (including that of the ego-  
153 vehicle itself) that may influence the target road user’s behavior,  
154 i.e. to avoid hazards or to minimize hindrance. For example, [5],  
155 [6], [17], [19], [20] use basic kinematics properties, such as  
156 relative distances and velocities, and the expected point of  
157 closest approach.

### 158 B. Motion models

159 Models for human motion trajectory estimation can be  
160 subdivided into physics-based, pattern-based and planning-  
161 based methods [10]. As motivated earlier, we focus here on  
162 physics-based methods, which represent motion by explicitly  
163 defined dynamic equations of one or more underlying dynamical  
164 models. Simple motion dynamics can be modeled by Linear  
165 Dynamical Systems (LDS), which commonly assume a linear  
166 relationship between states and measurements with Gaussian  
167 noise. Under these assumptions, the Kalman Filter (KF) [21] is  
168 an optimal filtering algorithm, which has been widely applied  
169 to pedestrian and vehicle tracking [8], [22].

170 In the scope of collision analysis, motion models play  
171 a role for predicting paths of targets such as a potentially

172 crossing pedestrian and the ego-vehicle. The probabilistic  
 173 models described here allow to extrapolate observed behaviors  
 174 into the future while accounting for uncertainties in the assumed  
 175 dynamics and observations.

176 Since traffic behavior may change at any time, a common  
 177 approach is to treat the complex dynamics by switching between  
 178 or combining multiple motion models at each prediction step,  
 179 e.g., by using Switching LDS (SLDS). SLDS can be extended  
 180 by dynamical models to incorporate contextual cues for path  
 181 prediction [6], [16]. Li *et al.* [23] combine the path prediction  
 182 output of Kooij *et al.* [6] with a sequence-to-sequence trajectory  
 183 generation method to leverage the complementary advantages  
 184 of hand-crafted models and data-driven methods.

185 Different methods have been introduced to predict the  
 186 paths of multiple interacting road users, e.g., Social Force  
 187 models for human-human interactions [24]. For pedestrian-  
 188 vehicle encounters, e.g., Kooij *et al.* [6] assume that the vehicle  
 189 does not change motion dynamics, while Braeuchle *et al.* [25]  
 190 use a Bayesian Network to find an appropriate vehicle motion  
 191 model which minimizes pedestrian injury risk. The pedestrian  
 192 motion model is fixed based on initial velocity. Gupta *et al.* [26]  
 193 simulate actions (speed up, slow down) of a self-driving vehicle  
 194 within a negotiation cycle with a crossing/yielding pedestrian  
 195 to optimize traffic throughput.

### 196 III. JOINT VEHICLE AND PEDESTRIAN PATH PREDICTION

#### 197 A. Overview and Main Contributions

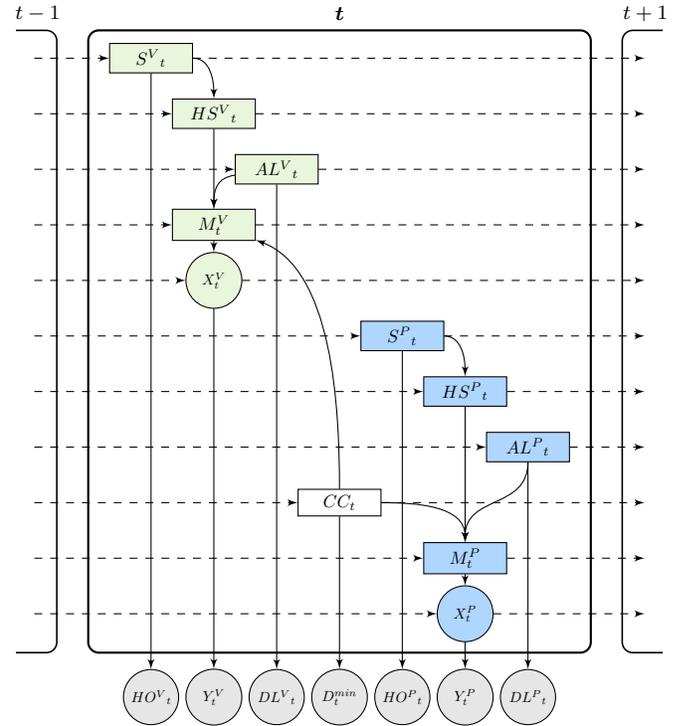
198 Kooij *et al.* [6] note that a pedestrian's decision to continue  
 199 walking or to stop in a crossing scenario is mainly influenced  
 200 by the presence of an approaching vehicle on collision course,  
 201 the pedestrian's awareness thereof, and the position of the  
 202 pedestrian with respect to the curbside. This knowledge is  
 203 encoded in a context-based SLDS (a special DBN), where  
 204 latent discrete states control the switching probabilities between  
 205 the continuous states dynamics of walking and standing.

206 In this work, we are interested in vehicle-pedestrian collision  
 207 risk, thus we extend the prediction component to the ego-  
 208 vehicle. We analogously argue that the vehicle's outcome  
 209 of continue moving or stopping is mainly influenced by the  
 210 presence of an approaching pedestrian on collision course, the  
 211 driver's awareness thereof and the distance of the vehicle to the  
 212 pedestrian's crossing location. We model pedestrian and vehicle  
 213 motion with two SLDSes which are linked to each other by a  
 214 shared latent state, which captures the motion coupling between  
 215 the two objects. The proposed DBN is shown in Figure 2 (see  
 216 Table I for the corresponding node descriptions).

217 Our main contributions are:

- 218 • We present a method for joint path prediction and collision  
 219 risk estimation of vehicle and pedestrian using observed  
 220 kinematics, mutual awareness, and environment cues.
- 221 • We provide an ablation study of the effect of various  
 222 context cues on situations where an intervention of either  
 223 road user is needed to avoid a collision.
- 224 • We apply our method on real sensor data from a vehicle.

225 Compared to our earlier work [17], we add collision risk anal-  
 226 ysis and perform more extensive evaluations (incl. *estimated*  
 227 head pose and *estimated* eye gaze in addition to invasively  
 228 measured head pose [17]) on a new and larger dataset.



229 Fig. 2: Graphical model representation of the Dynamic  
 230 Bayesian Network (DBN). Discrete nodes are rectangular,  
 231 continuous nodes are circular. Grey nodes represent observable  
 232 variables while the other nodes represent latent states. Dashed  
 233 lines depict temporal connections between latent context states  
 234 in subsequent time instances. Driver-related nodes are shaded  
 235 in green while pedestrian-related nodes are shaded in blue.  
 Context state description and purpose are provided in Table I.

#### 236 B. DBN

237 The DBN consists of two sub-graphs, one for the pedestrian  
 238 and one for the vehicle. The pedestrian sub-graph is congruent  
 239 with the DBN of Kooij *et al.* [6]. The vehicle sub-graph  
 240 displays analogous behavior for the vehicle, by encoding driver  
 241 awareness by driver gaze and braking manifestation by being  
 242 close to the crossing location of the pedestrian.

243 1) *Pedestrian-related context states:* The pedestrian  $P$  can  
 244 exhibit one of two motion types: *walking* ( $M_t^P = m_{move}^P$ ,  
 245 constant velocity) and *standing* ( $M_t^P = m_{stop}^P$ , constant  
 position). The motion state of the pedestrian contains two-  
 dimensional positions and velocities:  $X_t^P = [x_t, y_t, \dot{x}_t, \dot{y}_t]^T$ .  
 This results in the linear state transformation matrices:

$$246 A^{(m_{move}^P)} = \begin{bmatrix} 1 & 0 & \Delta t & 0 \\ 0 & 1 & 0 & \Delta t \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, A^{(m_{stop}^P)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (1)$$

247 The vehicle observes pedestrian world positions  $Y_t^P \in \mathbb{R}^2$   
 248 without velocities, resulting in the corresponding observation  
 249 matrix  $C^P = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$ .

250 For the context-based SLDS, the switching state  $M_t^P$  of the  
 251 pedestrian motion model is encoded in the DBN as a categorical  
 252 distribution  $M_{t+1}^P = \text{Cat}(M_t^P, AL_{t+1}^P, HS_{t+1}^P, CC_{t+1})$  as  
 253 shown in Figure 2. The pedestrian awareness context  $S_t^P$

TABLE I: Latent context states, their associated observation and the purpose within the DBN structure. States are grouped by vehicle/driver (common superscript  $V$ ), pedestrian (superscript  $P$ ) and shared contexts.

Latent State	Abbr.	Observation	Abbr.	Purpose
driver-sees-pedestrian	$S^V$	driver-head-orientation (gaze)	$HO^V$	encodes driver's awareness of the pedestrian
driver-has-seen-pedestrian	$HS^V$	-	-	memorizes driver's (past) awareness of the pedestrian
vehicle-at-location	$AL^V$	vehicle-distance-to-location	$DL^V$	manifests typical location of braking (ped. crossing location)
vehicle-motion-model	$M^V$	-	-	switches between <i>driving</i> and <i>braking</i> LDS
vehicle-position-state	$X^V$	vehicle-position	$Y^V$	LDS for vehicle state estimation
pedestrian-sees-vehicle	$S^P$	pedestrian-head-orientation	$HO^P$	encodes pedestrian's awareness of the driver/vehicle
pedestrian-has-seen-vehicle	$HS^P$	-	-	memorizes pedestrian's (past) awareness of the driver/vehicle
pedestrian-at-location	$AL^P$	pedestrian-distance-to-location	$DL^P$	manifests typical location of stopping (curb)
pedestrian-motion-model	$M^P$	-	-	switches between <i>walking</i> and <i>standing</i> LDS
pedestrian-position-state	$X^P$	pedestrian-position	$Y^P$	LDS for pedestrian state estimation
collision-course	$CC$	minimum-future-distance	$D^{min}$	separates early crossings from critical crossing

243 models whether the pedestrian sees the approaching vehicle.  
 244 Head orientation  $HO^P_t$  forms the evidence. The context  
 245 variable  $HS^P_t$  memorizes whether the pedestrian has seen the  
 246 vehicle in the past, acting as a logical *OR* between previous  
 247  $HS^P_{t-1}$  and current  $S^P_t$ . The environment context  $AL^P_t$   
 248 models whether the pedestrian is near the curb, thus encoding  
 249 where a pedestrian would normally stop to yield for oncoming  
 250 traffic.

2) *Vehicle-related context states*: The vehicle motion state is  $X_t^V = [x_t, y_t, \dot{x}_t, \dot{y}_t]^T$ . It uses a constant velocity model while driving, and a velocity decay model for braking:

$$A^{(m_{\text{move}}^V)} = \begin{bmatrix} 1 & 0 & \Delta t & 0 \\ 0 & 1 & 0 & \Delta t \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, A^{(m_{\text{stop}}^V)} = \begin{bmatrix} 1 & 0 & \Delta t & 0 \\ 0 & 1 & 0 & \Delta t \\ 0 & 0 & d & 0 \\ 0 & 0 & 0 & d \end{bmatrix} \quad (2)$$

251 The decay parameter  $d = \sqrt[10]{0.5} \approx 0.93$  is empirically  
 252 chosen to represent a velocity half-life of 0.5 s, i.e., the velocity  
 253 becomes  $d^{10} = 0.5$  of its initial value after 10 discrete time steps  
 254 (0.5 s). This results in a mean initial deceleration of  $\sim 4.2 \text{ m/s}^2$   
 255 over the first second, thus reflecting moderate braking. Also, the  
 256 vehicle  $V$  observes its own velocity, resulting in the observation  
 257 matrix  $C^V = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$ .

258 For the vehicle, the context-based SLDS' switching state  
 259  $M^V$  is encoded as a categorical distribution  $M_{t+1}^V =$   
 260  $\text{Cat}(M_t^V, AL_{t+1}^V, HS_{t+1}^V, CC_{t+1})$ . The driver awareness  
 261 context  $S^V_t$  models the driver's awareness of the pedestrian.  
 262 It is inferred from the attention eccentricity  $HO^V_t$ , i.e., the  
 263 absolute visual angle difference between the driver's center  
 264 of gaze (or head direction) and the pedestrian. The context  
 265 variable  $HS^V_t$  memorizes whether the driver has seen the  
 266 pedestrian analogous to  $HS^P_t$ . The static environment context  
 267  $AL^V_t$  indicates whether the vehicle is at a distance from the  
 268 pedestrian's crossing location where the driver can be expected  
 269 to yield, assuming he/she has the intention to do so.

270 3) *Shared context state*: Both pedestrian and vehicle dynam-  
 271 ics depend on  $CC_t$ , which indicates whether pedestrian and  
 272 vehicle are on a collision course. It uses the minimum distance  
 273  $D_t^{min}$  obtained when linearly extrapolating the trajectories with  
 274 their momentary estimated velocities [20].

### C. Inference

275 During inference the DBN states are propagated over time  
 276 by incorporating observations in a forward filtering procedure  
 277 (predict, update) following [6]. At each time step  $t$ , the entire  
 278 state of the DBN is represented by the 9 discrete latent states  
 279 (4 vehicle, 4 pedestrian, 1 shared) and two partially observable  
 280 continuous latent states ( $X_t^V, X_t^P$ ), see Figure 2. During the  
 281 predict step, the value of each discrete latent state changes  
 282 according to a fixed transition table, based on the values of  
 283 its input states, i.e., each state's input nodes in Figure 2,  
 284 including the state from the previous time step  $t-1$  (dashed  
 285 line). During the update step, observations are incorporated  
 286 based on the context likelihood distributions, see Figure 3. The  
 287 intermediate goal is to have the motion model switching states  
 288 for both vehicle ( $M_t^V$ ) and pedestrian ( $M_t^P$ ) which represent  
 289 the switching probability of the SLDS of each road user. The  
 290 two continuous latent states  $X_t^V, X_t^P$  are propagated over time  
 291 using observations ( $Y_t^V, Y_t^P$ ) by standard LDS means, i.e.,  
 292 Kalman filter. Prediction into future without observation follows  
 293 the same procedure, but without the update steps. Overall,  
 294 this results in predicted motion states including uncertainties  
 295 for both vehicle and pedestrian. To keep inference tractable,  
 296 we apply Assumed Density Filtering [27], resulting in the  
 297 probability distributions of  $X_t^V, X_t^P$  to be each modeled by a  
 298 Gaussian Mixture ( $K=2$ ).  
 299

## IV. PARAMETER ESTIMATION

300 We set the DBN model parameters by performing a data-  
 301 driven initialization step, followed by a gradient-based opti-  
 302 mization step, using the dataset we introduce in Section V.  
 303

### A. Model parameter initialization

304 Model parameters relate to motion dynamics and context.  
 305 They are initialized similar to Kooij *et al.* [6].  
 306

307 *Motion Dynamics*: The underlying motion models of  $M^V$   
 308 and  $M^P$  are represented by LDSes which model process  
 309 noise  $Q$  and observation uncertainty  $R$ . Process noise  $Q$  of  
 310 vehicle and pedestrian are set for both position and velocity  
 311 states and are limited to diagonal matrix entries. Values were  
 312 selected to reflect model uncertainty under typical velocity

changes of drivers and pedestrians [28], [29]. Observation noise  $R$  is set to reflect typical variance of measurement noise for pedestrian detection and vehicle movement observed on-board our vehicle, see Section V. The motion state transition matrices were obtained as follows. The vehicle motion state  $M^V$  was categorized as *braking* when such activity was detected, analogous to  $AL^V$ , and as *driving* otherwise. The pedestrian motion state  $M^P$  was categorized as *standing* in all scenarios where a pedestrian stops starting from three frames preceding  $TTE = 0$  (see Section V-B for definition of TTE), similarly to  $AL^P$  below. The motion state at all other time instants was categorized as *walking*. The motion state transition matrices were then obtained by counting and normalizing the occurrences of the respective transitions. The initial motion states assume the vehicle and pedestrian are driving and walking.

*Context:* To obtain the parameters for binary context states, we need to establish their ground truth values; we do so in a two-step approach. In the first step, ground truth values were roughly obtained by setting some states to the same values for the entire scenario based on its definition ( $S^P$ ,  $S^V$ ,  $CC$ ), by manual annotation ( $AL^P = 1 \iff TTE = 0$ ), or by an automatic observable criterion ( $AL^V = 1$  for all moments after first deceleration, i.e., pressing the brake pedal). This yields the context likelihood distributions as shown by the histograms in Figure 3. Parametric distributions were fitted by Maximum-Likelihood-Estimation and are shown by line plots. The parametric form of the distributions was chosen heuristically: Gaussian ( $DL^P$ ,  $DL^V$ ), Gamma ( $D^{min}$ ,  $HO^V$ ) or von-Mises distribution ( $HO^P$ ).

In a second step, more accurate ground truth values for the context states were obtained on the basis of the obtained context likelihood distributions. For context states  $AL^V$ ,  $AL^P$  and  $CC$ , the values were re-assigned based on maximum likelihood criterion (e.g.,  $CC = 1 \iff D^{min} < 2.6$  m, see Figure 3a). For  $S^P$  and  $S^V$ , re-assignment was done heuristically. We re-assigned  $S^P = 1 \iff HO^P \in [-30, 30]^\circ$  due to the largely overlapping distributions caused by miss-estimation of the head pose estimation algorithm. We re-assigned  $S^V = 1 \iff HO^V < 10^\circ$  whenever we use the head orientation and  $< 4^\circ$  otherwise for the eye gaze orientation. The transition matrices which represent the transition probabilities conditioned on the input states (i.e., incoming links in the DBN graph) were obtained by counting and normalizing the re-assigned binary context values between adjacent time steps. The transition probabilities of  $HS^V$  and  $HS^P$  are implemented as a binary OR in order to memorize the last state in accordance with their definition in Section III-B1.

The initial context states values were set conservatively at the beginning of each encounter: driver/pedestrian not looking, vehicle not near crossing location and pedestrian not at curb.

### B. Model parameter optimization

We employed the gradient-based method of Pool *et al.* [7] to obtain optimized model parameters. In short, the method performs back-propagation similar to neural networks on the DBN parameters on a differentiable loss function. We maximize the

observation log likelihood of the vehicle and pedestrian under their respective predicted Gaussian distributions, see Eq. (4). All intermediate time-steps up to the prediction horizon are incorporated into the loss function to enforce a consistent path. Measurements with time-to-event ( $TTE$ )  $\in [-2.5\text{ s}, 3.0\text{ s}]$  are considered for optimization, to cover periods of typical motion dynamics. Missing intermediate measurements are ignored for optimization.  $TTE$  is defined in Section V.

Optimization has been performed while enforcing properties of the DBN variables to keep the state representation interpretable, such as probabilities residing in  $[0, 1]$  and process and observation noises remaining positive definite. We also enforce the latter to be diagonal matrices with variability along elements of main direction of travel to reduce degrees of freedom and obtain more stable convergence in the optimization process.

The model parameters chosen for optimization are: process noises ( $Q$ ) of pedestrian and vehicle, transition probabilities, and context observation distribution parameters. The model was implemented in Python 3 using PyTorch 1.4 and was optimized using Adam [30].

## V. DATASET

### A. Scenarios

93 vehicle-pedestrian encounters with 4 trained drivers and 4 pedestrians were staged on two empty public roads. Each encounter consisted of a single pedestrian with the intention to cross the street in front of the approaching vehicle. The encounters represented nine disjoint scenarios (8-20 encounters each) with different combinations of situation criticality (collision course/sufficient time to cross), pedestrian behavior (stop at curb/cross), pedestrian awareness of the approaching vehicle (aware/unaware), vehicle behavior (brake/continue) and driver awareness of the approaching pedestrian (aware/unaware). The included scenarios are listed in the left of Table III.

All scenarios (except the anomalous scenario 9<sup>a</sup>) encode following behaviors:

- An aware pedestrian will yield to the vehicle. Pedestrian awareness is inferred from pedestrian head pose.
- An aware driver brakes for an inattentive pedestrian approaching the curb. Awareness is inferred from driver head or gaze orientation.
- In non-collision-course crossing scenarios, both participants continue walking/driving.
- Unaware participants continue walking/driving.

Scenarios 1 to 4 represent non-collision-course scenarios, meaning the pedestrian has sufficient time to cross. Scenarios 5 to 7 are safe through a change in behavior by either the driver or pedestrian due to awareness of the other participant. Scenario 8 represents a collision where both driver and pedestrian are unaware of each other's presence. Scenario 9<sup>a</sup> represents an anomalous scenario: the pedestrian crosses despite being aware of the approaching vehicle. The anomalous scenario is not considered for model parameter estimation.

Pedestrians were instructed to either “*continuously observe the vehicle*” or to “*keep facing forward and don't look at the vehicle*”. Drivers were instructed to either “*keep looking at*

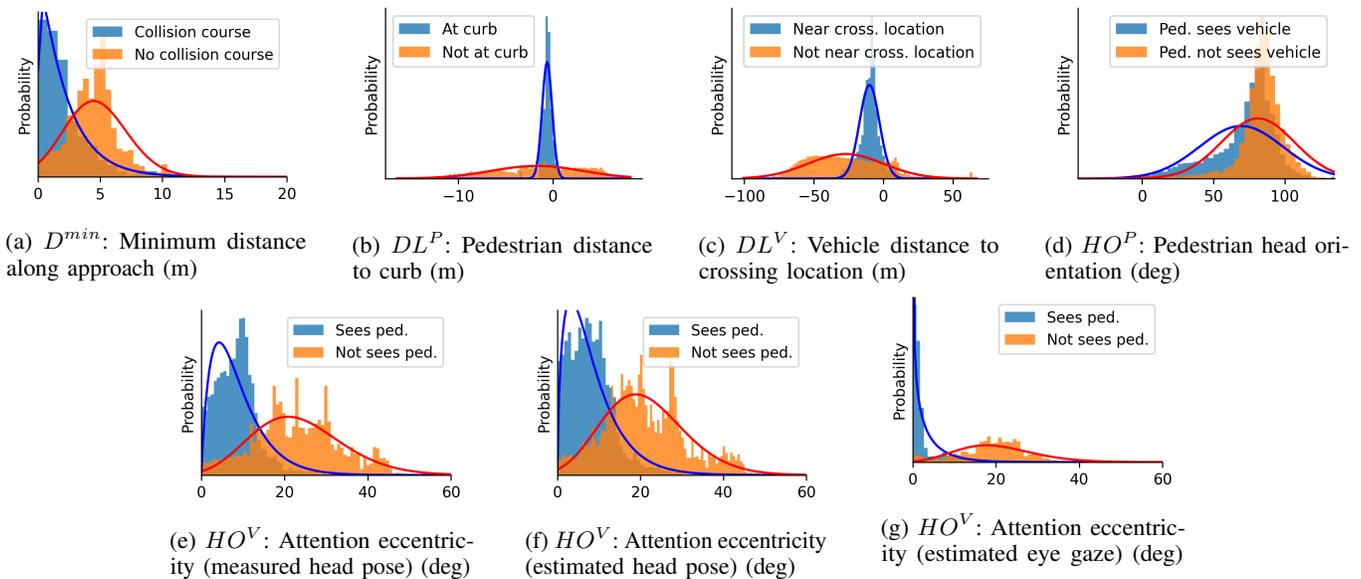


Fig. 3: Original and fitted context likelihood distributions. See Section IV-A for details.

424 the pedestrian” or to “avoid looking at the pedestrian” while  
 425 approaching the pedestrian.

426 While scenarios 8 and 9<sup>a</sup> represent collisions, naturally,  
 427 no actual collision took place during data collection. Instead,  
 428 the vehicle was brought to a full stop before colliding with  
 429 the pedestrian. The vehicle’s velocity and position data were  
 430 artificially replaced with a constant velocity model starting just  
 431 before the onset of braking.

432 To ensure safety, the road was overseen to halt the experi-  
 433 ments when other traffic entered the testing area. A co-driver  
 434 provided verbal instructions on when to brake. Target driving  
 435 speed was 20 km/h and pedestrians adopted their preferred  
 436 walking pace.

### 437 B. Instrumentation, measurements and ground truth

438 All data were collected with a TU Delft experimental  
 439 vehicle, whose instrumentation is described in further detail  
 440 in [31]. Vehicle position, orientation and velocity are ob-  
 441 tained from an ego-vehicle localization system which fuses  
 442 differential GNSS, IMU, steering wheel angle and wheel  
 443 ticks. We implement this by the Robot Operating System  
 444 (ROS) robot\_localization package [32] and gain the  
 445 transformations from vehicle frame to the world coordinate  
 446 frame, which is set to identity at the start of the system. The  
 447 GPS maintains a position accuracy of 4 cm while drift between  
 448 GPS updates is limited to 0.8% per unit of distance traveled.  
 449 The road was observed at 10 Hz using a forward-facing stereo  
 450 camera (baseline 22 cm, 1936 × 1216 px) mounted behind the  
 451 top-center of the windshield to obtain a dense stereo depth  
 452 image of the scene in front of the vehicle.

453 Driver head pose and gaze were recorded with two systems.  
 454 Estimated eye gaze and head pose were recorded with a high-  
 455 end commercial off-the-shelf eye tracker (Smarteye: 4-camera  
 456 Smart Eye Pro dx 5.0, software version 8.2, running at 60 Hz  
 457 with a gaze accuracy down to 0.5°). Secondly, measured head

pose is obtained by a head-worn infrared-reflective marker  
 458 tracked by an optical marker tracking system (Smarttrack)  
 459 mounted on the rear seat head rest [17], [33]. Additionally,  
 460 the driver was observed by a camera mounted above the  
 461 speedometer for visual verification purposes. All sensor data  
 462 were spatially calibrated and resampled to a target rate of  
 463 20 Hz.  
 464

465 Measured pedestrian positions on the ground plane were  
 466 obtained in three successive steps: (1) 2D pedestrian bounding  
 467 boxes were estimated from the forward facing camera by the  
 468 Single-Shot-Multibox-Detector (SSD) of Braun *et al.* [11].  
 469 (2) Distance to camera was found by median stereo dis-  
 470 parity [34] of the 2D bounding box. (3) Transformation of  
 471 this car-relative pedestrian position to ground plane posi-  
 472 tions in world coordinate frame was performed via ego-vehicle  
 473 localization. The time between the first pedestrian detection  
 474 and the pedestrian reaching the curb was (min / max / mean =  
 475 1.3 s / 3.2 s / 2.9 s) over the various sequences. In that period,  
 476 the pedestrian detection recall was 83%.

477 Similarly to Kooij *et al.* [5], we infer pedestrian’s focus-  
 478 of-attention from pedestrian head orientation. We apply the  
 479 method of Braun *et al.* [35] to obtain a single yaw angle  
 480 representing pedestrian head orientation.

481 In order to temporally compare prediction performance  
 482 among the various scenarios, a semantically meaningful event  
 483 was manually annotated for each sequence, as in [5], [15].  
 484 For scenarios where the pedestrian crosses, it represents the  
 485 first frame where the pedestrian’s foot crosses the curb. For  
 486 scenarios where the pedestrian stops, it represents the moment  
 487 where the last foot is placed on the ground near the curb. This  
 488 implicitly defines time-to-event (*TTE*) for each time-step of  
 489 each sequence (negative *TTE*: before event).

490 For each encounter, we obtained ground truth of the pedes-  
 491 trian position in the world coordinate frame. The pedestrian’s  
 492 path of travel is defined in the world coordinate frame as a

straight line and corresponds to the path the participants were instructed to move along. The pedestrian ground plane location is then obtained by the intersection of the annotated path of travel with the vertical plane spanned by the image column of the hip point which we manually annotated in each frame. We employ map information and ego-vehicle localization to estimate the location of the curb side.

## VI. RESULTS

To evaluate the incremental benefits of the DBN model components for an intelligent collision warning system, we compare six models with varying access to the used context cues on their joint prediction performance of vehicle- and pedestrian-path and collision risk. We adopt two evaluation metrics: the ability to predict driver and pedestrian location 1.5 s into the future, and collision risk across multiple prediction horizons. Evaluation is performed using 5-fold cross validation.

### A. Evaluation metrics

For each time  $t$ , each model creates a predictive distribution  $\tilde{P}_{t \rightarrow t+t_p}(X_t)$  for state  $X_t$  and prediction horizon  $t_p$ . Based on the predictive distributions of both vehicle and pedestrian, we evaluate individual path prediction performance and combined collision risk.

*Path prediction performance:* Two performance metrics are used to evaluate path prediction performance [5] [15]: (a) Euclidean distance error between predicted expected position and future ground truth position  $\text{GT}_{t+t_p}$ :

$$\text{error}(t_p|t) = \left| \mathbb{E} \left[ \tilde{P}_{t \rightarrow t+t_p}(X_t) \right] - \text{GT}_{t+t_p} \right| \quad (3)$$

and (b) the log likelihood of the future ground truth position  $\text{GT}_{t+t_p}$  under the predictive distribution:

$$\text{loglik}(t_p|t) = \log \left[ \tilde{P}_{t \rightarrow t+t_p}(\text{GT}_{t+t_p}) \right] \quad (4)$$

*loglik* encapsulates both the spatial error and certainty about the position observation. Larger *loglik* values denote better prediction performance.

*Collision risk:* We determine the probability for a collision by taking the integral of the predictive distributions over a collision area, which is defined by all possible intersections between vehicle and pedestrian locations. Let  $\tilde{P}_{t \rightarrow t+t_p}(X_t) = \mathcal{N}(\mu_{t \rightarrow t+t_p}, \sigma_{t \rightarrow t+t_p}^2)$  be a single Gaussian predictive position of either pedestrian P or vehicle V. The *combined* predictive position is then defined as  $\tilde{P}_{t \rightarrow t+t_p}^\phi(X_t^P, X_t^V) = \mathcal{N}(\mu_{t \rightarrow t+t_p}^P - \mu_{t \rightarrow t+t_p}^V, (\sigma_{t \rightarrow t+t_p}^P)^2 + (\sigma_{t \rightarrow t+t_p}^V)^2)$ . The collision risk predicted from  $t$  for  $t+t_p$  is given by:

$$\text{CR}(t_p|t) = \int_{A^\phi} \tilde{P}_{t \rightarrow t+t_p}^\phi(X_t^P, X_t^V) dX_t^P dX_t^V \quad (5)$$

with  $A^\phi$  being the combined spatial extent of vehicle and pedestrian. If the predictive distributions for the vehicle and the pedestrian are represented as Gaussian Mixtures (SLDS and DBN variants), the overall collision risk is given by the weighted pairwise collision risk between the Gaussian Mixture components. This extends the collision risk estimation method of Braeuchle *et al.* [25].

Context cue	LDS	SLDS	DBN.p [6]	DBN.pv	DBN.pvh	DBN.pvg
Pedestrian at-curb	-	-	x	x	x	x
Pedestrian awareness	-	-	x	x	x	x
Collision course	-	-	x	x	x	x
Vehicle near-crossing	-	-	-	x	x	x
Driver awareness	-	-	-	-	head pose	eye gaze
# Ped. motion models	1	2	2	2	2	2
# Veh. motion models	1	2	2	2	2	2

TABLE II: Context cues and number of motion models per road user used in the models. DBN suffixes denote used context: p: pedestrian [6]; v: vehicle ( $AL^V$ ); h: driver head pose; g: driver gaze. E.g., *DBN.pvg* uses pedestrian, vehicle and driver eye gaze awareness context.

For the application of collision risk warning, collision probability has to be classified into collision or no collision, and classification performance requires a ground truth for collision outcome. We define collision ground truth as true for any time instance where the vehicle and pedestrian ground truth overlap given their position and spatial extent. In order to assess the collision risk prediction performance at various prediction horizons, we select a fixed false positive rate (FPR) and find the attainable true positive rate (TPR) for each prediction horizon  $t_p$ .

### B. Model variants

We evaluate four context-aware models, including the method of Kooij *et al.* [6], which differ in their access to pedestrian and vehicle context, and compare them to two context-agnostic models. An overview of the used context cues of the models is given in Table II. All models were optimized individually as described in Section IV.

*Context-agnostic LDS:* Both linear dynamical systems for pedestrian and vehicle path prediction are instantiated by constant velocity motion models.

*Context-agnostic SLDS:* Vehicle and pedestrian motion are both modeled by context-agnostic SLDSes with the same underlying motion models as the context-aware models (driving/braking, walking/standing) described below.

*Context-aware models with varying pedestrian- and vehicle-context:* We analyze four variants of the model presented in Figure 2 which take different amounts of context into account: *DBN.p* represents the context-based pedestrian path prediction method of Kooij *et al.* [6]. The method is driver-agnostic and models the vehicle dynamics as a context-agnostic SLDS. *DBN.pv* is vehicle-aware and extends *DBN.p* with vehicle static environment cues but remains driver-agnostic. It includes proximity of the vehicle to the crossing location of the pedestrian ( $AL^V$ ). *DBN.pvh* additionally uses driver head pose as an awareness cue ( $S^V$ ). *DBN.pvg* uses driver eye gaze instead of driver head pose.

### C. Path prediction

Table III depicts average path prediction performance over various encounters of a certain scenario in terms of *loglik* and Euclidean distance error of both pedestrian and vehicle

565 for a prediction horizon  $t_p = 1.5$  s averaged over periods  
 566 where typical changes in dynamics occur (pedestrian: TTE  
 567  $\in [-0.5, 2.0]$  s, vehicle: TTE  $\in [-0.5, 3.0]$  s; TTE ranges  
 568 define times where predictions are made for). Let us consider  
 569 three scenario types.

### 570 1) Normal scenarios with no motion change

571 We first consider the normal scenarios where no motion  
 572 change occurs for a certain road user (i.e. scenarios 1-4 and  
 573 7-8 for the pedestrian, and scenarios 1-6 and 8 for the vehicle;  
 574 the respective average performances are listed in two separate  
 575 rows of Table III).

576 We see that the *LDS* for that road user has a comparatively  
 577 poor *loglik* overall (-4.4 and -13.9, resp.), as the uncertainty  
 578 region of its single-Gaussian state representation is large to  
 579 account for possible motion changes. On the other hand, its  
 580 maximum likelihood estimate is comparatively accurate: the  
 581 Euclidean distance error is smaller than that of other models  
 582 (65 cm and 52 cm, for pedestrian and vehicle resp.); this is to  
 583 be expected as its linear model precisely fits the actual motion.

584 We also observe that context-aware models are at least on-par-  
 585 with their context-agnostic (multi-motion) counterparts; cases of  
 586 outperformance suggest that the context in the former provides  
 587 more selective guidance when a motion change is probable.  
 588 Specifically, models that incorporate pedestrian context (all  
 589 DBN variants) are on-par-with (outperform) *SLDS* in terms of  
 590 the *loglik* (Euclidean distance error) metric for the pedestrian.  
 591 Models that incorporate vehicle context (*DBN.pv*, *DBN.pvh*  
 592 and *DBN.pvg*) are on-par-with *SLDS* in terms of the *loglik* and  
 593 Euclidean distance error metric for the vehicle.

### 594 2) Normal scenarios with motion change

595 Let us now consider the normal scenarios where motion  
 596 change occurs for a certain road user (i.e. scenarios 5-6 for the  
 597 pedestrian, and scenario 7 for the vehicle; the respective average  
 598 performances are listed in two separate rows of Table III).

599 We see that the context-aware models for a road user mostly  
 600 outperform their context-agnostic counterparts (*LDS* and *SLDS*)  
 601 in terms of *loglik* and Euclidean distance error for that road  
 602 user. We observe that having the full context of a road user  
 603 does not necessarily improve performance for that road user  
 604 as opposed to using only partial context (e.g. for the vehicle,  
 605 *DBN.pvh* and *DBN.pvg* underperform *DBN.pv* on Euclidean  
 606 distance error.)

607 We also observe that adding context related to the other  
 608 user does not improve performance for the original road user  
 609 (e.g. adding vehicle context *DBN.pvh* and *DBN.pvg* does not  
 610 outperform pedestrian prediction performance by *DBN.p*). An  
 611 outperformance might have been expected, as a motion change  
 612 indicates an interaction between the road users, where such  
 613 other road user context could be helpful. Apparently, the motion  
 614 coupling by means of the *CC* state variable in the DBN is  
 615 (too) weak, and is possibly overshadowed by data issues (e.g.  
 616 measurement noise, insufficient data).

617 Figure 4 shows a temporal analysis of vehicle path prediction  
 618 performance for sequences where the vehicle stops (scenario 7).  
 619 While the vehicle approaches the pedestrian with constant  
 620 velocity ( $TTE < -0.2$  s), the three compared models (*LDS*,  
 621 *SLDS*, *DBN.pvg*) show similar performance. As the vehicle

622 slows down, both *LDS* and *SLDS* increase in spread over  
 623 various sequences (shown by the standard deviations) and  
 624 gradually decrease in vehicle *loglik*. The *SLDS* model adapts  
 625 more quickly to the change of dynamics (switch from driving  
 626 to braking) compared to the *LDS*. The *DBN.pvg* model variant  
 627 anticipates the change in motion dynamics resulting in a higher  
 628 *loglik* and less uncertainty than the context-agnostic models,  
 629 therefore resulting in a better path prediction performance for  
 630 the vehicle.

### 631 3) Anomalous scenario

632 Finally, let us consider the anomalous scenario 9<sup>a</sup>. It is  
 633 anomalous as the pedestrian crosses despite seeing the vehicle.  
 634 We observe in Table III a lower prediction performance of  
 635 the context-aware models (all *DBN* variants) regarding the  
 636 pedestrian compared to the context-agnostic models (*SLDS* and  
 637 *LDS*). This is no surprise, as the context-aware models were  
 638 trained to expect stopping behaviour. Despite this, performance  
 639 degrades gracefully, since the measurements of the walking  
 640 pedestrian allow the context-aware models to infer decent  
 641 motion state estimates.

642 Figure 5 shows a comparison between driver gaze (*DBN.pvg*)  
 643 and driver head pose (*DBN.pvh*) as contextual cue for  $S^V$   
 644 (sees-pedestrian). For  $S^V = 1$ , driver gaze provides higher  
 645 classification confidence in  $HS^V$  (has-seen-pedestrian) com-  
 646 pared to head pose. For  $S^V = 0$ , both models incorrectly  
 647 believe that the driver has seen the pedestrian for a similar  
 648 fraction of sequences. However, this classification accuracy did  
 649 not yield a better vehicle path prediction performance when  
 650 comparing *DBN.pvg* to *DBN.pvh* in Table III. We attribute this  
 651 to the memorizing effect of  $HS^V$ .

652 Measured driver head pose (Smarttrack) provided virtually  
 653 identical results to estimated head pose (Smarteye) on all  
 654 scenarios, and was therefor excluded from analysis.

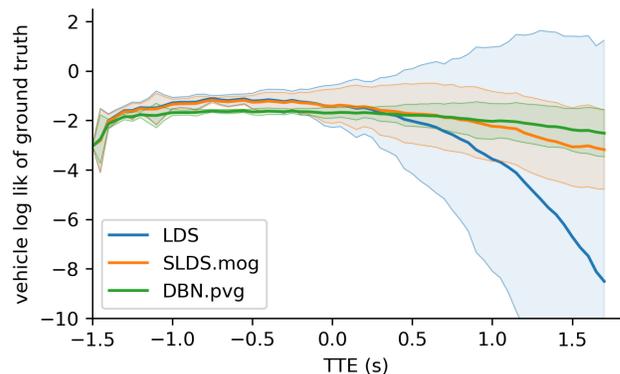


Fig. 4: *Loglik* and standard deviation over time for a braking vehicle (scenario 7) for a prediction horizon  $t_p = 1.5$  s, and drawn at the moment for which the prediction was created (i.e., the values shown at TTE = 0.0 s were predicted from measurements of TTE = -1.5 s). The vehicle initiates braking for the crossing pedestrian between -1.8 s and 0.6 s, with most vehicles braking from 0.0 s onward.

TABLE III: Scenario decomposition (left), mean path prediction performance in terms of  $\loglik$  (center) and Euclidean distance error (right) of various models for a prediction horizon of  $t_p = 1.5$ s. The top and lower halves of the table capture the prediction performances of pedestrian and vehicle along the dimension of main travel (i.e. lateral and longitudinal vs. vehicle main axis). See Section VI-B for model definitions. Higher  $\loglik$  and lower Euclidean distance error denote better prediction performance. Bold numbers denote best-performing model per scenario. Grey rows denote scenarios with a change in dynamics of the respective road user.

Scen.	CC	Ped. stops	Ped. sees	Veh. stops	Driver sees	LDS	SLDS	DBN p [6]	DBN pv	DBN pvh	DBN pvg	LDS	SLDS	DBN p [6]	DBN pv	DBN pvh	DBN pvg
						Pedestrian 1.5s $\loglik$						Pedestrian 1.5s Euclidean error (cm)					
1	0	0	0	0	0	-3.3	-2.2	<b>-2.1</b>	<b>-2.1</b>	-2.2	-2.2	64	99	<b>48</b>	51	52	51
2	0	0	0	0	1	-2.8	-2.7	-2.5	-2.5	<b>-2.4</b>	<b>-2.4</b>	<b>83</b>	140	112	110	110	111
3	0	0	1	0	0	-9.2	-3.5	<b>-3.1</b>	<b>-3.1</b>	-3.6	-3.7	77	133	<b>68</b>	71	77	73
4	0	0	1	0	1	-9.0	-2.3	-2.3	<b>-2.2</b>	-2.3	-2.3	54	73	55	50	<b>46</b>	49
5	1	1	1	0	1	-4.0	-2.4	<b>-1.8</b>	<b>-1.8</b>	-2.2	-2.2	122	131	<b>84</b>	86	91	91
6	1	1	1	0	0	-4.2	-2.5	<b>-1.7</b>	<b>-1.7</b>	-1.8	-1.8	114	131	<b>83</b>	87	87	87
7	1	0	0	1	1	-1.1	-1.5	-1.9	-1.8	-1.7	-1.7	<b>58</b>	90	71	70	70	70
8	1	0	0	0	0	-1.0	-1.3	-2.0	-1.9	-1.9	-1.9	<b>52</b>	74	63	61	63	63
9 <sup>a</sup>	1	0	1	0	0	-1.5	-1.8	-2.1	-2.0	-2.0	-2.0	<b>63</b>	100	79	77	73	73
non-anomalous, motion change (5-6)						-4.1	-2.5	<b>-1.8</b>	<b>-1.8</b>	-2.0	-2.0	118	131	<b>84</b>	87	89	89
non-anomalous, no motion change (1-4, 7-8)						-4.4	<b>-2.3</b>	<b>-2.3</b>	<b>-2.3</b>	-2.4	-2.4	<b>65</b>	102	70	69	70	70
						Vehicle 1.5s $\loglik$						Vehicle 1.5s Euclidean error (cm)					
1	0	0	0	0	0	-6.2	<b>-2.2</b>	-2.8	-2.8	-2.8	-2.8	54	53	<b>46</b>	52	55	55
2	0	0	0	0	1	-38.0	-7.4	-8.8	<b>-6.0</b>	-6.1	-6.1	60	62	<b>49</b>	53	55	55
3	0	0	1	0	0	-31.2	<b>-6.1</b>	-7.9	-7.9	-7.0	-7.0	48	52	<b>39</b>	44	51	50
4	0	0	1	0	1	-12.9	<b>-2.8</b>	-3.7	-3.6	-3.7	-3.8	63	66	<b>55</b>	56	58	58
5	1	1	1	0	1	-4.5	<b>-1.5</b>	-2.4	-2.1	-2.0	-2.0	<b>48</b>	54	<b>48</b>	117	69	69
6	1	1	1	0	0	-3.4	<b>-1.4</b>	-2.0	-2.0	-1.8	-1.8	43	52	<b>40</b>	103	61	61
7	1	0	0	1	1	-7.8	-2.7	-2.6	<b>-2.1</b>	-2.2	-2.2	245	189	195	<b>149</b>	175	175
8	1	0	0	0	0	-1.0	<b>-1.0</b>	-1.6	-1.7	-1.6	-1.6	46	47	<b>39</b>	81	45	45
9 <sup>a</sup>	1	0	1	0	0	-1.1	<b>-1.1</b>	-1.6	-1.8	-1.7	-1.7	38	47	<b>34</b>	78	45	45
non-anomalous, motion change (7)						-7.8	-2.7	-2.6	<b>-2.1</b>	-2.2	-2.2	245	189	195	<b>149</b>	175	175
non-anomalous, no motion change (1-6, 8)						-13.9	<b>-3.2</b>	-4.2	-3.7	-3.6	-3.6	52	55	<b>45</b>	72	56	56

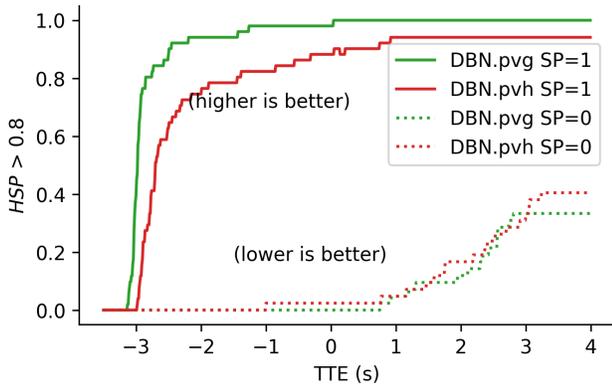


Fig. 5: Classification performance of  $DBN.pvg$  and  $DBN.pvh$  on the hidden  $HS^V$  state on sequences where driver is instructed to be attentive ( $S^V = 1$ ) and inattentive ( $S^V = 0$ ).

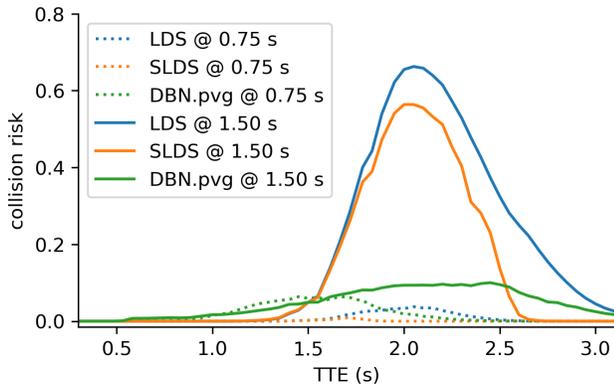
655 **D. Collision risk estimation**

656 We first compare how collision risk estimates evolve over  
 657 time for the  $LDS$ ,  $SLDS$  and  $DBN.pvg$  models on two exemplary  
 658 sequences with changing vehicle dynamics (scenario 7) and  
 659 collision (scenario 8), followed by an assessment of overall  
 660 collision risk prediction performance as function of prediction  
 661 horizon.

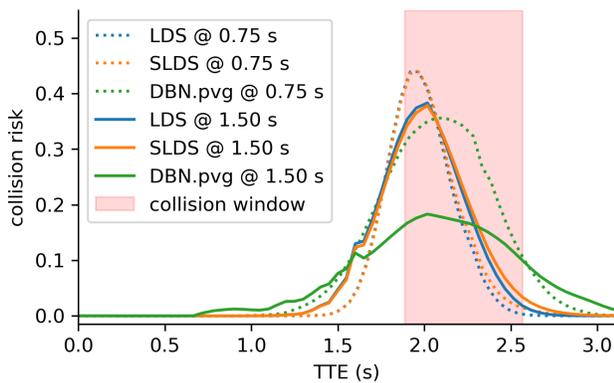
662 *Scenario-based collision risk:* Figure 6a shows collision  
 663 risk prediction for a sequence from scenario 7, where the  
 664 vehicle brakes due to an aware driver. Thus, a low predicted  
 665 collision risk is expected. For a prediction horizon  $t_p = 0.75$  s,  
 666 all models predict a negligible collision risk (dashed lines).  
 667 Predicting  $t_p = 1.5$  s into future, the  $LDS$  and  $SLDS$  models  
 668 anticipate a collision risk of 66% and 56% respectively  
 669 while the  $DBN.pvg$  model keeps a collision risk below 10%  
 670 throughout the sequence.

671 Figure 6b shows collision risk over time for one sequence  
 672 from the collision scenario (scenario 8), where both the vehicle  
 673 and the pedestrian continue their respective motion, being  
 674 unaware of each other. The *collision window* depicts all time  
 675 instances defined as a collision in accordance with Section VI-A,  
 676 i.e., where the geometries of vehicle and pedestrian overlap.  
 677 Predicting 0.75 s into the future, all compared models ( $LDS$ ,  
 678  $SLDS$ ,  $DBN.pvg$ ) depict similar maxima of collision risk within  
 679 the collision window. With increasing prediction horizon, each  
 680 model becomes less certain, resulting in a lower predicted  
 681 collision risk value.

682 The maxima are above 18% within the collision window for  
 683 the exemplarily depicted sequence. Figure 6 further shows  
 684 that only for  $DBN.pvg$ , there exists a range of collision  
 685 risk thresholds (10%–18%) for which a collision warning is  
 686 triggered in the collision sequence (Figure 6b) but not in the  
 687 non-collision sequence (Figure 6a).



(a) Sequence from scenario 7. Lower collision risk denotes better performance.



(b) Sequence from collision scenario 8. Higher collision risk denotes better performance. The collision window  $CW$  is shaded in red.

Fig. 6: Collision risk estimates obtained from different models for a braking vehicle (top) and collision (bottom) sequence. TTE indicates the time for which the predictions were made. Values are shown for prediction horizons  $t_p$  of 0.75 s and 1.5 s.

688 *Overall collision risk prediction:* To examine how collision  
 689 risk prediction performance changes with prediction horizon  
 690  $t_p$ , we select a FPR of 1% and evaluate the attainable TPR as  
 691 a function of  $t_p$ , see Figure 7. One observes that the context-  
 692 agnostic models ( $LDS$  and  $SLDS$ ) significantly under-perform  
 693 the context-aware models ( $DBN$  variants). For a prediction  
 694 horizon up to 0.75 s, all  $DBN$  variants achieve a TPR close  
 695 to 1.0. They continue to perform similarly until a prediction  
 696 horizon of about 1.3 s, after which point the driver aware  
 697 models  $DBN.pvh$  and  $DBN.pvg$  obtain a small edge. Towards  
 698 a horizon of 2.0 s, the TPR of the models drops towards 10%.

## 699 VII. DISCUSSION

700 We evaluated path prediction performance in three scenario  
 701 types within a time interval of a few seconds around a potential  
 702 motion change: in normal scenarios with no motion change,  
 703 in normal scenarios with motion change and in an anomalous  
 704 scenario. We did so as reporting aggregate performance would  
 705 not have been very insightful. This is because in reality, the  
 706 time steps in which “normal” scenarios apply with no motion  
 707 changes vastly outnumber the two other scenario types. Just

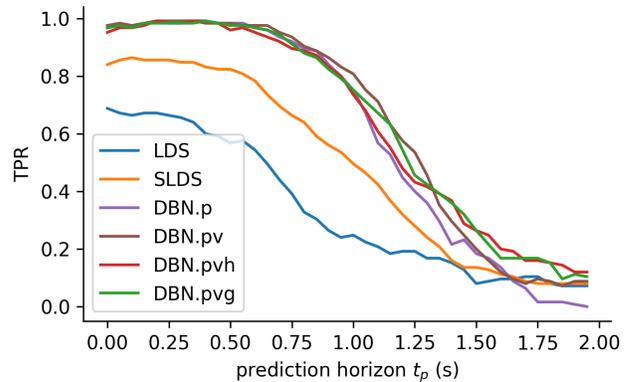


Fig. 7: Collision risk TPR of different models obtained under a 1% FPR for various prediction horizons. Higher values denote better performance.

708 considering aggregate performance would strongly favor simple  
 709 models like the  $LDS$  (or a parameter setting of a more complex  
 710 model that essentially implements such a simple model).  
 711 However, the time instants involving motion changes should  
 712 arguably carry more weight, as they might strongly induce  
 713 changes in collision risk. Listing separate performance values  
 714 for various scenario types allows to side-step this weighting  
 715 issue for now.

716 For normal scenarios with no motion change, the single-  
 717 motion model  $LDS$  performs best in terms of Euclidean distance  
 718 error, albeit with by far the worst  $loglik$  performance of all  
 719 models. Context-aware models ( $DBN.p$ ,  $DBN.pv$ ,  $DBN.pvh$ ,  
 720  $DBN.pvg$ ) were at least on-par-with their context-agnostic  
 721 (multi-motion) versions ( $SLDS$ ). They remained competitive  
 722 with the  $LDS$  on Euclidean distance error. The normal scenarios  
 723 with motion changes are those settings where the context-aware  
 724 models can potentially shine. Indeed, we found the context-  
 725 aware models to mostly outperform their context-agnostic  
 726 counterparts ( $LDS$  and  $SLDS$ ). Anomalous situations which defy  
 727 the anticipated motions, but still occur in real-world traffic,  
 728 provide a challenge to a context-aware model. They might  
 729 contradict the expert knowledge encoded in the  $DBN$  structure  
 730 or will not adhere to the parameters estimated on a training  
 731 set. Fortunately, the probabilistic modeling allows for softer  
 732 decisions: the switch of motion dynamics not only depends on  
 733 the pre-conditioning context, but also on the current positional  
 734 observations. Indeed, the performances of context-aware models  
 735 were shown to remain competitive with that of context-agnostic  
 736 counterparts.

737 Overall, one observes that the models using both pedestrian  
 738 and vehicle context ( $DBN.pv$ ,  $DBN.pvg$ ,  $DBN.pvh$ ) performed  
 739 best over the three time scenario types. Full context was not  
 740 shown to improve path prediction performance (i.e.  $DBN.pvg$   
 741 and  $DBN.pvh$  not outperforming  $DBN.pv$ ). While  $DBN.pv$ ,  
 742  $DBN.pvh$  and  $DBN.pvg$  encode typical vehicle braking locations,  
 743 variation in braking behavior seems to limit the predictive  
 744 value of the driver awareness cue. Contrary to our expecta-  
 745 tions, measuring driver gaze ( $DBN.pvg$ ) yielded similar path  
 746 prediction and collision risk estimation performance compared

to measuring driver head pose (*DBN.pvh*), i.e. see Figure 7. However, when multiple road users or driving distractions are introduced, it is likely that driver awareness will be disambiguated more accurately from gaze compared to head-pose. Other fixation-related metrics may provide further insights in driver awareness, such as number of fixations, total fixation duration and angle of first saccade landing within  $2^\circ$  of the pedestrian [36], though such evaluations would require natural as opposed to instructed viewing behavior, and other spatial regions competing for attention.

In this paper, we chose to model mutual awareness and interaction between vehicle and pedestrian loosely, by means of the shared context state *CC* (collision course) of the respective DBN sub-graphs. This has the advantage that we could easily scale-up to multiple road users, as their DBN sub-graphs can be designed and optimized individually, and the number of dependencies grow linearly. On the other hand, some limitations result from this loose motion coupling. The driver-aware models (*DBN.pvh*, *DBN.pvg*) encode the following: if one road user *A* is aware of the other *B*, this influences the motion of *A* which affects the shared collision course latent state *CC*, which in turn influences the motion of *B*. Not modeling the dependency between awareness of *A* and motion of *B* directly might lead to decreased performance. Consider the path prediction performance of the vehicle in scenarios 5 and 7. In both scenarios, the driver sees the pedestrian, however, only in scenario 7 the vehicle stops (due to the unaware pedestrian). The fact that the vehicle motion in the driver-aware models is not directly influenced by the pedestrian’s awareness might contribute to why *DBN.pvh* and *DBN.pvg* are not the best performing models for scenario 7.

DBNs provide a versatile structure to model expert knowledge. Dependencies amongst pairs of road users could be added, but limiting them to close spatial proximity, to remain scalable with increasing number of road users. Additional cues could be integrated, such as “exchanged” awareness [37], i.e. modeling the driver’s belief about the pedestrian’s awareness in addition to the driver’s awareness of the pedestrian’s presence.

One of the main insights of this paper is that context cues can help. However, simply using more complex motion models with additional context cues does not necessarily help prediction performance, if those context cues are not sufficiently informative or they cannot be reliably inferred from sensor measurements. Differences in path prediction performance between context cues can be very subtle and might also not materialize due to small data sample effects and due to errors in the estimation of ground truth.

## VIII. CONCLUSION

We presented a novel method for vehicle-pedestrian path prediction that takes into account the awareness of the driver and the pedestrian towards each other. The method jointly modeled the paths of a vehicle and a pedestrian within a single Dynamic Bayesian Network (*DBN*). Subsequently, collision risk was estimated by a probabilistic intersection operation. Overall, this work demonstrated an integrated system from on-board sensing up to collision warning.

We evaluated the incremental benefits of pedestrian- and vehicle-context in six models with varying access to the used context cues, namely Linear Dynamical System (*LDS*, one motion model), Switching Linear Dynamical System (*SLDS*, two motion models), *DBN.p* (pedestrian aware), *DBN.pv* (vehicle-aware and driver-agnostic), *DBN.pvg* (driver-gaze as awareness cue) and *DBN.pvh* (driver head pose as awareness cue).

For normal scenarios with no motion change, the single-motion model *LDS* performed best in terms of Euclidean distance error, albeit with the worst *loglik* performance by far of all models. Context-aware models (*DBN.p*, *DBN.pv*, *DBN.pvh*, *DBN.pvg*) were at least on-par-with their context-agnostic (multi-motion) versions (*SLDS*). They remained competitive with the *LDS* on Euclidean distance error. On the normal scenarios with motion changes we found the context-aware models to mostly outperform their context-agnostic counterparts (*LDS* and *SLDS*). Even in an anomalous scenario, the performances of context-aware models were shown to remain competitive with that of context-agnostic counterparts. Overall, models using both pedestrian and vehicle context (*DBN.pv*, *DBN.pvg*, *DBN.pvh*) performed best on path prediction. This was also reflected in collision risk estimation performance. For example, the collision risk warning true positive rate (TPR) was raised from 18% (pedestrian-aware model *DBN.p* of Kooij *et al.* [6]) to 27% for *DBN.pvg* for a prediction horizon of 1.5 s and a false positive rate (FPR) of 1% over the dataset.

Future work could involve improved pedestrian localization (e.g. sensor data fusion), additional and more realistic motion models within the *SLDS*, and more sophisticated context modeling (e.g. driver awareness by fixation cues). Tests are needed on large naturalistic datasets, consisting of a rich set of traffic scenarios with possibly multiple road users.

## ACKNOWLEDGMENT

We thank Ewoud Pool for sharing his code and expertise on DBN. We also thank Markus Braun for providing pedestrian detections and head pose on our dataset. This project was partially funded by the NWO-TTW Foundation, the Netherlands, within the IAVTRM project (#13712).

## REFERENCES

- [1] WHO, “Global Status Report on Road Safety,” *World Health Organization*, pp. 1–20, 2018.
- [2] European Commission, “Pedestrians and Cyclists,” *European Commission, Directorate General for Transport*, pp. 1–44, February 2018.
- [3] C. G. Keller, T. Dang, H. Fritz, A. Joos, C. Rabe, and D. M. Gavrila, “Active pedestrian safety by automatic braking and evasive steering,” *IEEE Trans. on Intelligent Transportation Systems*, vol. 12, no. 4, pp. 1292–1304, 2011.
- [4] H. Winner, “Fundamentals of Collision Protection Systems,” in *Handbook of Driver Assistance Systems*. Springer, 2016, pp. 1149–1176.
- [5] J. F. P. Kooij, N. Schneider, F. Flohr, and D. M. Gavrila, “Context-based Pedestrian Path Prediction,” in *European Conf. on Computer Vision (ECCV)*, 2014, pp. 618–633.
- [6] J. F. P. Kooij, F. Flohr, E. A. I. Pool, and D. M. Gavrila, “Context-Based Path Prediction for Targets with Switching Dynamics,” *Int. Journal of Computer Vision (IJCV)*, vol. 127, no. 3, pp. 239–262, 2019.
- [7] E. A. I. Pool, J. F. P. Kooij, and D. M. Gavrila, “Crafted vs. Learned Representations in Predictive Models - A Case Study on Cyclist Path Prediction,” *IEEE Trans. on Intelligent Vehicles*, DOI:10.1109/TIV.2021.3064253, 2021.

- 863 [8] S. Lefèvre, D. Vasquez, and C. Laugier, "A survey on motion prediction  
864 and risk assessment for intelligent vehicles," *ROBOMECH Journal*, vol. 1,  
865 no. 1, p. 1, 2014.
- 866 [9] D. Ridet, E. Rehder, M. Lauer, C. Stiller, and D. Wolf, "A Literature  
867 Review on the Prediction of Pedestrian Behavior in Urban Scenarios,"  
868 in *IEEE Int. Conf. on Intelligent Transportation Systems (ITSC)*, 2018,  
869 pp. 3105–3112.
- 870 [10] A. Rudenko, L. Palmieri, M. Herman, K. M. Kitani, D. M. Gavrila, and  
871 K. O. Arras, "Human motion trajectory prediction: a survey," *The Int.*  
872 *Journal of Robotics Research*, vol. 39, no. 8, pp. 895–935, 2020.
- 873 [11] M. Braun, S. Krebs, F. Flohr, and D. Gavrila, "EuroCity Persons: A Novel  
874 Benchmark for Person Detection in Traffic Scenes," *IEEE Trans. on*  
875 *Pattern Analysis and Machine Intelligence*, vol. 41, no. 8, pp. 1844–1861,  
876 2019.
- 877 [12] A. Palffy, J. Dong, J. F. P. Kooij, and D. M. Gavrila, "CNN based road  
878 user detection using the 3d radar cube," *IEEE Robotics and Automation*  
879 *Letters (RAL)*, vol. 5, no. 2, pp. 1263–1270, 2020.
- 880 [13] J. R. van der Sluis, E. A. I. Pool, and D. M. Gavrila, "An Experimental  
881 Study on 3D Person Localization in Traffic Scenes," in *IEEE Intelligent*  
882 *Vehicles Symposium (IV)*, 2020.
- 883 [14] M. Roth, D. Jargot, and D. M. Gavrila, "Deep End-to-end 3D Person  
884 Detection from Camera and Lidar," in *IEEE Int. Conf. on Intelligent*  
885 *Transportation Systems (ITSC)*, 2019, pp. 521–527.
- 886 [15] C. G. Keller and D. M. Gavrila, "Will the Pedestrian Cross? A Study on  
887 Pedestrian Path Prediction," *IEEE Trans. on Intelligent Transportation*  
888 *Systems*, vol. 15, no. 2, pp. 494–506, 2014.
- 889 [16] R. Quintero, I. Parra, D. F. Llorca, and M. A. Sotelo, "Pedestrian  
890 Intention and Pose Prediction through Dynamical Models and Behaviour  
891 Classification," in *IEEE Int. Conf. on Intelligent Transportation Systems*  
892 *(ITSC)*, 2015, pp. 83–88.
- 893 [17] M. Roth, F. Flohr, and D. M. Gavrila, "Driver and Pedestrian Awareness-  
894 based Collision Risk Analysis," in *IEEE Intelligent Vehicles Symposium*  
895 *(IV)*, 2016, pp. 454–459.
- 896 [18] E. A. I. Pool, J. F. P. Kooij, and D. M. Gavrila, "Using Road Topology to  
897 Improve Cyclist Path Prediction," in *IEEE Intelligent Vehicles Symposium*  
898 *(IV)*, 2017, pp. 289–296.
- 899 [19] S. Neogi, M. Hoy, K. Dang, H. Yu, and J. Dauwels, "Context Model  
900 for Pedestrian Intention Prediction Using Factored Latent-Dynamic  
901 Conditional Random Fields," *IEEE Trans. on Intelligent Transportation*  
902 *Systems*, pp. 1–12, 2020.
- 903 [20] S. Pellegrini, A. Ess, K. Schindler, and L. van Gool, "You'll never walk  
904 alone: Modeling social behavior for multi-target tracking," in *IEEE Int.*  
905 *Conf. on Computer Vision (ICCV)*, 2009, pp. 261–268.
- 906 [21] G. Welch and G. Bishop, "An Introduction to the Kalman Filter," *In*  
907 *Practice*, vol. 7, no. 1, pp. 1–16, 2006.
- 908 [22] N. Schneider and D. M. Gavrila, "Pedestrian Path Prediction with  
909 Recursive Bayesian Filters: A Comparative Study," in *German Conf. on*  
910 *Pattern Recognition (DAGM GPCR)*, 2013, pp. 174–183.
- 911 [23] Y. Li, X.-Y. Lu, J. Wang, and K. Li, "Pedestrian Trajectory Prediction  
912 Combining Probabilistic Reasoning and Sequence Learning," *IEEE Trans.*  
913 *on Intelligent Vehicles*, pp. 1–1, 2020.
- 914 [24] D. Helbing and P. Molnár, "Social force model for pedestrian dynamics,"  
915 *Physical Review E*, vol. 51, no. 5, pp. 4282–4286, 1995.
- 916 [25] C. Braeuchle, J. Ruenz, F. Flehmig, W. Rosenstiel, and T. Kropf,  
917 "Situation analysis and decision making for active pedestrian protection  
918 using Bayesian networks," in *6. Tagung Fahrerassistenzsysteme, TÜV*  
919 *SÜD*, 2013.
- 920 [26] S. Gupta, M. Vasardani, and S. Winter, "Negotiation between Vehicles  
921 and Pedestrians for the Right of Way at Intersections," *IEEE Trans. on*  
922 *Intelligent Transportation Systems*, vol. 20, no. 3, pp. 888–899, 2019.
- 923 [27] T. P. Minka, "A family of algorithms for approximate Bayesian inference,"  
924 Ph.D. dissertation, Massachusetts Institute of Technology, 2001.
- 925 [28] H. Saptoadi, "Suitable Deceleration Rates for Environmental Friendly  
926 City Driving," *Int. Journal of Research in Chemical, Metallurgical and*  
927 *Civil Engineering*, vol. 4, no. 1, pp. 2–5, 2017.
- 928 [29] M. I. Nazir, K. M. A. Al Razi, Q. S. Hossain, and S. K. Adhikary,  
929 "Pedestrian Flow Characteristics At Walkways In Rajshahi Metropolitan  
930 City Of Bangladesh," in *Int. Conf. on Civil Engineering for Sustainable*  
931 *Development*, 2014, pp. 978–984.
- 932 [30] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization,"  
933 in *International Conference on Learning Representations*, 2015.
- 934 [31] L. Ferranti, B. Brito, E. Pool, Y. Zheng, R. M. Ensing, R. Happee,  
935 B. Shyrokau, J. F. P. Kooij, J. Alonso-Mora, and D. M. Gavrila,  
936 "SafeVRU: A Research Platform for the Interaction of Self-Driving  
937 Vehicles with Vulnerable Road Users," in *IEEE Intelligent Vehicles*  
938 *Symposium (IV)*, 2019, pp. 1660–1666.
- [32] T. Moore and D. Stouch, "A Generalized Extended Kalman Filter  
939 Implementation for the Robot Operating System," in *Int. Conf. on*  
940 *Intelligent Autonomous Systems*, 2016, pp. 335–348.
- [33] M. Roth and D. M. Gavrila, "DD-Pose - A large-scale Driver Head Pose  
942 Benchmark," in *IEEE Intelligent Vehicles Symposium (IV)*, 2019, pp.  
943 927–934.
- [34] H. Hirschmüller, "Stereo Processing by Semi-Global Matching and  
945 Mutual Information," *IEEE Trans. on Pattern Analysis and Machine*  
946 *Intelligence*, vol. 30, no. 2, pp. 328–341, 2008.
- [35] M. Braun, Q. Rao, Y. Wang, and F. Flohr, "Pose-RCNN: Joint Object  
948 Detection and Pose Estimation Using 3D Object Proposals," in *IEEE Int.*  
949 *Conf. on Intelligent Transportation Systems (ITSC)*, 2016, pp. 1546–1551.
- [36] J. Stapel, M. E. Hassnaoui, and R. Happee, "Measuring Driver Perception:  
951 Combining Eye-Tracking and Automated Road Scene Perception,"  
952 *Human Factors*, 2020.
- [37] Y. Wang, Y. Ren, S. Elliott, and W. Zhang, "Enabling Courteous Vehicle  
954 Interactions through Game-based and Dynamics-aware Intent Inference,"  
955 *IEEE Trans. on Intelligent Vehicles*, vol. 5, no. 2, pp. 217–228, 2020.



**Markus Roth** received the Diploma degree in computer science from Karlsruhe Institute of Technology, Germany, in 2014. Since then he is working toward the Ph.D. degree at TU Delft, The Netherlands. He is also currently with Mercedes-Benz R&D in the Environment Perception department, Stuttgart, Germany. His research interests include machine learning and video analysis for driver analysis, with a focus on driver head pose estimation and joint awareness between driver and pedestrians.



**Jork Stapel** received the MSc degree in control and simulation at the faculty of Aerospace Engineering at the TU Delft, The Netherlands, in 2015. He received the Ph.D. degree investigating on-road assessment of driver state and driver behavior in automated driving in 2021 at the same university.



**Riender Happee** received the Ph.D. degree from TU Delft, The Netherlands, in 1986 and 1992, respectively. He investigated road safety and introduced biomechanical human models for impact and comfort at TNO Automotive (1992-2007). Currently, he investigates the human interaction with automated vehicles focussing on safety, comfort and acceptance at the Delft University of Technology, the Netherlands, where he is an Associate Professor.



**Darius M. Gavrila** received the Ph.D. degree in computer science from Univ. of Maryland at College Park, USA, in 1996. From 1997 until 2016, he was with Daimler R&D, Ulm, Germany, where he became a Distinguished Scientist. In 2016, he moved to TU Delft, where he since heads the Intelligent Vehicles group as a Full Professor. His current research deals with sensor-based detection of humans and analysis of behavior in the context of self-driving vehicles. He was awarded the Outstanding Application Award 2014 and the Outstanding Researcher Award 2019, both from the IEEE Intelligent Transportation Systems Society.