

# Modeling Drivers’ Situational Awareness from Eye Gaze for Driving Assistance

Anonymous Author(s)

Affiliation

Address

email

## Abstract:

Intelligent driving assistance can alert drivers to objects in their environment; however, such systems require a model of drivers’ situational awareness (SA) (what aspects of the scene they are already aware of) to avoid unnecessary alerts. Moreover, collecting the data to train such an SA model is challenging: being an internal human cognitive state, driver SA is difficult to measure, and non-verbal signals such as eye gaze are some of the only outward manifestations of it. Traditional methods to obtain SA labels rely on probes that result in sparse, intermittent SA labels unsuitable for modeling a dense, temporally correlated process via machine learning. We propose a novel interactive labeling protocol that captures dense, continuous SA labels and use it to collect an object-level SA dataset in a VR driving simulator. Our dataset comprises 20 unique drivers’ SA labels, driving data, and gaze (over 320 minutes of driving) which will be made public. Additionally, we train an SA model from this data, formulating the object-level driver SA prediction problem as a semantic segmentation problem. Our formulation allows all objects in a scene at a timestep to be processed simultaneously, leveraging global scene context and local gaze-object relationships together. Our experiments show that this formulation leads to improved performance over common sense baselines and prior art on the SA prediction task.

**Keywords:** driver awareness, driving assistance, situational awareness

## 1 Introduction

Future Advanced Driving Assistance Systems (ADAS) might include driver assistance systems that warn users about objects in their environment that they should pay attention to. Imagine a system that runs on your intelligent vehicle while you drive, tracking important traffic objects like vehicles and pedestrians [1]. Such a system could conceivably warn you about objects that are likely to be in your path or are otherwise dangerous, improving safety for everyone on the road. However, you are not very likely to adopt such a system if it alerts you about every object on the road regardless of your awareness of it — a well documented phenomenon known as “alert fatigue” [2]. To address this gap, we tackle the real-time object-level modeling of drivers’ Situational Awareness (SA) [3], specifically the set of traffic objects (vehicles, pedestrians, and two-wheelers) in the world that the driver is aware of at any given time.

Drivers’ eye gaze is closely linked to their situational awareness [4, 5, 6]. However, inferring situational awareness from eye gaze is not as simple as just counting gazed-at objects, since we regularly use our peripheral vision and memory to build and maintain situational awareness while driving [4]. Additionally, drivers can ostensibly “gaze” at objects without gaining situational awareness, due to effects like inattentive blindness or saccading over objects without fixating on them [7].

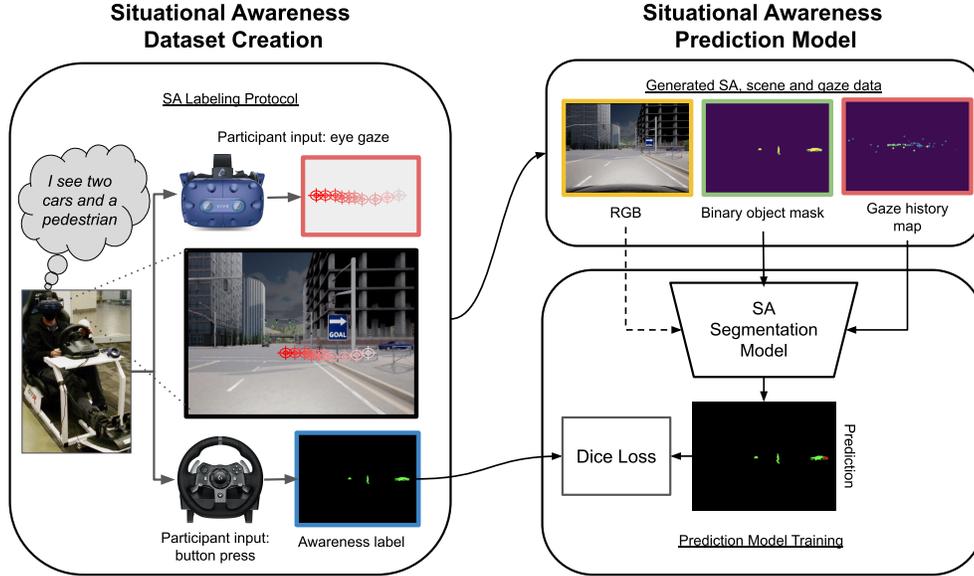


Figure 1: We collect drivers’ object-level situational awareness (SA) data via a novel interactive protocol in a VR driving simulator. We use the generated data to train a driver SA predictor from visual scene context and driver eye gaze. Casting this as a semantic segmentation problem allows our model to use global scene context and local gaze-object relationships together, processing the whole scene at once regardless of the number of objects present.

37 Thus, we aim to learn a supervised model for predicting a driver’s situational awareness from their  
 38 eye gaze and the scene context. However, training such a model requires a driving dataset with  
 39 explicitly labeled drivers’ object-level situation awareness. This dataset should be a collection of  
 40 sequences of driving events comprising the scene context, the driver eye gaze history over the scene,  
 41 and labels of the drivers’ situational awareness over each traffic object.

42 To be useful for machine learning and the downstream assistance tasks, there are a few key desiderata  
 43 for these awareness labels: 1. Labels should explicitly denote the start of the drivers’ awareness over  
 44 each object and hence be continuous. This is important since the transition of driver awareness is  
 45 crucial for determining when it is appropriate to alert the driver to the object. 2. Labels should be  
 46 dense over the set of traffic objects, i.e. we want a label for every traffic object that enters the driver’s  
 47 field of view. 3. Labels should be collected in a way that does not affect the normal gaze behavior  
 48 of the driver to avoid distribution shift between training and deployment gaze behavior.

49 Obtaining object-level awareness labels with *all* the aforementioned properties simultaneously is  
 50 challenging for a few reasons. Most current SA labeling efforts collect data either intermittently  
 51 or sparsely [8, 9, 10, 11, 12]. For instance, the common Situation Awareness Global Assessment  
 52 Technique (SAGAT) [13, 6] involves freezing and blanking the screen during occasional pauses in  
 53 simulated driving, followed by probing the driver about traffic objects present in the scene. These  
 54 collected labels are intermittent — they are valid for the moment when the simulation was paused,  
 55 but do not tell us when a driver first becomes aware of an object. Furthermore, these labels are  
 56 sparse, as the driver is only probed about objects within certain parts of the scene.

57 In this work, we introduce a novel SA labeling protocol (Sec. 3) that produces continuous and dense  
 58 object-level SA labels. As a part of our protocol, drivers indicate their awareness of all objects in  
 59 their field-of-view, by pressing directional buttons on the steering wheel controller (Fig. 1). We  
 60 collect a dataset of 80 episodes using our protocol. In each episode, drivers are instructed to drive to  
 61 an in-world goal inside a VR driving simulator [14] while following the SA labeling protocol. We  
 62 record their driving actions, eye gaze, and SA labeling button presses along with the simulator state.

63 Further, we use the aforementioned dataset to learn a model that predicts a drivers’ object-level SA  
 64 status given the scene context and a history of the driver’s eye gaze (Sec. 4). We cast this problem

65 as a semantic segmentation problem and show that it performs better than a common-sense gaze-  
66 intersection baseline and prior work that uses handcrafted features [6]. Our formulation allows us to  
67 process a variable number of objects in the scene in a single inference step as opposed to prior work  
68 which processes each object in a scene separately, repeating global computations.

69 In summary, our contributions (Fig. 1) are the following:

- 70 - **SA Labeling Protocol:** an interactive protocol for obtaining continuous and dense SA labels  
71 for on-road agents in a driving scene, without disrupting the driving task
- 72 - **SA Data Collection:** a driving dataset with continuous object-level SA labels, traffic object  
73 states, and driver eye gaze collected using our protocol in a VR driving simulator with 20  
74 drivers
- 75 - **SA Prediction Model:** a learned gaze-based driver situational awareness model which predicts  
76 SA over the scene on an object-level basis

77 Our code and dataset will be released publicly upon acceptance.

## 78 2 Related Work

79 **Measuring Situational Awareness:** Determining a driver’s internal awareness of the environment  
80 and traffic objects (vehicles, two-wheelers and pedestrians) is challenging due to our use of periph-  
81 eral vision and behaviors like intentional blindness or saccading [15]. Prior approaches for extracting  
82 information about a driver’s internal awareness involve collecting data intermittently or sparsely. An  
83 example of this is the Situation Awareness Global Assessment Technique (SAGAT), used by prior  
84 work to collect dense object-level SA labels from drivers [6]. This involved periodically pausing  
85 the simulated driving scenario, blanking the screen, and then asking the driver a series of ques-  
86 tions about their awareness of individual objects in the scene. Another approach, called Daze [16],  
87 mitigates some SAGAT issues by posing real-time queries about recent events without pausing the  
88 simulation. However, it does not yield dense object-level labels and requires looking away from  
89 the driving scene to answer affecting natural eye-gaze behavior. An influential indirect technique  
90 is the Situation Present Assessment Method (SPAM) [9], which uses real-time verbal probes about  
91 past, present, and future situations to indirectly measure SA based on response accuracy and latency.  
92 SPAM importantly also uses response times as an index of how readily this information is available.  
93 For our requirements, verbal queries have the same label sparsity issue as Daze as well as requiring  
94 manual post-processing to get machine readable annotations from verbal responses.

95 **Driver Situational Awareness Models:** Using eye gaze to infer driver attention and awareness are  
96 not new ideas, with preliminary studies having been around since at least the 1906s [17]. However,  
97 using these signals together with outward scene context for driver assistance is a relatively new  
98 area enabled by advances in sensor quality, form factors, and onboard computation —with the first  
99 papers appearing in the late-2000s [18]. Initial work used signals such as gaze direction in discrete  
100 traffic-facing zones as a crude proxy for driver attention to determine if traffic objects were causing  
101 distracted gaze. More recently, the paradigm has been to match driver gaze to objects in the traffic  
102 scene to determine whether the driver has noticed them and raise an alert when necessary [15].

103 We will focus our discussion on the process of matching gaze to traffic objects to determine which  
104 ones the driver is aware of. A naive solution is to simply count objects whose bounding boxes  
105 contain driver gaze points [19]. However objects can be perceived without being directly gazed at  
106 and 3D gaze direction estimation can have errors [20]. More recently, hand-designed feature based  
107 learning methods have emerged [13] that predict the driver’s attention given a history of their gaze  
108 relative to traffic objects. Some such methods even account for concepts of working memory from  
109 psychology [6]. However, evaluating these methods against one another is challenging. Some of  
110 these methods were evaluated qualitatively without any objective ground truth being present (SA  
111 ground truth is hard to collect as discussed in the previous section) [21]. Other methods have only  
112 been evaluated offline and on data collected using SAGAT, meaning they are evaluated on singular  
113 snapshots rather than a stream of driving data [13, 6] which prevents important aspects like aware-  
114 ness transition points to be represented in the data. Their data and models are also not publicly

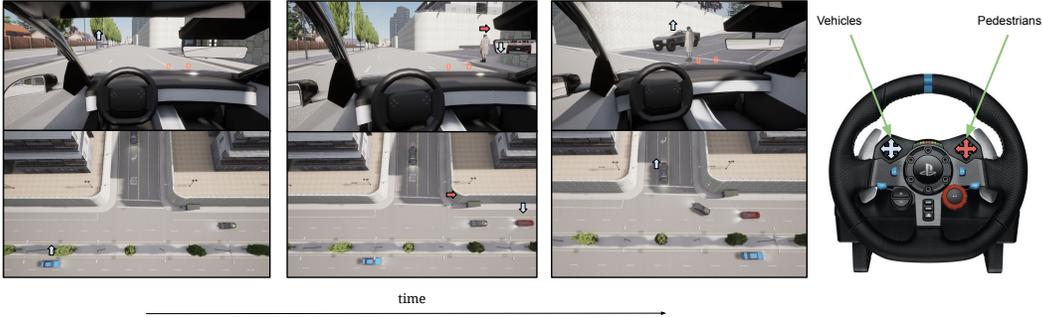


Figure 2: Example sequence of right hand turn with object-level driver responses. The top row shows the scene from the driver view and the bottom row shows the same scene via a birds-eye view. Labels are shown as colored arrows above the respective traffic object. Labels correspond to buttons on the steering wheel (right). Blue corresponds to vehicle labels and red to pedestrian labels.

115 available, making comparative evaluation difficult. To help mitigate this issue for future research,  
 116 we will release our continuously-labeled SA dataset publicly.

### 117 3 Situational Awareness Data Collection

118 We collected our driver object-level SA dataset in a VR driving simulator (DReyeVR [14]). Drivers  
 119 were asked to drive safely following a series of directional goal signs (see RGB image in Fig. 1)  
 120 along scripted routes. The drives were instructed to simultaneously follow the SA labeling protocol  
 121 to record object-level SA labels.

122 **Situational Awareness Labeling Protocol:** Under our proposed SA labeling protocol, drivers are  
 123 instructed to push a button on their steering wheel as soon as they perceive a vehicle, pedestrian,  
 124 or two-wheeler (collectively, traffic objects). For each new traffic object they perceive, they are  
 125 instructed to press one of four buttons to indicate their awareness (see Fig. 2). The button to be  
 126 pressed is determined by the relative position of the target object to the ego-vehicle. For instance, if  
 127 there is an object in front of the vehicle, the forward button should be pressed. The steering wheel  
 128 used has two sets of four buttons; the set of buttons on the left is used for vehicles and the right one  
 129 is used for 2-wheelers+pedestrians. An example sequence of traffic objects and their corresponding  
 130 button presses is shown in Fig. 2.

131 The awareness labels are generated by associating button clicks with target objects. The direction  
 132 is used to associate button presses with target objects. An object in a scene is considered ‘unaware’  
 133 until it is associated with a button press, after which it’s status is flipped to ‘aware’. More details  
 134 about how the awareness labels are generated can be found in the supplementary material.

135 **Route & traffic design:** Each route consists of a predefined source, destination, and path. Each  
 136 route also contains in-world navigational goal signs to direct the drivers along the path. Routes were  
 137 designed to have an average drive time of about 4 minutes. Each route was driven by a maximum of  
 138 8 drivers and a minimum of 4 drivers and there were a total of 15 unique routes. Participants were  
 139 pre-assigned routes so each route would be seen equally but some chose to terminate early due to  
 140 VR-induced nausea, causing an imbalance in the final number of routes.

141 At least one safety critical scenario such as a jaywalking pedestrian was included in each route. We  
 142 did so to ensure that driver gaze before and during safety critical scenarios was also represented in the  
 143 dataset. More details on the scenarios can be found in the supplementary material. The traffic along  
 144 each route was randomly generated. However, multiple objects appearing in the scene from any  
 145 single direction could lead to ambiguities in associating objects with button clicks. Hence, we limit  
 146 the number of new objects of each type appearing simultaneously at intersections in each direction  
 147 to one. Note that having different sets of buttons for vehicles and pedestrians(+two wheelers) allows  
 148 us to disambiguate between object types appearing in the same direction.

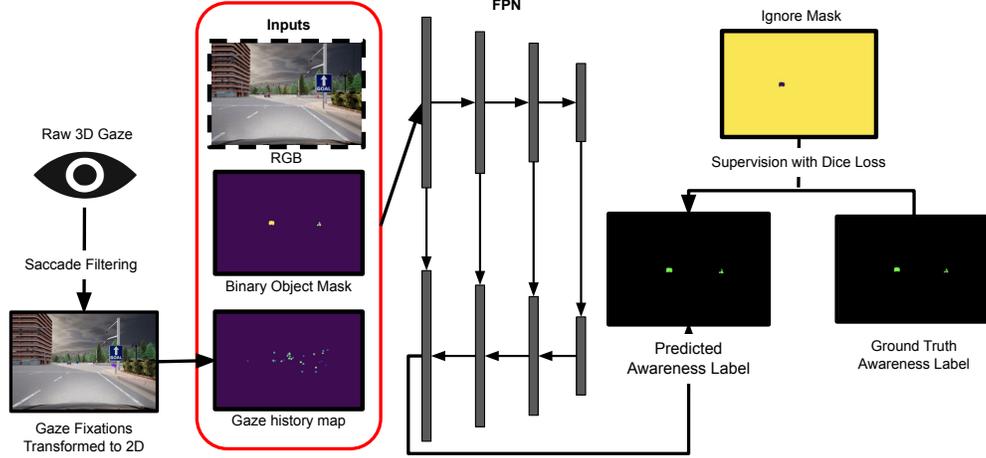


Figure 3: Object-wise SA prediction algorithm. A history of raw driver gaze is filtered to exclude saccades and then transformed to 2D pixels in the current camera position. These are used to create a gaze history map which is input together with an object segmentation of the scene (or optionally, RGB). The Feature Pyramid Network (FPN) then produces a 3 class segmentation (unaware, aware, background). During training, loss is ignored for objects which entered into the driver awareness prior to the gaze history window.

149 **Data collection details:** We ran our SA protocol with 20 participants, each with 1+ year of holding  
 150 a valid US or international driver’s license. Each participant was given a set of scripted instructions  
 151 and were first given time to interact and familiarize themselves with the interface and the simulator.  
 152 Once they were comfortable with driving in the simulator, they were introduced to the secondary  
 153 labeling task and asked to perform it while completing a trial route. Participants saw a maximum  
 154 5 non-trial routes each, but some participants did not complete all 5 routes due to the onset of  
 155 discomfort from VR cybersickness. We collected a total of 80 routes worth of data which resulted  
 156 in about 340 minutes of recorded driving time. The data collection was approved by the university’s  
 157 IRB. Some additional details about the data collection are provided in the supplementary material.

## 158 4 Modeling Driver SA

159 In modeling driver situational awareness, our goal is to predict a driver’s awareness status over all  
 160 dynamic traffic objects in the scene at a given time using scene information in conjunction with the  
 161 driver’s gaze. Specifically, for any given traffic object  $obj$ , we would like to produce a prediction of  
 162 the binary awareness status  $A_{obj}$  where  $A_{obj} \in \{aware, unaware\}$ .

163 **Problem formulation:** We cast the problem of driver SA modeling as a segmentation problem,  
 164 where the input is a visual representation of the scene in front of the user and the user’s gaze, and  
 165 the output is a prediction of the objects in the scene that the driver is aware of.

166 The scene is represented by a binary object mask indicating the location of objects in the scene (see  
 167 “Visual scene representation” below for details); the user’s eye-gaze history is input as an additional  
 168 channel in the same spatial coordinates (see “Gaze history map” in Figure 3). Under our formulation,  
 169 each timestep  $t$  represents a data point where the observations are an object mask of the scene and a  
 170 gaze map:  $O_t = (I_t^{obj} \in \mathbb{R}^{600 \times 800}, I_t^{gaze} \in \mathbb{R}^{600 \times 800})$ . The output of our model is a segmentation  
 171 map with 3 classes: *aware*, *unaware*, & *background*. Object-level awareness labels are then derived  
 172 from the output segmentation by taking the mode class of the pixels corresponding to each object  
 173 while ignoring the background class, giving us  $A_{obj}$  for each object that is visible in  $O_t$ .

174 Alternative formulations could see this posed as a classification problem, where each object is a data  
 175 point and the neural network is trained to predict a single object-level awareness label instead. In  
 176 contrast, our formulation requires one forward pass per timestep, rather than once per target object  
 177 in a timestep. This avoids repeated computations since the objects share their global context.

178 **Gaze representation:** The gaze history map  $I_t^{gaze}$  is obtained from a sequence of 3D gaze over a  
 179 historical window of length  $W$  seconds. If we sample gaze at a rate of  $s$  Hz over this window, we  
 180 obtain  $N_g = s \times W$  samples of gaze. Each gaze sample is a 3D ray  $G_i$  pointing in the direction  
 181 of the driver’s gaze, which we project into the camera coordinates to convert to a 2D pixel location.  
 182 The 3D point on this ray we project into 2D is the first point of intersection with the world while  
 183 ignoring the ego vehicle mesh (since the ego-vehicle windshield is not the point of interest). We  
 184 transform the gaze into 2D pixel coordinates  $g_i = M_t G_i \forall i \in \{1, 2, \dots, N_g\}$ , where  $M_t$  represents  
 185 a transform from world coordinates to the coordinates of the camera used at timestep  $t$ . Note that this  
 186 transformation accounts for the current pose of the ego-vehicle at time  $t$  such that the historical 3D  
 187 gaze points are transformed into pixels corresponding to their location at that previous timestep. This  
 188 means that sometimes older gaze points are out of the frame due to the traffic object’s subsequent  
 189 motion. In our experiments, we use a gaze window of  $W = 10$  s.

190 Gaze pixel locations are represented as a fixed size dot (see “Gaze history map” in Figure 3). We  
 191 also perform an ablation with a heatmap-based representation as is common with other literature  
 192 (e.g. [22]) but found this to perform worse (see Sec. 5). To include a sense of temporality in the  
 193 gaze, we fade the value of the gaze dot linearly from 255 to 10 across the window so that the most  
 194 recent gaze dots are the brightest. Additionally, since drivers cannot gain new awareness during gaze  
 195 saccades (see saccadic suppression, Ch 2. [23]), we perform gaze event detection using the I-BMM  
 196 classifier [24] and exclude saccades from the gaze map.

197 We also use an additional “ignore mask” to zero out losses from traffic objects that entered the user’s  
 198 awareness more than  $W$  seconds ago. Consider a vehicle that entered the user’s awareness 15 s prior  
 199 to the current timestep. If we use a history window  $W = 10$  s, the driver gaze correlated with  
 200 awareness of that vehicle is no longer represented, though the vehicle is still labeled as *aware*. If we  
 201 penalize the network during training for mis-classifying that object, we are penalizing a prediction  
 202 for which the network has incomplete information.

203 **Visual scene representation:** The visual scene representation uses a binary object mask to represent  
 204 the scene; the mask indicates the location of relevant dynamic traffic objects: vehicles, pedestrians,  
 205 and two-wheelers. We choose to use a fixed size ( $600 \times 800$ ) image representation from a viewpoint  
 206 in front of the ego-vehicle to control the scope of experiments. However, due to our formulation as  
 207 a segmentation problem, our model can deal with arbitrarily sized inputs. This can be useful, for  
 208 instance, when using wider aspect ratio visual inputs to represent the wide field of view that human  
 209 drivers naturally have. The binary object mask was obtained directly from CARLA, but could be  
 210 replaced by any off-the-shelf vehicle/pedestrian segmentation algorithm.

211 **Model and training details:** We used a Feature Pyramid Network [25] segmentation model with  
 212 a MobileNetV2 [26] backbone (pre-trained on ImageNet). The backbone was chosen for its low  
 213 number of parameters ( $2M$ ) and runtime efficiency. While our dataset contained a similar number  
 214 of aware to unaware objects, unaware objects usually were further from the ego-vehicle and occu-  
 215 pied much smaller portions of the input images. We calculated the ratio of the unaware pixels to  
 216 aware pixels in the label masks as about 1:20 and used an unaware class weight of 20 (background  
 217 weight= $10^{-5}$ ). We trained with the Dice loss due to its ability to handle class imbalanced data [27].

## 218 5 Evaluation & Discussion

219 **Baselines:** We compare our method to three baselines: the majority class, a common-sense gaze  
 220 intersection baseline, and a prior art baseline using handcrafted features. The “**majority class**”  
 221 baseline simply predicts the majority class in the test set (“unaware”: 53% share). The “**gaze inter-**  
 222 **section**” baseline performs a simple check: if the driver’s gaze is within the segmentation mask of a  
 223 traffic object (vehicle, pedestrian, or 2-wheeler) in the past  $T$  seconds, it assigns the *aware* label to  
 224 it (others assigned *unaware*). We use  $T = 10$ , matching the other baselines.

225 The prior art baseline (“**handcrafted features**”) is an SVM model that takes several handcrafted  
 226 features as input and produces a binary label output [6]. We re-implemented their model based on  
 227 the paper description (authors’ code or data were not publicly available). The original work lists

Model	inf. cmplx.	Acc. (↑)	Prec. (↑)	Recall (↑)	Model Ablation	Acc.	Prec.	Recall
Majority class	1	52.99%	0.53	<b>1</b>	No ignore mask	71.07%	0.79	0.62
Gaze intersection	1	46.87%	0.41	0.54	Raw gaze	73.69%	0.84	0.61
Handcrafted features [6]	N	65.47%	0.66	0.69	Gaze heatmap	76.13%	0.85	0.65
Ours	1	<b>79.21%</b>	<b>0.83</b>	0.77	No gaze fading	77.22%	0.85	0.69
					Gaze 20s hist.	74.05%	0.83	0.60
					Gaze 5s hist.	78.62%	<b>0.87</b>	0.70
					RGB	59.92%	0.83	0.30
					Ours (Full)	<b>79.21%</b>	0.83	<b>0.77</b>

(a) Performance of our model &amp; baselines

(b) Ablations for our model

Table 1: Prediction performance of models and baselines on the SA prediction task. Our model outperforms the non-trivial baselines on all 3 metrics and scales better as objects in the scene increase. [inf. cmplx. = inference time complexity with N objects, Acc. = Accuracy, Prec. = Precision]

228 5 sets of features, computed across a 10s analysis window (similar to the gaze history window in  
 229 our method): *Gaze point-based*, *Human visual sensory dependent*, *Object spatial-based*, *Object*  
 230 *property-based*, and *Human short-term memory-based*. We implemented the first 3 of these feature  
 231 sets and the object type feature (vehicle vs pedestrian) from the “Object property-based” set. Most  
 232 of the “Object property-based” features were excluded since they were difficult to compute and  
 233 required privileged scene information (*e.g.* one feature required the state of the corresponding traffic  
 234 light for every traffic object in scene; another was manually annotated). Human short-term memory-  
 235 based features were also excluded since they were difficult to compute and did not contribute much  
 236 (< 1% point) to overall performance in the original evaluation [6]. The original SVM was trained  
 237 on 1078 training samples. Since neither the trained model nor code were available, we trained our  
 238 implementation of the SVM on a subset of our training data. We trained the SVM on 10 episodes in  
 239 our train set, which is about  $3\times$  the training data used in the original work. SVM implementations  
 240 generally cannot handle very large datasets since the entire dataset is loaded into memory during  
 241 training and mini-batch SVM training is non-trivial. To train the SVM, we used a machine with  
 242 128GB RAM but could only use 15% of the training set.

243 **Experimental settings:** Our dataset contains 80 episodes of which we used 64 (80%) for training.  
 244 10% of the training episodes were used as the validation set. The test set was a separately held out  
 245 set of 16 episodes. It was partitioned so that participants were disjoint between the train and test set.  
 246 This is important since we want to test the generalization to new users; it would be impractical to  
 247 put every new driver through the SA protocol when deploying such a system.

248 We use 3 metrics to evaluate and compare methods: object-level accuracy, precision, and recall.  
 249 For precision and recall, the positive class is the “unaware” class. This is because downstream  
 250 applications such as driver assistance systems which alert the driver will care about how well our  
 251 system can predict which traffic objects the driver is not aware of. “Precision” is thus a measure  
 252 of how often our prediction of an object being unaware is correct — errors are “aware” objects  
 253 classified as “unaware.” This type of error can lead to alert fatigue for an end-user [2]. “Recall,” on  
 254 the other hand, indicates how many of the “unaware” objects in the dataset were correctly predicted  
 255 — these are objects that the driver wasn’t aware of but our system predicted that they were.

256 **Results & Discussion** Our quantitative evaluation results can be found in Table 1. The naive gaze-  
 257 intersection baseline, as expected, performs the worst, confirming that it is not enough to simply  
 258 count which objects were “gazed-at”. The prior art handcrafted features baseline performs better  
 259 but significantly worse than our method. In terms of runtime, the prior art baseline has 2 expensive  
 260 parts: computing features per object and doing SVM inference (this can be batched across objects).  
 261 On average each part takes 5 ms, resulting in a total average runtime of  $(5N + 5)$ ms on an AMD  
 262 5955WX CPU (for N objects in scene). In contrast, our network takes 11ms total for a forward pass  
 263 (on a 4090 GPU) and does not scale with the number of scene objects. Some qualitative comparisons  
 264 of model outputs can be seen in Fig. 4.

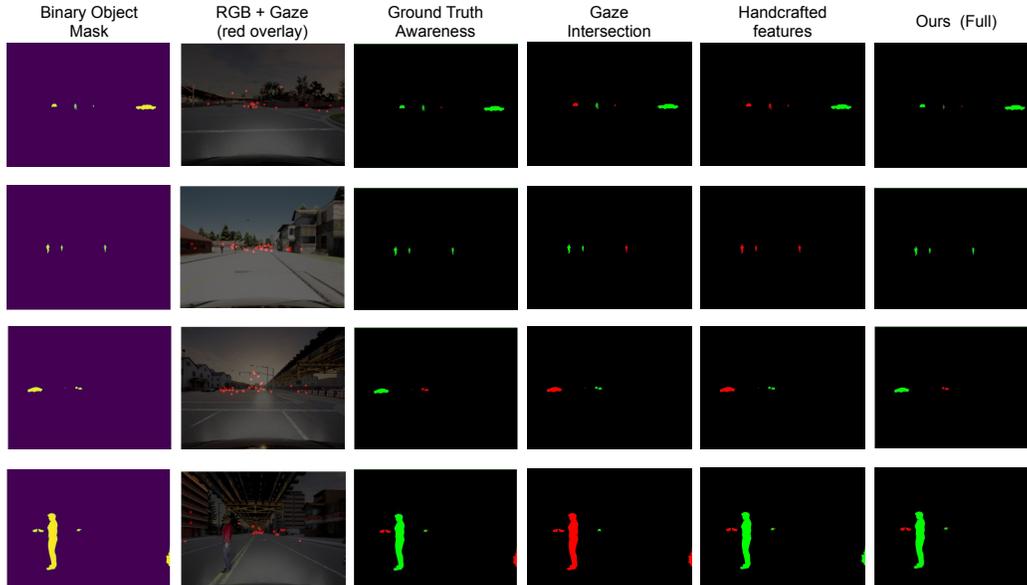


Figure 4: Qualitative results for our model and baselines. Each row represents an independent driving scene. The RGB image shows the most recent 10s of gaze overlaid as red dots.

265 Our ablations (Table 1, right) show the performance impact of several design choices described in  
 266 Sec. 4. In terms of gaze representation, the ignore mask (used to avoid penalizing mispredictions  
 267 of awareness transitions outside the gaze history window) was the most important during training  
 268 — responsible for an 8% accuracy drop when removed. Using saccade filtered gaze instead of raw  
 269 gaze was the next most important. We also investigated the use of gaze heatmaps as the gaze repre-  
 270 sentation similar to previous work [22, 28], in which each gaze point is represented by an isometric  
 271 2D Gaussian that could accumulate in weight at fixations; this performed about 3% worse than us-  
 272 ing fixed sized dots. This is similar to the issue of representing corrective clicks in an interactive  
 273 segmentation task, where a similar result has been found [29]. The results indicate that the use of  
 274 gaze fading was only responsible for about 2% of the model’s performance. This suggests that the  
 275 presence and location of a gaze point within the gaze history window contains most of the informa-  
 276 tion about awareness rather than the exact temporal order of the gaze. Finally, using an RGB image  
 277 as input resulted in 20% worse accuracy with the same model size (except the initial layer), as the  
 278 model now has to simultaneously perform segmentation and SA modeling.

279 **Limitations:** Our proposed SA labeling protocol is mainly limited by the fact that some traffic  
 280 configurations can lead to ambiguity in assigning a button — whenever there is more than one  
 281 new object of the same type (vehicle or pedestrian) from the same cardinal direction relative to the  
 282 driver. We created an interface for manual annotation to resolve ambiguities post-hoc. The biggest  
 283 limitation of our model is its static, memoryless nature. Since SA is inherently a temporal signal,  
 284 improvements can probably be achieved by performing temporal modeling. Currently, our method  
 285 treats each timestep as independent and would require an external module to implement memory.

## 286 6 Conclusion & Future Work

287 We proposed a new interactive protocol to record human drivers’ object-level situational awareness  
 288 that produces continuous and dense awareness labels. We use it to record a SA dataset with 20  
 289 drivers in a VR driving simulator. Additionally, we use this dataset to train a driver object-level SA  
 290 model by casting it as a semantic segmentation problem. Our model outperforms baselines and prior  
 291 work while scaling better to arbitrary numbers of objects in the scene. In the future, we plan to use  
 292 our driver SA model in the inner loop of a driver assistance system that provides intelligent alerts or  
 293 interventions in safety-critical situations and evaluate this in a simulator-based user study. We also  
 294 commit to releasing our code and data publicly upon acceptance in the hope that it will facilitate  
 295 more work in the domain.

## References

- 296
- 297 [1] P. Gupta, A. Biswas, H. Admoni, and D. Held. Object importance estimation using counter-  
298 factual reasoning for intelligent driving. *IEEE Robotics and Automation Letters*, 2024.
- 299 [2] J. S. Ancker, A. Edwards, S. Nosal, D. Hauser, E. Mauer, R. Kaushal, and W. the HITEC Inves-  
300 tigators. Effects of workload, work complexity, and repeated alerts on alert fatigue in a clinical  
301 decision support system. *BMC medical informatics and decision making*, 17:1–9, 2017.
- 302 [3] R. Kaufman, D. Kirsh, and N. Weibel. Developing situational awareness for joint action with  
303 autonomous vehicles. *arXiv preprint arXiv:2404.11800*, 2024.
- 304 [4] L. Gugerty et al. Situation awareness in driving. *Handbook for driving simulation in engineer-  
305 ing, medicine and psychology*, 1:265–272, 2011.
- 306 [5] T. Victor, M. Dozza, J. Bärgrman, C.-N. Boda, J. Engström, C. Flannagan, J. D. Lee, and  
307 G. Markkula. Analysis of naturalistic driving study data: Safer glances, driver inattention, and  
308 crash risk. Technical report, 2015.
- 309 [6] H. Zhu, T. Misu, S. Martin, X. Wu, and K. Akash. Improving driver situation awareness pre-  
310 diction using human visual sensory and memory mechanism. In *2021 IEEE/RSJ International  
311 Conference on Intelligent Robots and Systems (IROS)*, pages 6210–6216. IEEE, 2021.
- 312 [7] Y. Wang, Y. Wu, C. Chen, B. Wu, S. Ma, D. Wang, H. Li, and Z. Yang. Inattentive blind-  
313 ness in augmented reality head-up display-assisted driving. *International Journal of Human-  
314 Computer Interaction*, 38(9):837–850, 2022.
- 315 [8] M. R. Endsley. Design and evaluation for situation awareness enhancement. In *Proceedings of  
316 the Human Factors Society annual meeting*, volume 32, pages 97–101. Sage Publications Sage  
317 CA: Los Angeles, CA, 1988.
- 318 [9] F. T. Durso, C. A. Hackworth, T. R. Truitt, J. Crutchfield, D. Nikolic, and C. A. Manning.  
319 Situation awareness as a predictor of performance for en route air traffic controllers. *Air Traffic  
320 Control Quarterly*, 6(1):1–20, 1998.
- 321 [10] N. Martelaro, D. Sirkin, and W. Ju. Daze: a real-time situation awareness measurement tool  
322 for driving. In *Adjunct Proceedings of the 7th International Conference on Automotive User  
323 Interfaces and Interactive Vehicular Applications*, pages 158–163, 2015.
- 324 [11] R. M. Taylor. Situational awareness rating technique (sart): The development of a tool for  
325 aircrew systems design. In *Situational awareness*, pages 111–128. Routledge, 2017.
- 326 [12] J. C. de Winter, Y. B. Eisma, C. Cabrall, P. A. Hancock, and N. A. Stanton. Situation awareness  
327 based on eye movements in relation to the task environment. *Cognition, Technology & Work*,  
328 21(1):99–111, 2019.
- 329 [13] H. Kim, S. Martin, A. Tawari, T. Misu, and J. L. Gabbard. Toward real-time estimation of  
330 driver situation awareness: An eye-tracking approach based on moving objects of interest. In  
331 *2020 IEEE Intelligent Vehicles Symposium (IV)*, pages 1035–1041. IEEE, 2020.
- 332 [14] G. Silvera, A. Biswas, and H. Admoni. Dreyevr: Democratizing virtual reality driving simu-  
333 lation for behavioural & interaction research. In *Proceedings of the ACM/IEEE International  
334 Conference on Human-Robot Interaction*, pages 639–643, 2022.
- 335 [15] I. Kotseruba and J. K. Tsotsos. Behavioral research and practical models of drivers’ attention.  
336 *arXiv preprint arXiv:2104.05677*, 2021.
- 337 [16] D. Sirkin, N. Martelaro, M. Johns, and W. Ju. Toward measurement of situation awareness  
338 in autonomous vehicles. In *Proceedings of the 2017 CHI Conference on Human Factors in  
339 Computing Systems*, pages 405–415, 2017.

- 340 [17] N. A. Kaluger and G. Smith Jr. *Driver eye-movement patterns under conditions of prolonged*  
341 *driving and sleep deprivation*. PhD thesis, Ohio State University, 1969.
- 342 [18] M. M. Trivedi, T. Gandhi, and J. McCall. Looking-in and looking-out of a vehicle: Computer-  
343 vision-based enhanced vehicle safety. *IEEE Transactions on Intelligent Transportation Sys-*  
344 *tems*, 8(1):108–120, 2007. doi:10.1109/TITS.2006.889442.
- 345 [19] T. Bär, D. Linke, D. Nienhüser, and J. M. Zöllner. Seen and missed traffic objects: A traffic  
346 object-specific awareness estimation. In *2013 IEEE Intelligent Vehicles Symposium (IV)*, pages  
347 31–36. IEEE, 2013.
- 348 [20] A. Biswas and H. Admoni. Characterizing drivers’ peripheral vision via the functional field  
349 of view for intelligent driving assistance. In *2023 IEEE Intelligent Vehicles Symposium (IV)*,  
350 pages 1–8. IEEE, 2023.
- 351 [21] J. Schwehr and V. Willert. Multi-hypothesis multi-model driver’s gaze target tracking. In *2018*  
352 *21st International Conference on Intelligent Transportation Systems (ITSC)*, pages 1427–1434.  
353 IEEE, 2018.
- 354 [22] A. Palazzi, D. Abati, F. Solera, R. Cucchiara, et al. Predicting the driver’s focus of attention:  
355 the dr (eye) ve project. *IEEE transactions on pattern analysis and machine intelligence*, 41(7):  
356 1720–1733, 2018.
- 357 [23] D. Irwin. Eye movements and perception. 2001.
- 358 [24] E. Tafaj, G. Kasneji, W. Rosenstiel, and M. Bogdan. Bayesian online clustering of eye move-  
359 ment data. In *Proceedings of the symposium on eye tracking research and applications*, pages  
360 285–288, 2012.
- 361 [25] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid  
362 networks for object detection. In *Proceedings of the IEEE conference on computer vision and*  
363 *pattern recognition*, pages 2117–2125, 2017.
- 364 [26] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen. Mobilenetv2: Inverted  
365 residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision*  
366 *and pattern recognition*, pages 4510–4520, 2018.
- 367 [27] R. Zhao, B. Qian, X. Zhang, Y. Li, R. Wei, Y. Liu, and Y. Pan. Rethinking dice loss for medical  
368 image segmentation. In *2020 IEEE International Conference on Data Mining (ICDM)*, pages  
369 851–860. IEEE, 2020.
- 370 [28] A. Biswas, B. A. Pardhi, C. Chuck, J. Holtz, S. Niekum, H. Admoni, and A. Allievi. Gaze  
371 supervision for mitigating causal confusion in driving agents. In *Proceedings of the 23rd*  
372 *International Conference on Autonomous Agents and Multiagent Systems*, pages 2159–2161,  
373 2024.
- 374 [29] K. Sofiuk, I. A. Petrov, and A. Konushin. Reviving iterative training with mask guidance for  
375 interactive segmentation. In *2022 IEEE International Conference on Image Processing (ICIP)*,  
376 pages 3141–3145. IEEE, 2022.