



European Research Council
Established by the European Commission

Self-assessment Oracles for Anticipatory Testing

TECHNICAL REPORT: TR-Precrime-2019-02

Andrea Stocco, Michael Weiss, Marco Calzana, Paolo Tonella
Misbehaviour Prediction for Autonomous Driving Systems

Project no.: 787703
Funding scheme: ERC-2017-ADG
Start date of the project: January 1, 2019
Duration: 60 months

Technical report num.: TR-Precrime-2019-02
Date: August 2019
Organization: USI Università della Svizzera Italiana
Authors: Andrea Stocco, Michael Weiss, Marco Calzana, Paolo Tonella
Dissemination level: Public
Revision: 1.0

Disclaimer:

This Technical Report is a pre-print of the following publication:

Andrea Stocco, Michael Weiss, Marco Calzana, Paolo Tonella: *Misbehaviour Prediction for Autonomous Driving Systems*. Proceedings of the ACM International Conference on Software Engineering (ICSE), Seoul, South Korea, May 23–29, 2020

Please, refer to the published version when citing this work.





Università della Svizzera Italiana (USI)

Principal investigator: Prof. Paolo Tonella
E-mail: paolo.tonella@usi.ch
Address: Via Buffi, 13 – 6900 Lugano – Switzerland
Tel: +41 58 666 4848
Project website: <https://www.pre-crime.eu/>

Abstract

Deep Neural Networks (DNNs) are the core component of modern autonomous driving systems. To date, it is still unrealistic that a DNN will generalize correctly in all driving conditions. Current testing techniques consist of offline solutions that identify adversarial or corner cases for improving the training phase, and little has been done for enabling online healing of DNN-based vehicles.

In this paper, we address the problem of estimating the confidence of DNNs in response to unexpected execution contexts with the purpose of predicting potential safety-critical misbehaviours such as out of bound episodes or collisions. Our approach SELFORACLE is based on a novel concept of self-assessment oracle, which monitors the DNN confidence at runtime, to predict unsupported driving scenarios in advance. SELFORACLE uses autoencoder and time-series-based anomaly detection to reconstruct the driving scenarios seen by the car, and determine the confidence boundary of normal/unsupported conditions.

In our empirical assessment, we evaluated the effectiveness of different variants of SELFORACLE at predicting injected anomalous driving contexts, using DNN models and simulation environment from Udacity. Results show that, overall, SELFORACLE can predict 77% misbehaviours, up to 6 seconds in advance, outperforming the online input validation approach of DeepRoad by a factor almost equal to 3.

Contents

1	Introduction	1
2	Background	2
3	Problem formulation	3
4	Approach	4
4.1	Training of SELFORACLE under Nominal Driving Behaviour	5
4.1.1	Reconstructor	5
4.1.2	Probability Distribution Fitting	5
4.2	Usage Scenario	7
4.2.1	Time-aware Anomaly Score Prediction	7
5	Empirical Evaluation	7
5.1	Research Questions	7
5.2	Self Driving Car Models	7
5.3	Simulation Platform	7
5.3.1	Unexpected Context Generator	8
5.3.2	Collision/OBE Detection System	8
5.4	Procedure	9
5.4.1	Data Generation (Training Set)	9
5.4.2	SDC Model Setup & Training	9
5.4.3	Evaluation Set	9
5.4.4	SELFORACLE’s Configurations	10
5.5	Evaluation of Simulation Results	10
5.5.1	Labelling of Anomaly and Normal Windows in Evaluation Data	10
5.5.2	Metrics used for Analysis	11
5.6	Results	11
5.7	Threats to Validity	13
6	Related Work	13
7	Conclusions and Future Work	14

1 Introduction

Self-driving cars are one of the emerging technologies nowadays, and possibly the standard way of transportation in the future. Such autonomous driving systems receive data from a multitude of sensors, and analyze them in real time using deep neural networks (DNNs) to determine the driving parameters for the actuators.

To test such complex software systems, companies perform a limited number of expensive in-field tests, driving a car on real world streets, or within closed-course testing facilities [10]. This provides detailed sensor data of the vehicle that are recorded, played back, and recreated within a simulator to obtain comprehensive test scenarios. Simulation-based test scenarios allow re-testing new autopilot releases on a large numbers of nominal conditions, as well as challenging (e.g., adverse weather) and dangerous circumstances (e.g., a pedestrian suddenly crossing the road), at a low cost [10].

The potentially unlimited number of testable driving scenarios, combined with the lack of human interpretability of the internal functioning of DNNs [2], makes it difficult to predict the vehicle's (mis-)behaviour with respect to unforeseen edge-case scenarios. Misbehaviours span a wide range of situations, associated with different degrees of severity, from cases where the car does not drive smoothly (e.g., excessively high derivative of the steering angle over time), to safety-critical failures and casualties [33, 32, 11, 12]. In this respect, promptly detecting unexpected, untested behaviours is of paramount importance, so as to make sure that human-driven or self-healing corrective actions take place to ensure safety.

In this paper, we tackle the *self-assessment oracle problem* for autonomous driving system, i.e., the problem of monitoring the confidence level of a DNN-based autonomous driving system in order to timely predict the occurrence of future misbehaviours. The problem is critical because a failure in detecting an unexpected condition may have severe consequences (i.e., a fatal crash), whereas false alarms, even if not dangerous, may cause driver's discomfort and negatively affect the driving experience. Creating a self-assessment oracle that evaluates the confidence of a DNN *at runtime*, and predicts whether the system is within a low-confidence zone, is a largely unexplored research problem.

Challenges arise because unexpected driving conditions are, by definition, unknown at training time, otherwise they would be used to train a better DNN [9]. As a consequence, the problem being addressed belongs to the unsupervised class of data analysis, and we have to infer the unexpected only by looking at the normal driving scenarios. Moreover, the ensemble of possible misbehaviours for a DNN-based system is vast and necessarily domain-dependent, being associated with deviations from the functional requirements.

In recent years, researchers have proposed solutions for testing autonomous driving systems software [14, 13, 26, 34, 42, 1, 4, 5, 14]. A number of approaches propose input generation techniques that produce corner/adversarial cases, used to improve the robustness of self-driving car modules by re-training [26, 34, 42, 18]. Other works target test case generation to expose faults for extreme conditions, such as the vehicle colliding with a pedestrian [5], or driving off the road [13]. All these approaches concern *offline* solutions for improving the robustness and reliability of DNNs, achieved by enhancing the training data and the autopilot module, which are extended to include underrepresented critical scenarios. To the best of our knowledge, no existing work is focused on online confidence estimation for self-driving vehicles with the aim of anticipating the occurrence of misbehaviours so as to enable self-healing procedures that can avoid future failures.

In this paper, we propose a novel self-assessment oracle for autonomous vehicles based on *confidence estimation*, *probability distribution fitting*, and *time series analysis*. Our technique is implemented in a tool called SELFORACLE, which leverages reconstruction-based techniques from the deep learning (DL) field to analyze spatiotemporal historical driving information. The reconstruction error is used as a black-box confidence estimation for the DNN. During probability distribution fitting, SELFORACLE captures the behaviour of the self-driving car under nominal conditions, and fits a Gamma distribution to the observed data. Analytical knowledge of Gamma's parameters allows SELFORACLE to estimate an optimal confidence threshold, as a trade-off between prediction of all misbehaviours (the stricter the threshold, the better), and probability of false alarms (reduced when threshold is higher). By observing a decreasing confidence trend over time, SELFORACLE can anticipate a misbehaviour by recognizing unexpected conditions timely enough to enable healing actions (e.g., manual or automated disengagement).

We have evaluated the effectiveness of SELFORACLE on the Udacity simulator of self-driving cars [37], using DNNs available from the literature. We have modified the simulator to being able to inject unexpected driving conditions in a controllable way (i.e., day/night cycle, rain, or snow). Such injected conditions are by construction disjoint from those used to train the DNNs. In our experiments on 72 simulations, SELFORACLE is able to safely anticipate 77% out of bound episodes/crashes, up to 6 seconds in advance. The

corresponding false alarm rate in nominal driving conditions is 1%. A comparative experiment with the online input validation strategy of DeepRoad [42] shows that SELFORACLE achieves substantially superior misbehaviour prediction on all the considered effectiveness metrics (e.g., twice as high in terms of AUC-PR and almost three times higher in terms of FPR).

Our paper makes the following contributions:

Technique An unsupervised technique for misbehaviour prediction based on confidence estimation, probability distribution fitting and time series analysis, implemented in the tool SELFORACLE, which is available [35].

Simulator An extension of the Udacity simulator to inject unexpected driving conditions dynamically during the simulation.

Evaluation An empirical study showing that the reconstruction error used by SELFORACLE for time series analysis is a promising confidence metric for misbehaviour prediction, outperforming the online input validation approach of DeepRoad.

Dataset A dataset of 765 labeled simulation-based collision and out-of-bound episodes that can be used to evaluate the performance of prediction systems for autonomous driving cars.

2 Background

DNN-based Autonomous Vehicles. Self-driving cars (SDC, hereafter) have benefited from many technological advancements both in hardware and in software. Data gathered by LIDAR sensors, cameras, and GPS are analyzed in real time by advanced Deep Neural Networks (DNN) which govern over the actual maneuvers of the car (i.e., steering, braking, acceleration). In order to manage a wide variety of driving scenarios, SDCs necessitate a large amount of driving data, combining nominal and adversarial scenarios [6].

To date, it is still unlikely for a DNN to generalize correctly to the plethora of driving situations met everyday by human drivers. As such, a component monitoring the *confidence* of the DNN may promptly detect when the SDC is entering a low-confidence zone, and activate healing strategies that bring the vehicle to a safe state.

In self-driving cars, depending on the level of autonomy, the self-healing procedure can either involve the human driver, or can be delegated to an automated system.¹ At both levels, early and accurate misbehaviour prediction is an essential precondition to enable safe healing, and an overall pleasant driving experience.

Confidence Measures in DNNs. The prediction must consider DNN uncertainties originating from the measurements in response to possible adverse environmental conditions in which the SDC operates. The confidence level of a SDC can be measured through white-box or black-box techniques. White-box monitoring of the internal behaviour of a DNN component. For simple classifiers, measuring softmax probabilities, or information theoretic metrics such as entropy, and mutual information [27] may suffice. For more complex networks such as those that operate on a SDC, softmax probabilities and entropy are known to be unreliable confidence estimators [41]. Moreover, white-box metrics require a transparent access to the network, and substantial domain-knowledge for the creation of nontrivial probabilistic models that approximate the network’s uncertainty.

Black-box techniques, differently, model the SDC uncertainty by monitoring the relation between the current input (images) and the input data used during training. For instance, consider a SDC which has been trained only with images representing highways. If images representing a narrow city street are given to the DNN, the model will still output steering angles, but ideally we would like to warn the SDC of a drop in the confidence level. The main advantages of black-box confidence metrics consist in being independent from the specific SDC architecture, requiring no modifications to the existing DNN model because they use information which is readily available for analysis, and in being, therefore, highly generalizable. In this paper, we focus on black-box confidence estimation, because softmax probabilities and entropy are known to be unreliable confidence estimators to model complex DNN networks [41]. We next describe autoencoders and time series analysis, which are the main building blocks of our approach.

Autoencoders. An autoencoder (AE) is a DNN designed to reconstruct its own input. It consists of two sequentially connected components (an encoder, and a decoder) that are arranged symmetrically. The simplest form of autoencoder (SAE) is a three-layer DNN: the input layer, the hidden layer, and the output layer. The hidden layer encodes any given input $x \in \mathcal{R}^D$ to its internal representation (*code*) $z \in \mathcal{R}^Z$

¹<https://www.nhtsa.gov/technology-innovation/automated-vehicles-safety>

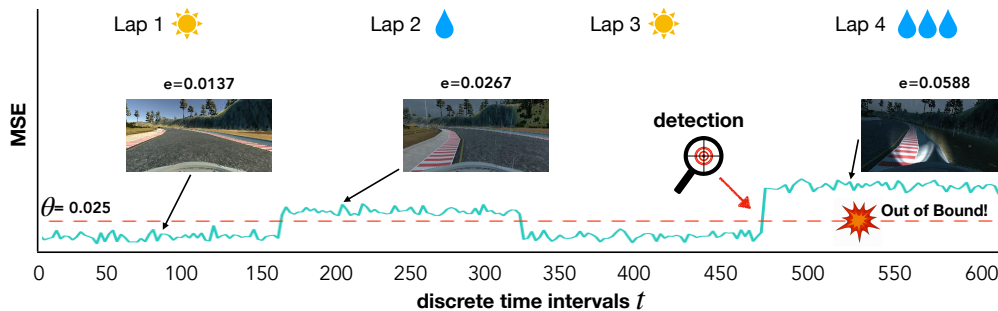


Figure 1: Confidence levels of NVIDIA’s DAVE-2 [7] in response to changing driving scenarios. The picture shows frames captured during the execution of the SDC along one of our testing tracks, under different conditions: (Lap 1) sun, (Lap 2) light rain, (Lap 3) sun, and (Lap 4) heavy rain. The picture juxtaposes the reconstruction error by the anomaly detector of SELFORACLE, which is used as a proxy for the DNN confidence. We can notice that the reconstruction error is low when the car drives under sunny conditions (i.e., conditions similar to those observed during training), whereas the error increases moderately with adverse conditions (the car does no longer follow the center of the road), and grows above a given threshold when facing heavy rainy conditions at night time (which cause the SDC to drive off the road).

with a function $f(x) = z$. Usually $Z \ll D$. The output layer (decoder) decodes the encoded input with a reconstruction function $g(z) = x'$, where x' is the reconstructed input x . The autoencoder minimises a loss function $\mathcal{L}(x, g(f(x)))$, which measures the distance between the original data and its low-dimensional reconstruction. A widely used loss function in autoencoders is the Mean Squared Error (MSE).

The input and output layers of autoencoders have the same number of nodes. If multiple hidden layers are used, the architecture is referred to as deep autoencoder (DAE). The surge of novel kinds of DNNs has correspondingly produced variants of autoencoders based on such architectures. For example, convolutional autoencoders [21] allow learning powerful spatial-preserving relationship within images at a lower training time with respect to fully dense layers. Another interesting proposal are variational autoencoders (VAE) that are able to model the relationship between the latent variable z and the input variable x by learning the underlying probability distribution of observations using variational inference [3].

Time Series Analysis. Traditional feedforward DNNs assume that all inputs and outputs are independent of each other. However, learning temporal dependencies between inputs or outputs is important in tasks involving continuous streams of data. Thus, a predictive model can take advantage of information from the previous inputs/outputs, to enhance its predictive capability.

Time series analysis can be applied to the output sequence produced by a DNN to identify the trend and predict future values. Among the numerous models available for time series analysis, the most widely used ones are autoregressive (AR), integrated (I) and moving average (MA) models, along with their combinations. An AR model of order k predicts the next value x_t as a linear combination of past values x_{t-1}, \dots, x_{t-k} :

$$x_t = \alpha_0 + \sum_{i=1}^k \alpha_i x_{t-i} + \epsilon_t \quad (1)$$

where coefficients $\alpha_0, \dots, \alpha_k$ can be estimated by the least square method and ϵ_t represents the error term. Processing of a sequence of inputs can be achieved by recurrent neural networks (RNN) equipped with long short-term memory (LSTM) [15], which is capable of dealing with both short and long range dependencies. In LSTM, outputs are influenced not only by the current input but also by the state of the RNN, which encodes the entire history of past inputs. A *forget gate* layer of LSTM decides what information to keep/remove from the previous network state whereas an *additive gate* decides how to update the state based on the current input.

3 Problem formulation

We focus on SDCs that perform *behavioural cloning*, i.e., the DNN learns the *lane keeping* [13] behaviour from a human driver. Models such as the ones by NVIDIA [7] or the Udacity self-driving challenge [38] are

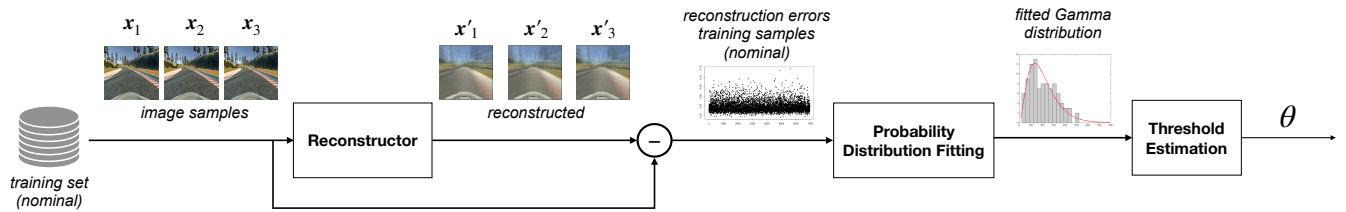


Figure 2: Model Training under Nominal Driving Behaviour.

trained with visual inputs (i.e., images) from car-mounted cameras that record the driving scene, paired with the steering angles from the human driver. The DNN then “learns how to drive” by discovering underlying patterns within the training images representing the shape of the road, and predicting the corresponding steering angle.

For a classification problem (e.g., hand-written digit recognition), a misbehaviour can be defined as an input that can be confidently labeled by humans while it is misclassified by a DNN. Differently, for a regression problem such a definition is more troublesome, because there is no expected outcome for an individual output of the DNN and it is only the overall behaviour resulting from the DNN predictions that may or may not be acceptable, depending on the specific application domain.

In steering angle prediction, it is challenging to decide if the steering angle produced by a DNN is wrong, because the optimal steering angle is generally unknown for a new test scenario and even if it were known, the amount of difference between predicted and expected steering angle that qualifies as an error is difficult to decide a priori. It is instead more realistic to figure out whether a chain of inaccurate predictions, regardless of the amount of deviation from the optimal angles, ultimately leads to a misbehaviour of the system.

In fact, the definition of misbehaviour should be decoupled from the notion of correct/wrong DNN output, being instead linked to the ability of a DNN of abstracting from the training examples and learning how to drive in different ways/conditions. It is the task of the DNN to generalize the training knowledge to make the model robust w.r.t. slightly different conditions from those observed in the training set. For instance, if the SDC is facing a 90-degree bend, the DNN may decide to steer at a certain steering angle θ if the car speed is 60 mph, or $\frac{\theta}{2}$ if the speed is 30 mph.

Definition 1 (Misbehaviour of a DNN) *A DNN exhibits a misbehaviour in a given test scenario if the overall system that contains the DNN does not respect its requirements due to the outputs produced by the DNN.*

In the autonomous driving domain, there are many possible misbehaviours, associated with the different requirements that such systems are supposed to realize. Safety requirement violations are by far the most critical requirements, as a misbehaviour in the steering component may cause a crash of the vehicle with potential casualties. However, in general, a SDC might violate also other driving requirements, e.g., related to ride comfort [8], such as excessive derivative of the steering angle, unstable movement around the centerline, or excessive deceleration. All these could also be considered as misbehaviours by our definition.

In this paper, we focus on the prediction of two safety-critical misbehaviours: (1) collisions, and (2) out-of-bound episodes (OBES). The rationale for this choice are as follows. First, they represent the vital requirement to be satisfied and thoroughly tested (i.e., the car should stay in lane and avoid whatsoever collision), without which autonomous driving vehicles would be hardly accepted in production. Moreover, leveraging a simulation environment such as Udacity’s [37], allows us to: (1) safely test such critical scenarios, (2) precisely define, observe and measure them, in order to support crash analysis and reproduction.

Thus, the *problems* we want to address in this paper are: (1) recognizing when a self-driving car is within a low-confidence area because of an unexpected execution context, and (2) predicting such situation timely enough so as to take countermeasures before the vehicle may crash or drive off road.

4 Approach

The goal of our approach is to monitor the confidence level of a SDC as it runs and to promptly predict whether drops in confidence correlate with potential future misbehaviours. Our approach works in an end-to-end fashion, analyzing directly the input data as retrieved by the car (an image from the center

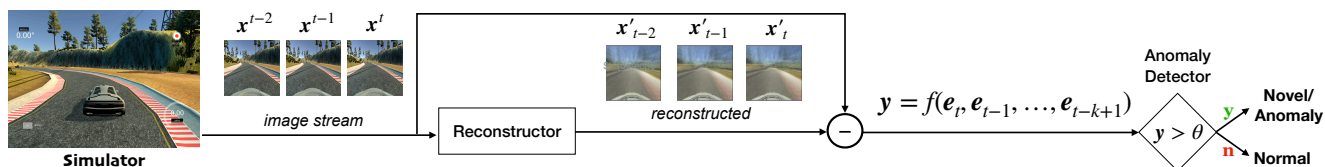


Figure 3: Usage Scenario of SELFORACLE.

camera), making the approach independent from the specific architecture of the self-driving component, requiring no modifications to the existing DNN model, and being therefore, highly generalizable.

The main working assumption is that a prediction model trained on normal data should learn the normal time series patterns. When the model is used on a SDC in the field, it should worsen its performance as the car approaches previously unseen regions as compared to normal, known regions (Figure 1). Then, what's needed is a way to set up a good decision boundary in order to timely alert the human driver (NHTSA Level 4) or the main self-driving component (NHTSA Level 5) about triggering self-healing.

We now detail each step of our approach, which consists of two main phases: (1) *model training under nominal driving behaviour*, and (2) *field usage of the trained model*.

4.1 Training of SELFORACLE under Nominal Driving Behaviour

Figure 2 illustrates the training phase of our approach, which consists of several steps.

4.1.1 Reconstructor

The first step consists in retrieving a model of normality from the training driving scenarios. Thus, in the training set, we capture the visual input stream of the SDC under nominal situations.

Then, we train our driving scenario reconstructor with such “normal” instances. Let us consider a training set $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$ of n image frames, where the index $i \in [1 : n]$ of $x_i \in \mathbf{X}$ represents the discrete time t .

Depending on the considered architecture, a reconstructor can be *singled-image* or *sequence-based*. For singled-image reconstructors, only one image frame is considered at a time. When the discrete time is $t = i$, x_i is the input and the reconstructor recreates it into x'_i . For sequence-based reconstructors, assuming k image frames preceding x_i are used to reconstruct x_i , the sequence $\langle x_{i-k}, \dots, x_{i-1} \rangle$ is the input used to output x'_i , a prediction of the actual current frame x_i . For instance, for $k = 3$ and $i = 4$, the reconstructor considers the sequence $\langle x_1, x_2, x_3 \rangle$ in order to predict the current frame x_4 .

At the end of this task, each reconstruction error $e_i = d(x_i, x'_i)$ can be computed, where d is a proper distance function (e.g., Euclidean distance). This results in the set of reconstruction errors $\mathbf{E} = \{e_1, e_2, \dots, e_n\}$, available for all elements in the training set $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$.

4.1.2 Probability Distribution Fitting

We build a model of normality for the reconstruction errors $\mathbf{E} = \{e_1, e_2, \dots, e_n\}$ collected in nominal driving conditions. We use probability distribution fitting to obtain a statistical model of normality and using such model we determine a threshold θ that brings the expected false alarm rate in nominal conditions below some acceptable, configurable level ϵ (e.g., $\epsilon = 10^{-3}$ or $\epsilon = 10^{-4}$).

The reconstruction error $e = d(x, x')$ can be computed by comparing the individual pixels of the images x and x' and taking the mean pixel-wise squared error. Assuming images have width W , height H and C channels (usually, 3 RGB channels for colour images), the reconstruction error is defined as follows:

$$d(x, x') = \frac{1}{WHC} \sum_{i=1, j=1, c=1}^{W, H, C} (x[c][i, j] - x'[c][i, j])^2 \quad (2)$$

With some reasonable approximation, we can assume that the pixel-wise error $e[c][i, j] = x[c][i, j] - x'[c][i, j]$ follows a normal distribution with pixel dependent variance: $e[c][i, j] \sim \mathcal{N}(0, \sigma_{c,i,j})$. Correspondingly, the sum of the squares of pixel-wise errors $e[c][i, j]$ will be a gamma distribution: $e = d(x, x') \sim \Gamma(\alpha, \beta)$. We get

a gamma instead of a χ^2 distribution because pixel-wise errors have different (channel/pixel dependent) and non-unitary variances.

Definition of Gamma Distribution. Gamma is a probability model for a continuous variable on $[0, \infty)$ which is widely used in engineering, science, and business, to model continuous variables that are always positive and have skewed distributions.

The *probability density function* of a random variable $x \sim \Gamma(\alpha, \beta)$ is:

$$f(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \quad x > 0; \alpha, \beta > 0 \quad (3)$$

where α is the *shape* parameter (which affects the shape of the distribution), β is the *rate* parameter (or inverse scale, which stretches/shrinks the distribution) and Γ is the gamma function. When α is large, the gamma distribution closely approximates a normal distribution with the advantage that the gamma distribution has non-zero density only for positive real numbers.

The *gamma function* Γ can be seen as a solution to the interpolation problem of finding a smooth curve that connects the points (n, m) with $m = (n - 1)!$ at any positive integer value for n . Such a definition was extended to all complex numbers with a positive real part by Bernoulli, as a solution to the following integral:

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx \quad \Re(z) > 0; \quad (4)$$

Fitting the Gamma Distribution. One powerful method to estimate the parameters of a (gamma) distribution that fit the data (reconstruction errors) best is by maximum likelihood estimation (MLE). The likelihood function reverses the roles of the variables: in Equation 3, the values of x are known, so they are the fixed constants, whereas the unknown variables are the parameters α and β . MLE involves calculating the values of these parameters so as to obtain the highest likelihood of observing the values of x when the given parameters are supplied to f .

Under the assumption of independence of the data, the likelihood of the data given the parameters of the distribution is conveniently defined as the logarithm of the joint probability of the data for a given choice of the parameters. In the case of a gamma distribution, with a dataset consisting of reconstruction errors $\mathbf{E} = \{e_1, e_2, \dots, e_n\}$, we get:

$$\mathcal{L}(\alpha, \beta; \mathbf{E}) = \frac{1}{n} \sum_{i=1}^n \log f(e_i | \alpha, \beta) = \quad (5)$$

$$n(\alpha - 1) \overline{\log e} - n \log \Gamma(\alpha) - n\alpha \log \beta - n\bar{e}/\beta \quad (6)$$

where \bar{e} ($\overline{\log e}$) is the mean (log) reconstruction error over \mathbf{E} . To find the values of parameters α and β that maximize \mathcal{L} we have to find a solution to the equations: (1) $\partial \mathcal{L} / \partial \alpha = 0$; (2) $\partial \mathcal{L} / \partial \beta = 0$. The second equation can be easily solved analytically, resulting in $\beta = \bar{e} / \alpha$. By substituting the value of β into the first equation we get an equation that unfortunately cannot be solved analytically. However, the Newton method can iteratively converge to the solution quite quickly. The output of such numerical estimation will be the pair of parameters α and β of the gamma distribution that best fit the reconstruction errors.

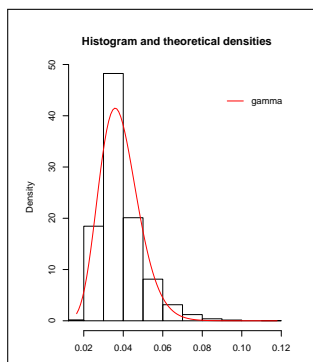


Figure 4: Fitted Gamma distribution of reconstruction errors from a VAE on the Dave2 dataset.

Example of Threshold Estimation. Let us consider a set of reconstruction errors. Figure 4 shows the histogram of those produced on the Dave2 dataset when the reconstructor is a VAE. On such dataset the MLE method estimates the following gamma parameters: $\alpha = 15$; $\beta = 392$. With these parameter values, the plot of the gamma distribution is the red line in Figure 4. We can notice a good agreement between the histogram and the estimated probability distribution.

Let us assume that we are willing to accept a false alarm rate $\epsilon = 10^{-2}$. The threshold θ with a probability mass above the threshold equal to 10^{-2} can be easily obtained as the inverse of the cumulative gamma distribution $F(x)$: $\theta = F^{-1}(1 - \epsilon)$. This ensures that the cumulative probability of values $\leq \theta$ is $1 - \epsilon$, leaving only a probability of ϵ to the tail following θ . We use the estimated θ as threshold to distinguish anomalous conditions (reconstruction error $\geq \theta$) from normal ones (reconstruction error $< \theta$).

4.2 Usage Scenario

Figure 3 shows how SELFORACLE is used online for misbehaviour prediction after model training (i.e., after fitting the gamma distribution and estimating the threshold θ). Misbehaviour prediction is executed online as the SDC drives. In this phase, the SDC generates data continuously and the reconstructor recreates the incoming stream of images. The sequence of reconstruction errors is passed through the autoregressive model f and the resulting, filtered error is compared against the threshold θ , which determines whether an anomaly is detected or not. In the former case, self-healing is triggered and the SDC is brought to a safe state.

4.2.1 Time-aware Anomaly Score Prediction

The reconstruction error e_t at time t might be susceptible to single-frame outliers, which are not expected to have a big impact on the driving of the car, but would indeed make the misbehaviour predictor falsely report an anomalous context. For this reason, we smooth such noisy oscillations by applying an autoregressive filter: the sequence of reconstruction errors is passed to a module that performs time series-analysis. In Figure 3, this corresponds to the AR filter f . The output of this filter, instead of the raw reconstruction error e_t , is compared with the threshold θ to recognize unexpected driving conditions. In our experiments we used a simple AR model (see Equation 1) with $\alpha_0 = 0$ and $\alpha_i = 1/k$ for $i = 1, \dots, k$.

5 Empirical Evaluation

5.1 Research Questions

We consider the following research questions:

RQ₁ (effectiveness): How effective is SELFORACLE in predicting anomalies for autonomous vehicles? What are the best reconstructors to use?

RQ₂ (prediction): How does the misbehaviour predictions of SELFORACLE change as we increase the reaction period (i.e., we anticipate the time of prediction)?

RQ₃ (comparison): How does SELFORACLE compare with DeepRoad’s [42] online input validation?

5.2 Self Driving Car Models

We evaluate our framework on three existing DNN-based SDCs: NVIDIA’s DAVE-2 [7], Epoch [31], and Chauffeur [30]. We choose these models because they are publicly available, thus they can be trained and evaluated on the simulator. Moreover, they have been objects of study of other testing works [26, 34]. DAVE-2 consists of three CNNs, followed by five fully-connected layers. Chauffeur uses a CNN to extract the features of input images, and a RNN to predict the steering angle from 100 previous consecutive images. Epoch consists of a single CNN model.

5.3 Simulation Platform

A major problem in anomaly detection research is the lack of labeled benchmark datasets [9], and the self-driving car domain is no exception. Unlike previous works [42, 34, 26], we cannot rely on existing driving

image datasets such as the ones released by Udacity [39], because they lack any episode of crash, or cars driving off road whatsoever. Moreover, our definition of misbehaviour (Section 3) requires the creation of a set of “controllable” unexpected conditions that may potentially cause them, along with a way to precisely record them. Thus, to investigate the effectiveness of our approach in predicting safety-critical misbehaviours, we evaluated SELFORACLE in the Udacity simulator [37].

The Udacity simulator is developed with Unity [40], a popular cross-platform game engine. The simulator provides two default tracks, for testing DNNs models. The simulator executes in two modes: (1) *training mode*, in which the user manually controls the car while the simulator records her actions, and (2) *autonomous mode*, in which the car is controlled by an external agent, such as a DNN-based autonomous driving system. Moreover, we added a third track [36] to the existing set, and we implemented two additional components, namely, an *unexpected context generator*, and a *collision/OBE detection system*.

5.3.1 Unexpected Context Generator

First, we developed a method to gradually inject unseen conditions during testing mode (i.e., conditions diverse from the training mode’s defaults). The first condition deals with illumination and introduces a *day/night cycle* component to gradually change the light condition of the track during the simulation. The effect is customizable and consists of smooth increase/decrease of brightness/darkness over a fixed period (60 s in our experiments). The second kind of unexpected condition deals with *weather*. We implemented rain, snow and mist effects, with a variable intensity during the simulation. For the implementation, we used a specific Unity component called Particle System [25] which can simulate the physics of a cluster of particles with high performance. Rain particles emission rate ranges between a minimum of 100 (light rain) to a maximum of 10,000 particles/s (heavy rain); fog between [100..2,000] particles/s, and snow between [100..800] particles/s. Figure 5 shows a few examples.

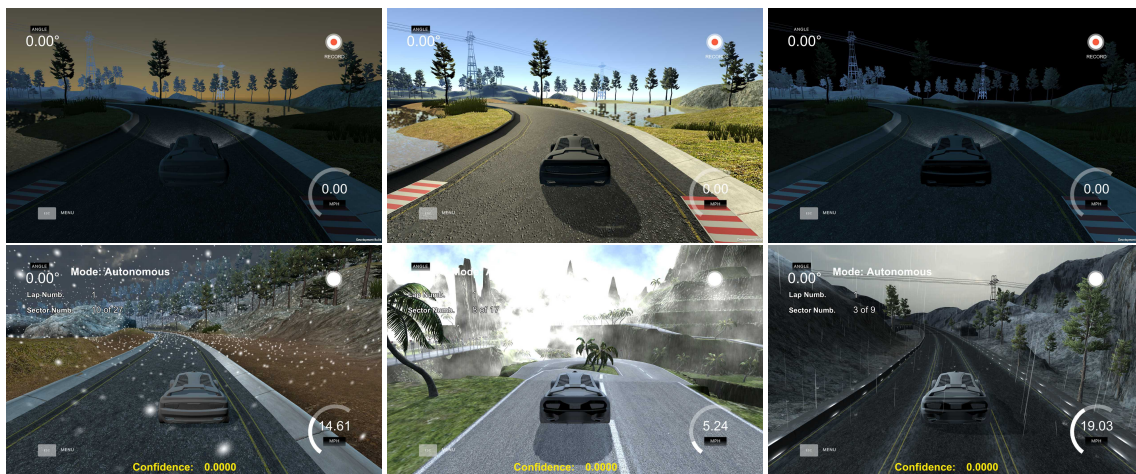


Figure 5: (top) Day/night cycle (sunrise, day, night) and (bottom) weather effects (snowy Lake Track, foggy Jungle Track, and rainy Mountain Track).

5.3.2 Collision/OBE Detection System

Following our definition of safety-critical misbehaviours, we implemented an automated collision/OBE detection system (ACODS) that records any unwanted interaction of the SDC with the environment, allowing us to experiment the effectiveness of SELFORACLE at anticipating such episodes during the occurrence of unexpected scenarios (Section 5.3.1).

We implemented ACODS based on *colliders*, which are consolidated building blocks of modern game engines to simulate the physical interaction between objects. We approximate the car body with a geometry mesh, and implemented a collider callback that informs the simulator of any physical interaction of the car with scene objects. When the car “hits” the road, then it means the car is actually on track, whereas when the car “collides” with any other object, the callback registers whether it is a crash against some object (Figure 6 (left)), or whether it is an OBE (Figure 6 (right)).

Secondly, we implemented an automatic restart mechanism that restores the SDC to a safe position after a crash/OBE, allowing us to record multiple simulations without the need for manual restart.



Figure 6: Simulator crash/OBE detection.

5.4 Procedure

5.4.1 Data Generation (Training Set)

Training data were collected by the authors in training mode by performing 10 laps on each track, following two different track orientations (normal, reverse). Overall, we obtained a dataset of 124,638 training images (at 10-13 fps), divided as follows: 32,243 for Track 1 (Lake), 51,422 for Track 2 (Jungle), and 40,973 for Track 3 (Mountain). Differences depend on the track lengths. To allow a smooth driving and a correct behaviour capture (i.e., lane keeping), the maximum driving speed was set to 30 mph, the default in the Udacity simulator.

5.4.2 SDC Model Setup & Training

All SDCs models were trained on 41,546 images from the central camera. We used data augmentation as a consolidated practice for building more reliable and generalizable SDCs, limiting the lack of image diversity in the training data. Specifically, 60% of the data was augmented through different image transformation techniques (e.g., flipping, translation, shadowing, brightness). We cropped the images to 66 x 200, and converted them from RGB to YUV colour space. All SDC models were trained for 500 epochs with batch size of 256 on a machine featuring an i9 processor, 32 GB of memory, and an Nvidia GPU 2080 TI with 11GB of memory. Basically, the training was meant to create solid models for testing, i.e., able to drive multiple laps on each track under nominal conditions without showing any misbehaviour in terms of crash/OBE.

5.4.3 Evaluation Set

To collect the evaluation data, we executed 72 simulations (2 laps each) in autonomous mode (3 SDC x 8 conditions x 3 tracks). As in data generation, the maximum speed was set to 30 mph. Specifically, for each SDC and for each track, we performed 1 simulation in the same normal conditions as the training set. This allows us to estimate the number of false alarms (false positives) in nominal conditions. Second, we performed 4 simulations activating in turn a single unexpected condition: day/night cycle, rain, snow, fog. Third, we performed 3 simulations activating in turn a combined condition: day/night cycle + rain, day/night cycle + snow, day/night cycle + fog.

In our experiments, we used a value of 60 s both for the day/night cycle and for the loop between the minimum and maximum intensity of the effects. This value was chosen empirically given the relatively short speed of the SDC, and the small length of the tracks, and allowed us to test the behaviour of the SDC on each of the track subsets under all possible conditions. For example, the bridge part of Track1 (Lake) has been driven on under both dawn/day/sunset/night conditions (day/night cycle) and minimal/maximal intensity of rain/fog/snow.

Overall, we obtained a dataset of 778,592 images, divided as follows: 188,032 for Track 1 (Lake), 260,064 for Track 2 (Jungle), and 208,064 for Track 3 (Mountain) with unknown conditions and 33,088 for Track 1, 46,272 for Track 2 and 43,072 for Track 3 in known conditions.

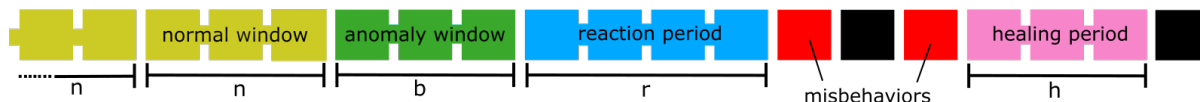


Figure 7: Labelling of Anomalous and Normal Windows in Driving Stream.

5.4.4 SELFORACLE’s Configurations

We used four autoencoders taken from existing guidelines [16]: (1) *SAE* (simple autoencoder with a single hidden layer), (2) *DAE* (deep five layers fully-connected autoencoder), (3) *CAE* (convolutional autoencoder alternating convolutional and max-pooling layers), and (4) *VAE* (*variational autoencoder*).

All autoencoders take as input a single image. We added images taken by the side cameras of the car (left, right) to allow better generalization, even though, during testing, autoencoders are not used for prediction. Lastly, we performed further data augmentation on 60% of the inputs, as described in Section 5.4.2.

As an additional sequence-based reconstructor, we also implemented an LSTM consisting of two LSTM-layers and one convolutional layer. For implementation details, we made our code publicly available in the replication package accompanying this paper [35].

Baseline. We use the input validation technique of DeepRoad [42] as baseline for SELFORACLE. Unfortunately, authors did not make their code available; therefore, we implemented our own version according to the description in the paper. For input validation, DeepRoad uses the pre-trained VGG19 ImageNet classifier [28] to extract style and feature vectors from a given image. Principal Component Analysis (PCA) is then used to reduce all style and feature vectors, concatenated into a matrix, to three dimensional representations, which support distance/similarity estimation. To allow a fair comparison, we integrated it within SELFORACLE as reconstructor. However, unlike autoencoders, DeepRoad is computationally very expensive and memory demanding (due to the size of the matrix supplied to PCA). In the paper, authors reduced their training set to 600 images, which were resized to 120x90. During input validation, the three dimensional representation of an online input image is compared to the nominal images by measuring the average of the top-100 minimum distances from the training set.

Our own implementation relaxed the restrictions above by considering a training set consisting of 3,000 randomly sampled images (i.e., $5\times$ improvement w.r.t. the original implementation described in the paper), resized to 224x224, which is the default input size for VGG19 [17]. Keeping the fraction of the training set ($\frac{1}{6}$) constant, we compute similarity based on the average of the top-500 minimal distances. For implementation details, we made our code publicly available in the replication package accompanying this paper [35].

5.5 Evaluation of Simulation Results

To allow a fair comparison with our baseline DeepRoad, we evaluated all approaches offline, by splitting the evaluation set of recorded images in *windows of consecutive frames*, which we labelled as either anomalous or normal (Figure 7). In anomalous windows, SELFORACLE is expected to predict the shortly-following misbehavior.

5.5.1 Labelling of Anomaly and Normal Windows in Evaluation Data

Let $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$ be the sequence of considered images (frames). Misbehaviors are represented as $m_j \in \{0, 1\}$, where $m_j = 1$ iff a misbehavior is recorded at $x_j \in \mathbf{X}$. We define a healing period as a sequence of h misbehaviour-free frames following a misbehaviour at time t . We define a reaction period as a sequence of r misbehaviour-free frames preceding a misbehaviour at time t' and not intersecting any healing period. We define an *anomalous window* as a consecutive misbehaviour-free frames followed by a reaction period that does not intersect any healing period. We define a *normal window* as b consecutive misbehaviour-free frames followed by an anomalous window, or a normal window that does not intersect any healing period. This is illustrated graphically in Figure 7. Formally, assuming any misbehavior recorded at time $t \in [1:n]$, i.e., $m_t = 1$:

- window x_{t+1} to x_{t+h} is labelled as healing period if $m_t = 1$ and $m_j = 0 \forall j \in [t+1:t+h]$;
- furthermore, the window x_{t+1} to x_{k-1} is also labelled as healing period if $m_k = 1$ with $k > t$ and $k - t - 1 < h$, and $m_j = 0 \forall j \in [t+1:k-1]$;

- window x_{t-r} to x_{t-1} is labelled as reaction period iff $m_t = 1, m_j = 0 \forall j \in [t-r:t-1]$ and no healing period contains any frame from x_{t-r} to x_{t-1} ;
- window x_i to x_{i+a-1} is labelled as anomaly window iff a reaction period starts at $x_{i+a}, m_j = 0 \forall j \in [i:i+a-1]$ and no healing period contains any frame from x_i to x_{i+a-1} ;
- window x_i to x_{i+b-1} is labelled as normal window iff an anomaly or a normal window starts at $x_{i+b}, m_j = 0 \forall j \in [i:i+b-1]$ and no healing period contains any frame from x_i to x_{i+b-1} .

Moreover, if $m_j = 0$ for all $j \in [k:n]$, all consecutive windows of size b starting within x_k to $x_{n-r-a-1}$ which do not intersect any healing period are labelled as normal. The labelling described above ensures that after the last misbehavior in a sequence, h healing images are ignored (i.e., not labeled) before another anomaly or normal window is defined. h must be chosen high enough that the car is back safely on the road when the next windows are labelled. Furthermore, r images occur in between an anomaly window in which the system is supposed to predict the upcoming misbehavior, and the actual misbehavior. This period would, in practice, be used by the self-healing system to execute countermeasures against the predicted future misbehavior. Intuitively, misbehavior prediction is expected to be much harder as the value of r increases. In our experiments, we set the value of $n = b = 30$ frames (i.e., normal/anomalous windows), which is $\approx 3s$.² The size of the healing window was set to $h = 60$ ($\approx 5s$) frames, and the size of the reaction window to $r = 50$ ($\approx 4s$) frames.

5.5.2 Metrics used for Analysis

If the loss score for an image is higher than the automatically estimated threshold θ (Section 4.2), SELFORACLE triggers an alarm. Consequently, a true positive is defined when SELFORACLE triggers an alarm during an anomalous window, early enough to predict a misbehavior. Conversely, a false negative occurs when SELFORACLE does not trigger an alarm during an anomalous window, thus failing at predicting a misbehaviour in time for triggering self-healing. A false positive represents a false alarm by SELFORACLE, whereas true negative cases occur when SELFORACLE detects correct detection of normality.

We assume that a single alarm immediately starts the self healing system, such that multiple consecutive alarms within the healing time have no effect. Correspondingly, once a FP occurs and the self healing system is running, additional consecutive FP windows have no effect in practice and are thus excluded from our analysis.

Our goal is to achieve (1) high recall, or true positive rate (TPR, defined as $TP/TP+FN$), i.e., true alarms, while (2) minimizing the complement of specificity, or false positive rate (FPR, defined as $FP/TN+FP$), i.e., labelling safe situations as unsafe. We are also interested in F1-score ($F_1 = 2 \cdot \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$) because, in practice, it is informative to have a high F1-score at a given threshold.

We also consider two widely adopted threshold-independent metrics for evaluating classifiers at various thresholds settings such as AUC-ROC (area under the curve of the Receiver Operating Characteristics), and AUC-PRC (area under the Precision-Recall curve).

5.6 Results

Effectiveness (RQ₁). Table 1 presents the effectiveness results on a per-SDC model basis. Columns 2 and 3 show threshold-independent measures of AUC-PRC and AUC-ROC. The remaining of the table shows the effectiveness metrics across two confidence thresholds that we found interesting for analysis and correspond to $\epsilon = 0.05$ and $\epsilon = 0.01$ (at lower values of ϵ , both FPR and TPR get close to zero, making the misbehaviour predictor useless).

Overall, LSTM and VAE are the best performing reconstructors on the AUC-PRC and AUC-ROC metrics. Columns 12 and 21 (Nominal/FPR) show FPR under conditions similar to those of the training set. Values are almost always near to zero and occasionally even equal to zero (this is indicated by omitting the decimals) across the variants of SELFORACLE. This is an empirical validation of the accuracy of the gamma distribution as a statistical model for the reconstruction errors. In fact, at $\epsilon = 0.05$ most FPR reported in column 12 are very close to the theoretical value, 0.05 (see, e.g., rows under *Totals*). At $\epsilon = 0.01$ some values drop to zero. This means that in those configurations SELFORACLE will raise no false alarm when the SDC drives in nominal conditions. For instance, with both thresholds, LSTM never raised false alarms within the 23,728 considered normal windows .

²In our setting, Udacity frame rate was approximately 10/12 fps.

Table 1: Evaluation Results for all variants of SELFORACLE across all SDCs. Best strategies are highlighted.

	AUC-PRC \uparrow	AUC-ROC \uparrow	$\epsilon = 0.05, 1 - \epsilon = 0.95$										$\epsilon = 0.01, 1 - \epsilon = 0.99$							
			Unexpected					Nominal					Unexpected				Nominal			
			TP	FP	TN	FN	TPR \uparrow	FPR \downarrow	F1	Prec.	FPR \downarrow	TP	FP	TN	FN	TPR \uparrow	FPR \downarrow	F1	Prec.	FPR \downarrow
DAVE-2																				
VAE	0.354	0.902	149	304	1,970	47	0.760	0.134	0.459	0.329	0.042	107	183	2,959	89	0.546	0.058	0.440	0.369	0.003
DAE	0.330	0.891	104	210	2,679	92	0.531	0.073	0.408	0.331	0.041	21	55	4,063	175	0.107	0.013	0.154	0.276	0.010
SAE	0.336	0.891	138	260	1,877	58	0.704	0.122	0.465	0.347	0.050	108	183	2,665	88	0.551	0.064	0.444	0.371	0.002
CAE	0.290	0.821	6	22	4,208	190	0.031	0.005	0.054	0.214	0.024	0	0	4,282	196	0	0	n.a.	n.a.	0
LSTM	0.357	0.903	16	34	3,990	177	0.083	0.008	0.132	0.320	0	7	12	4,119	186	0.036	0.003	0.066	0.368	0
DeepRoad	0.198	0.780	65	344	3,170	131	0.332	0.098	0.215	0.159	0.054	44	250	3,651	152	0.225	0.064	0.180	0.150	0.037
Epoch																				
VAE	0.391	0.904	158	331	1,952	51	0.756	0.145	0.453	0.323	0.049	106	169	2,858	103	0.507	0.056	0.438	0.386	0.001
DAE	0.399	0.895	112	188	2,720	97	0.536	0.065	0.440	0.373	0.042	24	34	3,653	185	0.115	0.009	0.180	0.414	0.010
SAE	0.386	0.883	147	284	2,026	62	0.703	0.123	0.459	0.341	0.050	120	175	2,838	89	0.574	0.058	0.476	0.407	0.002
CAE	0.310	0.822	6	23	3,661	203	0.029	0.006	0.050	0.207	0.020	0	0	3,731	209	0	0	n.a.	n.a.	0
LSTM	0.385	0.879	23	34	3,503	175	0.116	0.010	0.180	0.404	0.001	7	13	3,592	191	0.035	0.004	0.064	0.350	0
DeepRoad	0.213	0.807	70	308	2,917	139	0.335	0.096	0.239	0.185	0.053	43	201	3,240	166	0.206	0.058	0.190	0.176	0.040
Chauffeur																				
VAE	0.242	0.951	98	392	3,700	23	0.810	0.096	0.321	0.200	0.049	81	267	5,391	40	0.669	0.047	0.345	0.233	0.002
DAE	0.203	0.944	78	281	5,045	43	0.645	0.053	0.325	0.217	0.051	13	95	7,730	108	0.107	0.012	0.114	0.120	0.009
SAE	0.241	0.931	96	354	3,650	25	0.793	0.088	0.336	0.213	0.056	86	240	5,177	35	0.711	0.044	0.385	0.264	0.003
CAE	0.172	0.909	7	34	8,127	114	0.058	0.004	0.086	0.171	0.023	0	0	8,229	121	0	0	n.a.	n.a.	0
LSTM	0.240	0.945	11	41	7,879	111	0.090	0.005	0.126	0.212	0	4	12	8,035	118	0.033	0.002	0.058	0.250	0
DeepRoad	0.098	0.797	37	594	5,564	84	0.306	0.097	0.098	0.059	0.055	23	458	6,551	98	0.190	0.065	0.076	0.048	0.042
Totals																				
VAE	0.320	0.924	405	1027	7,622	121	0.770	0.119	0.414	0.283	0.046	294	619	11,208	232	0.559	0.052	0.409	0.322	0.002
DAE	0.301	0.911	294	679	10,444	232	0.559	0.061	0.392	0.302	0.045	58	184	15,446	468	0.110	0.012	0.151	0.240	0.009
SAE	0.312	0.907	381	898	7,553	145	0.724	0.106	0.422	0.298	0.052	314	598	10,680	212	0.597	0.053	0.437	0.344	0.002
CAE	0.255	0.864	19	79	15,996	507	0.036	0.005	0.061	0.194	0.022	0	0	16,242	526	0	0	n.a.	n.a.	0
LSTM	0.329	0.915	50	109	15,372	463	0.098	0.007	0.149	0.315	0	18	37	15,746	495	0.035	0.002	0.063	0.327	0
DeepRoad	0.159	0.799	172	1246	11,651	354	0.327	0.097	0.177	0.121	0.054	110	909	13,442	416	0.209	0.063	0.142	0.108	0.040

In terms of TPR (to be maximized) and FPR (to be minimized), the best reconstructors are VAE and SAE, with comparable overall performance: 77%, 11% (TPR, FPR of VAE) and 72%, 10% (SAE) at $\epsilon = 0.05$; 55%, 5% (VAE) and 59%, 5% (SAE) at $\epsilon = 0.01$. It can be noticed that FPR is higher than ϵ (10% vs 5% and 5% vs 1%) with both reconstructors. This is expected, since we are measuring FPR in tracks with injected anomalies. These tracks differ substantially from the nominal tracks even in conditions not so extreme as to cause a misbehaviour, hence inflating the FPR a bit. However, we can notice that even in such non nominal conditions, the FPR remains low and not too far from the theoretical prediction ϵ .

We can **answer to RQ1** by noticing that with reconstructors VAE and SAE we achieve a FPR in nominal conditions close to the theoretical expectations (resp. 5% and 1%, in the two considered configurations); that in anomalous conditions FPR increases by a moderate amount (resp., +5% and +4%); and that the achieved TPR is quite high (with VAE, SAE, resp. 77%, 72% and 55%, 59%).

Prediction (RQ2). Figure 8 shows the PRC-AUC of the various configurations of SELFORACLE over different reaction periods. The general trend is that predictions get harder when the SDC is far from a critical scenario, having a longer reaction period to prevent the misbehavior, but quite surprisingly there is no drop in performance as we move away from the misbehaviour. Our explanation of this unexpected finding is that the tracks used for the evaluation of the approach contain always a relatively high degree of anomalous features, which might trigger a self-healing reaction. Occasionally, the level of detected anomalies surpasses the threshold and the misbehaviour predictor raises an alarm. Correspondingly, although slightly reduced, the signals of an upcoming misbehaviour exist in images quite far (even 60 frames or around 6s) from the misbehaviour.

We can **answer to RQ2** by noticing that the performance of SELFORACLE degrades smoothly as we anticipate the prediction (AUC-PRC remains quite high even 6s before the misbehaviour), but we should also remark that this result must be taken with care and might be partially due to the characteristics of the considered tracks, which contain a continuously and smoothly increasing degree of injected anomalies by design.

Comparison (RQ3). In our experiments, SELFORACLE is constantly superior to DeepRoad at predicting misbehaviours. Results of AUC-PRC and AUC-ROC show significant improvements across all thresholds, and regardless of the technique being used and the reaction period considered (in Figure 8, DeepRoad is the lowest curve). With $\epsilon = 0.05$ (resp. $\epsilon = 0.01$), VAE and SAE (see Table 1) expose more than twice the misbehaviours exposed by DeepRoad, with a TPR = 77%, 72% vs 32% (resp., 55%, 59% vs 20%), at comparable false positive rate FPR = 11%, 10% vs 9% (resp., 5% vs 6%).

While reimplementing DeepRoad, we noticed its high computational cost that makes it quite unsuitable for online misbehaviour prediction. SDC manufacturers use much larger training datasets than the one we used for our empirical study and differently from autoencoders such as VAE and SAE, DeepRoad is quite sensitive to the size of the training set, which must be sub-sampled dramatically to make the approach

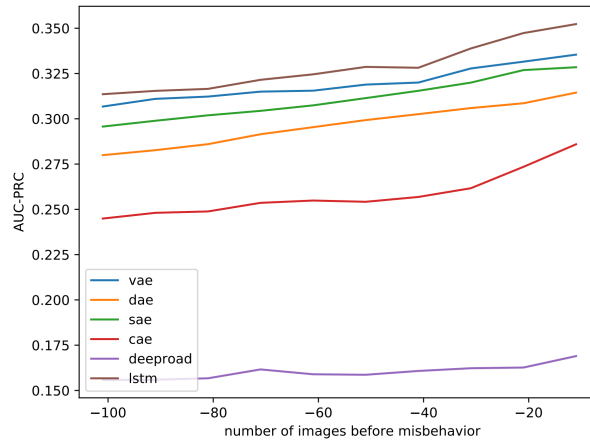


Figure 8: Misbehavior prediction capability over time.

applicable. On the contrary, autoencoders seem a very promising option, given their relatively simple architecture. In particular, SAE is very efficient, yet it has comparable performance as the more sophisticated VAE. Hence, as an input validation framework, DeepRoad has shown to be both computationally very expensive and inaccurate, which makes it less promising and desirable than SELFORACLE for online misbehaviour prediction.

Our **answer to RQ3** is that SELFORACLE outperforms DeepRoad in all respects: computational cost, accuracy of misbehaviour prediction (see TPR) and minimization of false alarms (see FPR and AUC-PRC).

5.7 Threats to Validity

Internal validity. We compared all variants of SELFORACLE and DeepRoad under identical parameter settings, and on the same evaluation set. The main threat to internal validity concerns our custom implementation of unexpected conditions within the simulator. However, this was a mandatory choice, since we are not aware of open source driving simulators that can inject unexpected execution contexts in a controllable way. Another possible threat may be the choice and the training of our own SDCs, which may exhibit a large number of misbehaviour if trained inadequately. We mitigated this threat by training and fine-tuning the best publicly available driving models. Our own implementation of DeepRoad may be another threat to internal validity, that we mitigated by developing an implementation which improves the original one by processing $5\times$ more information.

External validity. We used a limited number of self-driving systems in our evaluation, as well as tracks, which pose a threat in terms of generalizability of our results. We tried to mitigate this threat by choosing popular subjects developed with real-world frameworks (e.g, Keras).

Reproducibility. All our results, the source code of SELFORACLE, the simulator, and all subjects are available [35], making the evaluation repeatable and our results reproducible.

6 Related Work

Adversarial Input Generation. Adversarial generation approaches aim at generating inputs that trigger inconsistencies between multiple autonomous driving systems [26], or between the original and transformed driving scenarios [34, 42, 23]. These works exploit the well-known fragility of DNNs to adversarial examples. Therefore, their main use-case concerns the identification of underrepresented scenarios in the training data (e.g., snowy weather condition) to support re-training and better generalization after re-training. The only comparable technique is the online input validation of DeepRoad [42], for which we carried out an explicit comparison in our empirical study, finding poor performance when used for online misbehaviour prediction.

Despite the different goal (test generation vs misbehaviour prediction), we share with these works the problem of how to empirically validate the proposed technique in the absence of a precise oracle that defines the expected behaviour of a self-driving car. The prevalent choice in test generators [34, 42, 23, 26] is to address the oracle problem by *differential testing*, i.e., by comparing the behaviours of multiple DNNs, or by

metamorphic testing, i.e., by comparing the behaviour before and after applying a metamorphic transformation to the input. Approaches based on verification are also being under investigation [14]. In this paper, we adopt a precise definition of *DNN misbehaviour*, which gives us a very accurate *functional oracle*, with no need for differential testing, metamorphic testing, or verification.

Search-based Generation. Abdessalem et al. [1, 4, 5] combine genetic algorithms and machine learning to test a pedestrian detection system. Mullins et al. [22] use Gaussian processes to drive the search towards yet unexplored regions of the input space, whereas Gambi et al. [13] propose ASFAULT, a search-based test generator for autonomous vehicles based on procedural content generation. ASFAULT uses search operators which mutate and recombine road segments to construct road networks for testing the lane keeping functionality of self-driving cars. Their goal is to generate extreme and challenging roads, maximizing the number of observed OBEs, while our goal is to avoid OBEs by predicting misbehaviours.

Anomaly Detection for Time Series. Anomaly detection for self-driving vehicles has been studied from a security perspective. For example, Narayanan et al. [24] use a Hidden Markov Model to detect malicious behaviours from real vehicles, and issue alerts while a vehicle is in operation. Taylor et al. [29] use LSTM neural networks to detect car’s controller area network (CAN) bus attacks, whereas Marchetti et al. [20] used an information theoretic approach. Lin et al. [19] used the Mahalanobis distance between multiple sensor data to identify unusual events in Unmanned Aerial Vehicles (UAVs). However, their approach is not image-based and thus is not comparable with our work. Moreover, ours is the first approach that detects unexpected driving conditions for the prediction of misbehaviours of an autonomous vehicle to enable self-healing.

7 Conclusions and Future Work

In this paper, we studied the problem of estimating the confidence of the DNN-based autonomous in response to unexpected execution contexts. Our tool SELFORACLE was able to anticipate by several seconds many potentially safety-critical misbehaviours, such as out of bound episodes or collisions, with a low false alarm rate, outperforming the input validator of DeepRoad. Future work concerns devising novel metrics of DNN confidence, including white-box ones, with a potential for hybridization. It would be also interesting to characterize and predict other kind of misbehaviours (e.g., derivative of steering angle) as well as implementing confidence-guided self-healing within the simulator. We believe that our promising results in online misbehaviour detection, united with the availability of a labeled dataset of crashes and a simulation environment, can foster novel approaches for online prediction and self-healing of autonomous driving systems.

References

- [1] Raja Ben Abdesslem, Annibale Panichella, Shiva Nejati, Lionel C. Briand, and Thomas Stifter. Testing autonomous cars for feature interaction failures using many-objective search. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, ASE 2018*, pages 143–154, New York, NY, USA, 2018. ACM.
- [2] David Alvarez-Melis and Tommi S. Jaakkola. Towards robust interpretability with self-explaining neural networks. In *Annual Conference on Neural Information Processing Systems (NeurIPS)*, pages 7786–7795, 2018.
- [3] Jinwon An and Sungzoon Cho. Variational autoencoder based anomaly detection using reconstruction probability. 2015.
- [4] R. Ben Abdesslem, S. Nejati, L. C. Briand, and T. Stifter. Testing advanced driver assistance systems using multi-objective search and neural networks. In *2016 31st IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 63–74, Sep. 2016.
- [5] R. Ben Abdesslem, S. Nejati, L. C. Briand, and T. Stifter. Testing vision-based control systems using learnable evolutionary algorithms. In *2018 IEEE/ACM 40th International Conference on Software Engineering (ICSE)*, pages 1016–1026, May 2018.
- [6] BGR Media, LLC. Waymo’s self-driving cars hit 10 million miles. <https://techcrunch.com/2018/10/10/waymos-self-driving-cars-hit-10-million-miles>, 2018. Online; accessed 18 August 2019.
- [7] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseon Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. End to end learning for self-driving cars. *CoRR*, abs/1604.07316, 2016.
- [8] Georg Burkhard, S. Vos, N. Munzinger, E. Enders, and D. Schramm. Requirements on driving dynamics in autonomous driving with regard to motion and comfort. In Michael Bargende, Hans-Christian Reuss, and Jochen Wiedemann, editors, *18. Internationales Stuttgarter Symposium*, pages 683–697, Wiesbaden, 2018. Springer Fachmedien Wiesbaden.
- [9] Guilherme O. Campos, Arthur Zimek, Jörg Sander, Ricardo J. Campello, Barbora Micenková, Erich Schubert, Ira Assent, and Michael E. Houle. On the evaluation of unsupervised outlier detection: Measures, datasets, and an empirical study. *Data Min. Knowl. Discov.*, 30(4):891–927, July 2016.
- [10] Vinton G. Cerf. A comprehensive self-driving car test. *Commun. ACM*, 61(2):7–7, January 2018.
- [11] Electrek. A Google self-driving car caused a crash for the first time. <https://electrek.co/2016/07/01/understanding-fatal-tesla-accident-autopilot-nhtsa-probe/>, 2016. Online; accessed 18 August 2019.
- [12] Electrek. Tesla Model S driver crashes into a van while on Autopilot. <https://electrek.co/2016/05/26/tesla-model-s-crash-autopilot-video/>, 2016. Online; accessed 18 August 2019.
- [13] Alessio Gambi, Marc Mueller, and Gordon Fraser. Automatically testing self-driving cars with search-based procedural content generation. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2019*, pages 318–328, New York, NY, USA, 2019. ACM.
- [14] Divya Gopinath, Guy Katz, Corina S. Păsăreanu, and Clark Barrett. DeepSAFE: A data-driven approach for assessing robustness of neural networks. In Shuvendu K. Lahiri and Chao Wang, editors, *Automated Technology for Verification and Analysis*, pages 3–19, Cham, 2018. Springer International Publishing.
- [15] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November 1997.
- [16] keras. Building Autoencoders in Keras. <https://blog.keras.io/building-autoencoders-in-keras.html>. Online; accessed 21 August 2019.
- [17] Keras. VGG19. <https://keras.io/applications/#vgg19/>. Online; accessed 21 August 2019.

- [18] Jinhan Kim, Robert Feldt, and Shin Yoo. Guiding deep learning system testing using surprise adequacy. In *Proceedings of the 41st International Conference on Software Engineering, ICSE '19*, pages 1039–1049, Piscataway, NJ, USA, 2019. IEEE Press.
- [19] R. Lin, E. Khalastchi, and G. A. Kaminka. Detecting anomalies in unmanned vehicles using the mahalanobis distance. In *2010 IEEE International Conference on Robotics and Automation*, pages 3038–3044, May 2010.
- [20] M. Marchetti, D. Stabili, A. Guido, and M. Colajanni. Evaluation of anomaly detection for in-vehicle networks through information-theoretic algorithms. In *2016 IEEE 2nd International Forum on Research and Technologies for Society and Industry Leveraging a better tomorrow (RTSI)*, pages 1–6, Sep. 2016.
- [21] Jonathan Masci, Ueli Meier, Dan Cireşan, and Jürgen Schmidhuber. Stacked convolutional auto-encoders for hierarchical feature extraction. In Timo Honkela, Włodzisław Duch, Mark Girolami, and Samuel Kaski, editors, *Artificial Neural Networks and Machine Learning – ICANN 2011*, pages 52–59, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.
- [22] Galen E. Mullins, Paul G. Stankiewicz, R. Chad Hawthorne, and Satyandra K. Gupta. Adaptive generation of challenging scenarios for testing and evaluation of autonomous vehicles. *Journal of Systems and Software*, 137:197 – 215, 2018.
- [23] S. Müller, D. Hospach, O. Bringmann, J. Gerlach, and W. Rosenstiel. Robustness evaluation and improvement for vision-based advanced driver assistance systems. In *2015 IEEE 18th International Conference on Intelligent Transportation Systems*, pages 2659–2664, Sep. 2015.
- [24] Sandeep Nair Narayanan, Sudip Mittal, and Anupam Joshi. Obd securealert: An anomaly detection system for vehicles. In *IEEE Workshop on Smart Service Systems (SmartSys 2016)*, May 2016.
- [25] Unity3d particle system. <https://docs.unity3d.com/ScriptReference/ParticleSystem.html>, 2019.
- [26] Kexin Pei, Yinzhi Cao, Junfeng Yang, and Suman Jana. Deepxplore: Automated whitebox testing of deep learning systems. In *Proceedings of the 26th Symposium on Operating Systems Principles, SOSP '17*, pages 1–18, New York, NY, USA, 2017. ACM.
- [27] Claude Elwood Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379–423, 7 1948.
- [28] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. VGGNet, <https://gist.github.com/baraldilorenzo/07d7802847aaad0a35d3>.
- [29] A. Taylor, S. Leblanc, and N. Japkowicz. Anomaly detection in automobile control network data with long short-term memory networks. In *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pages 130–139, Oct 2016.
- [30] Team Chauffeur. Steering angle model: Chauffeur. <https://github.com/udacity/self-driving-car/tree/master/steering-models/community-models/chauffeur>, 2016. Online; accessed 18 August 2019.
- [31] Team Epoch. Steering angle model: Epoch. <https://github.com/udacity/self-driving-car/tree/master/steering-models/community-models/cg23>, 2016. Online; accessed 18 August 2019.
- [32] The Verge. A Google self-driving car caused a crash for the first time. <https://www.theverge.com/2016/2/29/11134344/google-self-driving-car-crash-report>, 2016. Online; accessed 18 August 2019.
- [33] The Verge. Tesla hit with another lawsuit over a fatal Autopilot crash. <https://www.theverge.com/2019/8/1/20750715/tesla-autopilot-crash-lawsuit-wrongful-death>, 2019. Online; accessed 18 August 2019.
- [34] Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. Deeptest: Automated testing of deep-neural-network-driven autonomous cars. In *Proceedings of the 40th International Conference on Software Engineering, ICSE '18*, pages 303–314, New York, NY, USA, 2018. ACM.

- [35] Mis-behaviour prediction for autonomous driving systems. <https://github.com/icse2020submission/misbehavior-prediction/>, 2019.
- [36] Unity3d snow mountain track. <https://assetstore.unity.com/packages/3d/environments/roadways/mountain-race-track-53775>, 2019.
- [37] Udacity. A self-driving car simulator built with Unity. <https://github.com/udacity/self-driving-car-sim>, 2017. Online; accessed 18 August 2019.
- [38] Udacity. Udacity self-driving car’s challenge. <https://github.com/udacity/self-driving-car/>, 2017. Online; accessed 18 August 2019.
- [39] Udacity. Udacity self-driving car’s datasets. <https://github.com/udacity/self-driving-car/tree/master/datasets>, 2017. Online; accessed 18 August 2019.
- [40] Unity3d. <https://unity.com>, 2019.
- [41] V. T. Vasudevan, A. Sethy, and A. R. Ghias. Towards better confidence estimation for neural models. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7335–7339, May 2019.
- [42] Mengshi Zhang, Yuqun Zhang, Lingming Zhang, Cong Liu, and Sarfraz Khurshid. Deeproad: Gan-based metamorphic testing and input validation framework for autonomous driving systems. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, ASE 2018*, pages 132–142, New York, NY, USA, 2018. ACM.