

System Safety Monitoring of Learned Components Using Temporal Metric Forecasting

SEPEHR SHARIFI, EECS, University of Ottawa, Canada

ANDREA STOCCO, Technical University of Munich, Germany and fortiss GmbH, Germany

LIONEL C. BRIAND, University of Ottawa, Canada, and Lero SFI Centre for Software Research, University of Limerick, Ireland

In learning-enabled autonomous systems, safety monitoring of learned components is crucial to ensure their outputs do not lead to system safety violations, given the operational context of the system. However, developing a safety monitor for practical deployment in real-world applications is challenging. This is due to limited access to internal workings and training data of the learned component. Furthermore, safety monitors should predict safety violations with low latency, while consuming a reasonable computation resource amount.

To address the challenges, we propose a safety monitoring method based on probabilistic time series forecasting. Given the learned component outputs and an operational context, we empirically investigate different Deep Learning (DL)-based probabilistic forecasting to predict the objective measure capturing the satisfaction or violation of a safety requirement (*safety metric*). We empirically evaluate safety metric and violation prediction accuracy, and inference latency and resource usage of four state-of-the-art models, with varying horizons, using autonomous aviation and autonomous driving case studies. Our results suggest that probabilistic forecasting of safety metrics, given learned component outputs and scenarios, is effective for safety monitoring. Furthermore, for both case studies, the Temporal Fusion Transformer (TFT) was the most accurate model for predicting imminent safety violations, with acceptable latency and resource consumption.

CCS Concepts: • **Software and its engineering** → **Software safety**; • **Computing methodologies** → **Artificial intelligence**; • **Computer systems organization** → **Robotics**.

Additional Key Words and Phrases: ML-enabled Autonomous System, Learned Component, System Safety Monitoring, Probabilistic Time Series Forecasting

ACM Reference Format:

Sepehr Sharifi, Andrea Stocco, and Lionel C. Briand. 2024. System Safety Monitoring of Learned Components Using Temporal Metric Forecasting. 1, 1 (October 2024), 43 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Autonomous systems are increasingly being empowered using learned components to perform perception, prediction, planning, and control tasks [81]. Since such components' behaviour is learned through training, as opposed to being expressed in source code or specification, ensuring the reliability of such systems through conventional software engineering practices is inadequate. These risks are particularly acute when autonomous systems are employed in safety-critical

Authors' addresses: Sepehr Sharifi, s.sharifi@uottawa.ca, EECS, University of Ottawa, 800 King Edward, Ottawa, Canada, K1N 6N5; Andrea Stocco, andrea.stocco@tum.de, Technical University of Munich, Boltzmannstraße 3, Munich, Germany, 85748 and fortiss GmbH, Guerickestraße 25, Munich, Bayern, Germany, 80805; Lionel C. Briand, lbriand@uottawa.ca, University of Ottawa, Canada, and Lero SFI Centre for Software Research, University of Limerick, Tierney building, Limerick, Ireland, V94 NYD3.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Association for Computing Machinery.

XXXX-XXXX/2024/10-ART \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

applications, e.g., autonomous driving [14], autonomous aviation [38], medical diagnosis [98] or disease prediction [99], as failures could directly jeopardize human safety. Recently, methods have been proposed to make reliable, robust, and accurate learned components through novel testing methods [31, 95]. Nevertheless, such components are never perfect and even systems comprised of reliable components are still prone to accidents [47]. For instance, some accidents are caused by unsafe component interactions [2, 47]. Thus, the impact of ML components on safety can only be studied in the context of the system they are integrated into and in a specific operational context [12, 47].

The specialized nature of learned components, i.e., trained on necessarily limited training data, necessitates the use of runtime assurance mechanisms [80], i.e., safety monitors [77]. Runtime safety monitors observe the system, its operational context, and the inputs and outputs of a component that cannot be fully trusted, such as a learned component, predicting if its outputs may lead the system toward a safety requirement violation. In such cases, a warning is raised by the safety monitor to prevent the outputs of the learned component from propagating to the rest of the system. For example, safety recovery measures include falling back on a less efficient but trustworthy component [80] or taking pre-designed safety recovery measures, such as an emergency stop in autonomous vehicles (AVs). Safety monitors must know the operational context of the system to determine whether the component might contribute to a hazard. For instance, a misclassification by an AV object detection component can lead to non-hazardous outcomes under certain system contexts, e.g., when an AV misidentifies a horse-drawn carriage in its front as a truck and maintains a safe distance from it. Thus, runtime monitoring of both the operational context and the learned components outputs is crucial in developing effective safety monitors that can identify transitions of the system from safe to hazard states, which can lead to safety requirement violations.

However, monitoring the impact of a learned component on learning-enabled system safety poses several significant challenges. First, many safety-critical learning-enabled autonomous systems are developed by system integrators who are developing the system using various components, including learned components, many of which are developed by third parties. Thus, system integrators often do not have access to the training or test data of the learned component, nor to white-box information such as their architecture or neuron weights. Second, safety monitors should be able to monitor not only the outputs of the learned component over time but also the operational context, which typically includes static parameters such as weather, and may also include dynamic parameters such as the trajectory of other vehicles in proximity to an AV. Third, in a safety-critical context, the safety monitor must predict a safety violation early enough to allow the system or a user sufficient time to mitigate it. As such, efficiency is a key requisite, which translates to the necessity of developing monitors that exhibit a low reaction latency and do not exceed the practical limits of onboard computing units, as opposed to cloud-based alternatives.

Currently, existing methods fail to address all of the above challenges as they monitor for learned component mispredictions, as opposed to system safety violations [25, 27, 29, 83, 85, 90, 96]. Furthermore, many of the proposed methods rely on internal information sources from the learned component [28, 63, 83].

To address the above challenges, we propose a safety monitoring method based on the idea of predicting the near-future values of a safety metric, i.e., the objective measure used to determine the satisfaction or violation of a safety requirement [8], given the history of learned component outputs and the operational context of the system.

Given the safety-criticality of learning-enabled autonomous systems, where the cost of not predicting safety violations at runtime is very high, instead of relying on single forecast values for each timestep, our method predicts the probability distribution of the safety metric and relies on its tail-end values to conservatively predict safety violations. We leverage Deep Learning (DL) based

probabilistic time series forecasters and empirically evaluate state-of-the-art models in terms of prediction accuracy as well as average inference latency and runtime computation resource usage.

Contributions. The contributions of this paper are as follows:

- A safety monitoring method that leverages time series forecasting of a safety metric to identify learned component behavior and system context that lead to system safety violations at runtime.
- An application of the safety monitoring method to widely used case studies in autonomous aviation [6, 7, 40, 42, 66], i.e., Autonomous Centerline Tracking (ACT), and in autonomous driving, i.e., an Autonomous Driving System (ADS) [83], including a dataset generated from system-in-the-loop simulations.
- A large-scale empirical evaluation (7500+ GPU hours and 42 calendar days of computation) with state-of-the-art DL-based probabilistic forecasting models targeting safety metric and safety violation prediction accuracy, inference latency, and runtime resource usage.

Key Findings. The key findings of our empirical evaluation as follows:

- Overall, the results of our study suggest that probabilistic forecasting of safety metrics, given learned component outputs and scenarios, is effective for safety monitoring.
- For our ACT case study, DL-based probabilistic forecasting methods, especially those with sequence-to-sequence architecture, yield low inference latency while consuming feasible computing resources in terms of model size and peak memory usage during inference.
- Using Temporal Fusion Transformers (TFT) for predicting *imminent* safety violations—where the hazard forecast horizon is equal to the minimum reaction time, for all lookback horizons—leads to the most accurate predictions with acceptable inference latency and reasonable computational resource usage.

Paper Structure. Section 2 provides the necessary background on time series forecasting models and the main DL-based architectures. Section 3 formally defines the safety metric forecasting problem and details its challenges. Section 4 discusses related work. Section 5 presents our proposed safety monitoring method in detail. Section 6 provides an empirical evaluation of our method and discusses the results. Section 7 concludes the paper and suggests future directions for research and improvement.

2 BACKGROUND

In this section, we discuss the main characteristics of time series forecasting methods as well as the main Deep Learning (DL)-based time series forecasting architectures.

2.1 Time-series Forecasting

Time series forecasting aims at predicting the future values of a time series. As described in Januschowski et al. [37], we can distinguish among forecasting methods along a number of dimensions such as *global vs. local*, *probabilistic vs. point*, *computational complexity and costs*, and *data-driven vs. model-based*.

Global vs. Local Forecasting. Local methods involve estimating model parameters independently for each time series, while global methods estimate parameters jointly using all available time series [11, 37]. This distinction is concerned with how model parameters are estimated and does not necessarily imply a specific dependency structure between the time series. For instance, a global model can still assume independence between forecasts for different time series for computational efficiency reasons, even though it estimates parameters jointly [37]. While traditional statistical

methods often adopt local approaches, global methods have been utilized in both the statistics and machine learning (ML) communities. Recent trends show that Deep Neural Networks (DNNs) which are trained as global models, surpass all forecasting models when used as local models [11, 57].

Probabilistic vs. Point Forecasting. Forecasting techniques can also be broadly categorized into probabilistic and point forecasting methods. While point forecasts offer a single best prediction, probabilistic forecasting methods quantify predictive uncertainty, allowing decision-makers to consider this uncertainty when using the forecast [37]. For a safety-critical application, e.g., predicting whether a system will experience a safety violation in the future and taking recovery actions, it is vital to be able to take the uncertainty associated with the forecasts into account. For instance, one can use the tail-end values of the predicted probability distribution to take into account the worst-case predictions as a basis for decision making.

Methods for handling uncertainty include Bayesian approaches and frequentist approaches like model ensembles and bootstrap sampling [11]. The use of Bayesian approaches in estimating parameter and model uncertainty are well studied in ML literature (for introductory and recent work references refer to the study by Januschowski et al. [37]). The predictive uncertainty for a time series is fully described by the predictive distribution, but probabilistic forecasting methods differ in how they enable users to access this distribution, often providing pointwise predictive intervals or Monte Carlo sample paths [11]. Some methods assume a parametric form of the distribution and return its parameters [11, 75]. Modern ML methods handle uncertainty by estimating quantile functions directly [37]. The results of the M4 competition have demonstrated the accuracy of prediction intervals obtained from ML methods, even though they may lack theoretical underpinnings [37, 56]. This highlights the effectiveness of ML approaches in handling uncertainty in forecasting.

Data-driven vs. Model-based Forecasting. Methods commonly associated with machine learning, such as deep neural networks, are characterized by their data-driven nature. These approaches excel at capturing intricate patterns from data without relying on strong structural assumptions. However, their flexibility comes at the cost of requiring large amounts of data to effectively tune the multitude of parameters they possess. For instance, recurrent neural networks (RNNs) can discern complex nonlinear patterns from data, as exemplified by their ability to predict time series with oscillating variance amplitudes [37]. Nevertheless, the risk of overfitting arises due to their capacity to memorize patterns, a challenge that regularization techniques like Dropout [82], aim to mitigate. In contrast, statistical models like AutoRegressive Integrated Moving Average (ARIMA) models and Generalized Linear Models (GLMs) are characterized by their parsimonious parameterization and reliance on assumptions to model patterns [16]. These models require less data to accurately estimate their parameters but are inherently more rigid due to the limitations imposed by their structural assumptions [37]. Furthermore, a study by Kolassa [45] shows that simpler models can sometimes outperform complex, correctly specified ones, showcasing the intricacies of model-driven approaches [37]. Furthermore, model-driven approaches require meticulous feature engineering and model specification. Conversely, data-driven models are often preferred for forecasting tasks that involve a large number of time series, from which complex patterns can be extracted [57]. Moreover, DL-based forecasting models can often be trained on large datasets without the need for problem-specific feature engineering [37].

2.2 DL-based Forecasting Architectures

DL-based forecasting models can be categorized into two main categories of architectures, namely *iterative* and *sequence-to-sequence*. The iterative architecture generates forecasts step by step, where the model predicts a one-time step based on the previous hidden state and current available information [11]. The process is repeated until the desired forecast horizon is reached. Iterative

models can easily be applied to any forecast horizon length. However, since the generated forecast at each time step has an error, the recursive structure of iterative models can potentially lead to large errors being accumulated over long forecast horizons [49]. RNN models such as long short-term memory networks (LSTMs) and gated recurrent units (GRUs) are commonly employed in iterative architectures [11]. On the other hand, sequence-to-sequence architectures operate by mapping an input sequence to an output sequence, potentially of different lengths. This architecture consists of two main components: an encoder and a decoder. The encoder transforms the input sequence into a fixed-size context vector, which is then used by the decoder to generate the output sequence of a predetermined length. A typical training instance in this approach includes the target and covariate (static and time-series features or embeddings [49]) values up to a specific time point t as input, while the neural network generates a set number of target values beyond time t .

3 PROBLEM AND CHALLENGES

In this section we cast the learned component safety monitoring problem as a safety metric forecasting problem and discuss its challenges.

3.1 Problem Definition

Safety-critical systems such as autonomous vehicles (AVs) or Unmanned Aircraft Systems (UASs) use learned components such as Deep Neural Networks (DNNs) to automate and inform perception, localization, and planning tasks. In this paper, we use as a running example an Autonomous Centerline Tracking (ACT) software, which is used to ensure accurate and safe UAS taxiing on a runway, by detecting and following certain reference points or a designated path without human intervention. The distance between the system position to such reference point or centerline is called Centerline Track Error (cte). The ACT uses a DNN to estimate the cte from camera images and steer the physical system, e.g., a vehicle or an aircraft, towards the centerline where $cte = 0$.

During its operation, the system must satisfy certain safety requirements such as “*the system shall stay within 5 meters of the centerline*”. Although the learned component and the system have to be thoroughly tested and validated before going into operation, during certain challenging or unexpected execution scenarios, the learned component could contribute to the system violating the safety requirement, with potentially life-critical consequences. Therefore, early run-time prediction of a safety violation is an important endeavor and a prerequisite for developing fallback measures and mitigation strategies [80], that include blocking the output of the learned component from being broadcast throughout the rest of the system. To measure the degree of satisfaction or violation of a safety requirement, safety metrics are used. For example, from the above requirement, the safety metric can be defined as the difference between the actual cte (measured by calculating the difference between the system and centerline GPS locations) of the system and a maximum safe cte threshold of 5 m. Note that the safety metric value varies over time and is therefore calculated at each time step.

More concretely, let s be the ACT system including the learned component m for image-based cte estimation, operating in its environment under an operational scenario x . The latter is represented by static and dynamic properties that exist during system operation, e.g., the angle of the sun, cloud cover, runway properties, or the initial position of the aircraft. For each time step t , the system takes an input $in_{s,t}$ from the camera, thus capturing the state of the environment, and provides a pre-processed (e.g., by drivers or information fusion) image $in_{m,t}$ ready to be consumed by m . m produces a real number $out_{m,t}$ which represents the cte estimate. s processes $out_{m,t}$, generates a steering command $out_{s,t}$, and applies it to the system. The state of the environment relative to s changes based on $out_{s,t}$, and the entire process repeats during the operation of s , whereas the next learned component inputs are partially determined by previously learned component outputs.

For a safety requirement r (e.g., “the system shall not deviate from the centerline more than 5 m”), we can measure the degree of safety violation of s at time t , denoted by $y_{r,t}$, with a continuous function $f_r(t)$ which determines at time t whether r has been violated ($y_{r,t} = f_r(t) \geq 0$) or how close it has come to violating it ($y_{r,t} = f_r(t) < 0$). Note, that the exact definition of f_r is context-dependent and varies based on the system and the safety requirement of interest. For the ACT system and the safety requirement above, we define f_r as denoted in Equation 1.

$$y_{r,t} = f_r(t) := |cte_{act}| - |cte_{thr}| \quad (1)$$

Whereas, cte_{act} and cte_{thr} are the actual and safety violation threshold values of the centerline track error, respectively. Based on the above context, let $x^{(n)}$ be a given set of environmental conditions in the space of all possible conditions, also referred to as *operational scenarios*. Let $o_{m,t-k:t}^{(n)}$ and $y_{t-k:t}^{(n)}$ be the sequence of observed m 's outputs and safety metric values from time $t-k$ to t (where k denotes the *lookback horizon*), given $x^{(n)}$.

Given a *hazard forecast horizon* h^1 , we want to predict the sequence of safety metric values from time $t+1$ to $t+h$, i.e., $\hat{y}_{t+1:t+h}^{(n)}$, using a prediction model g , as expressed in Equation 2, as accurately as possible.

$$\hat{y}_{t+1:t+h}^{(n)} = g(h, y_{t-k:t}^{(n)}, o_{m,t-k:t}^{(n)}, x^{(n)}) \quad (2)$$

As mentioned in Section 1, aside from the ACT system mentioned above, this paper additionally targets a second cases study, i.e., an Autonomous Driving System (ADS) which performs the *lane keeping* functionality autonomously, while relying only on image inputs from the camera. The formal problem definition of the problem provided in this section, especially Equation 2, equally apply to the ADS case study. We provide, in Section 6.1, complete details for both the ACT and ADS case studies evaluated in this paper.

3.2 Challenges

Given the context and the problem definition provided in Section 3.1, we observe multiple challenges. First, the development of learned components is often outsourced to third parties [30, 71], which are later integrated into the main autonomous system. Thus, the limited or lack of access to the learned component details inhibits the application of white-box methods for safety monitoring [85]. Such details include the training data and the model's architecture, weights, activation patterns, or gradients during the feed-forward pass.

Second, as mentioned in Section 1, evaluating the safety of a learning-enabled autonomous system relies both on the *static* operational context data (*scenario*) and *dynamic* time-series data related to the behavior of the learned component and the history of safety metric values over time. Thus, the safety monitor should be able to utilize both types of data to provide an accurate forecast of the safety metric values over the hazard forecast horizon.

Third, safety monitors are often developed for safety-critical cyber-physical systems with limited computation capabilities. Thus, it is paramount that the safety monitor introduces low latency and memory overhead to the system. Although many safety monitoring methods that rely on white-box confidence estimation techniques [25, 83] are more accurate than their black-box counterparts, their memory and computing overhead makes their adoption in resource-constrained settings impractical [84].

¹Hazard forecast horizon is the number of timesteps in the future [26], over which we want to predict the values of a safety metric such that safety violations (*hazards*) can be predicted and mitigated or avoided.

To address the above challenges, henceforth denoted C1 through C3, respectively, in this paper we evaluate time series forecasting methods, especially the ones based on DL, to forecast the safety metric values of a learning-enabled system over the hazard forecast horizon. DL methods have been shown to provide accurate forecasts while being amenable to multiple data types (static and time series) and a large number of samples [11, 37, 57]. In Section 5, we provide further details on the proposed safety metric forecasting solutions, whereas in Section 6 we discuss the experimental evaluation of different state-of-the-art forecasting models for the ACT case study.

4 RELATED WORK

This section discusses existing studies related to the problem of learned component safety monitoring. Some surveys [13, 53, 60, 70], distinguish between Out-Of-Distribution (OOD) detection and uncertainty estimation (quantification) methods. The former focuses on identifying learned component inputs that are not within its training distribution, while the latter estimates the uncertainty associated with the learned component outputs. Since these methods aim, at a high level, at a similar goal, i.e., identifying inputs that lead to uncertain and thus untrustworthy outputs, they should therefore be discussed here. Next, we discuss the main safety monitoring techniques in the literature, primarily categorized based on the type of system information access they assume, namely black-box and white-box approaches.

4.1 Black-box Methods

Black-box methods use information such as learned component inputs and outputs, as well as its training and test datasets,² to identify the shift in the distribution of inputs observed during operation from the training input distribution, which can lead to mispredictions during operation [101].

For example, Zhang et al. [97] proposed DeepRoad which was mainly designed for testing AV learned components by validating single input images according to their minimum distance from the training set based on the embeddings generated according to VGGNet [78] features. SelfOracle [85] is a black-box failure predictor that uses an autoencoder and time series-based anomaly detection to reconstruct the input images observed by the learned component and to use reconstruction loss to detect OOD inputs. Similar methods that utilize variational autoencoders (VAEs) to measure an anomaly score have also been proposed by other studies [15, 27, 33]. DeepGuard, proposed by Hussain et al. [33], uses the VAE reconstruction error to prevent roadside collisions with other vehicles, Borg et al. [15] proposed an OOD detector based on VAEs combined with object detection for an automated emergency braking system.

Moreover, some black-box methods quantify the uncertainty of the learned component outputs, to help practitioners identify the learned component inputs leading to unreliable outputs, by estimating probability distributions of the outputs given past system executions. These methods leverage Bayesian networks or their approximations [7, 24, 55], allowing them to incorporate expert domain knowledge in their Bayesian network models. In the context of autonomous aviation systems (similar to our ACT example), Asaadi et al. [7] used a non-parametric Bayesian-based uncertainty quantifier, i.e., Gaussian Process (GP) regressor, trained on a subset of the learned component training data, to estimate the uncertainty in learned component outputs given its inputs.

²The survey conducted by Riccio et al. [73] categorizes the methods that require access to learned component train and test dataset as *data-box* methods. However, to avoid confusion, we categorize them as black-box methods here as they do not use any internal information from the model itself.

4.2 White-box Methods

Unlike black-box methods, white-box methods take advantage of internal information sources from the learned component, e.g., model confidence [28], neuron activation patterns [94] or gradients [83], comparing their observations at runtime (during the learned component operation) against design-time (during the learned component training).

For example, Lakshminarayanan et al. [46] proposed the use of an ensemble of neural networks (*Deep Ensembles*), to effectively predict the uncertainty of perception component outputs at runtime. Kendall and Gal [43] proposed a Bayesian deep learning framework that captures uncertainties associated with both the learned component inputs, also referred to as *aleatoric* uncertainty, as well as the model itself, also known as *epistemic* uncertainty, for a perception component (image segmentation and depth regression). In the context of autonomous driving, Grewal et al. [25] evaluate different uncertainty quantification methods for the misbehavior prediction of failures. Hendrycks and Gimpel [28] used the learned component's own confidence, i.e., its softmax probability distributions, to measure uncertainty in learned component outputs. However, since learned components are prone to generating incorrect outputs (misprediction) with high confidence [63], many methods have leveraged other information sources to estimate uncertainty. ThirdEye [83] uses an eXplainable AI (XAI) technique, namely attention maps, to generate a confidence score for the learned component (in this case a DNN) based on input images and gradients of the DNN. The generated confidence score is then used to predict a failure by comparing it with a failure threshold learned from past system executions (simulations).

4.3 Limitations of Existing Methods

Although the white-box and black-box methods mentioned above are effective at evaluating the inputs to the learned component, they do not consider the effect of learned component outputs on system safety. Learned component inputs that can lead to inaccurate outputs (i.e., mispredictions) may not lead to system safety violations, depending on the system's operational context. As discussed in Section 3.2, a safety monitoring method must be able to predict the combinations of system context and learned component outputs that can lead to system-level safety violations.

In terms of information requirements, methods such as DeepRoad [97] and the Bayesian method proposed by Asaadi et al. [7], require access to training and test datasets. Furthermore, as mentioned in Section 4.2, to identify safety-violating inputs, white-box methods rely on internal information of the model [28, 83, 94]. As discussed in Section 3.2, system integrators often do not have access to such information, nor training and test datasets, as they are frequently developed by third parties.

Finally, both the black-box and white-box methods discussed in Sections 4.1 and 4.2, respectively, were not evaluated in terms of their inference latency and computation resource usage at runtime [7, 24, 55, 83, 85, 94, 97].

Different from the described black-box and white approaches, we evaluate time-series DL methods for safety monitoring, a previously unexplored topic. We empirically evaluate the effectiveness of such methods when using both the operational context of the system and learned component behavior while being computationally feasible for runtime monitoring. We then predict when system context and learned component behavior together lead to system-level safety violations.

5 TEMPORAL FORECASTING OF SAFETY METRICS

In this work, as mentioned in Section 3, we have cast the safety metric prediction problem as a safety metric forecasting problem. Recall that, as described in Section 3.1 and Equation 2, the inputs to the forecasting model are static operational scenario data (x), dynamic (time-dependent) learned component behavior ($o_{m,t-k:t}$), and past safety metric data of learning-enabled autonomous

systems ($y_{t-k:t}$). While the output of the forecasting model ³ is the safety metric forecasts over the hazard forecast horizon ($\hat{y}_{t+1:t+h}$). In the rest of this section, we evaluate existing time-series forecasting methods and propose potentially suitable candidates for the problem of predicting the value of the safety metric of a learning-enabled autonomous system over a hazard forecast horizon, as introduced in Section 3.1. Note that the hazard forecast horizon is set by the system developers and safety engineers according to the system properties, its intended mission and its corresponding set of operational contexts, also referred to as its Operational Design Domain (ODD) [19, 77].

We assume that for training the forecasting model, a dataset containing numerous samples from historical system executions under various scenarios has been collected through simulation testing, a common practice for learning-enabled autonomous systems deployed in safety-critical contexts [4, 52]. As mentioned in Section 2, model-based time-series forecasting models such as fitting ARIMA models are not suitable for forecasting when the dataset is large, highly dimensional, or contains non-linear relationships (between features and target), as the time required to fit them to data considerably increases and their prediction performance degrades. Classical machine learning (ML) models, especially tree-based methods (e.g., gradient-boosted tree methods), have been used widely in time series forecasting as they provide superior prediction performance to their statistical counterparts [11]. While they can be trained on a large number of samples, these models require substantial feature engineering, thus requiring expert knowledge of the system [11], which consumes significant time and effort.

Recently, deep learning time-series forecasting models (DL forecasters) have shown great potential for challenging forecasting problems. As discussed in Section 2, DL forecasters can be trained as global models on large time-series samples without requiring white-box knowledge of the system under test, or specific feature engineering. Thus, addressing challenge C1 described in Section 3.2 by relying only on black-box information related to the system under test. Moreover, as survery by Benidis et al. [11], some DL forecasters can handle samples that contain both static and dynamic data types with complex and non-linear relationships [11], addressing challenge C2 stated in Section 3.2. Last, DL forecasters, given real-valued times series inputs, consume limited computation resources and enable low inference latency, addressing challenge C3 discussed in Section 3.2 ⁴.

Finally, due to the safety-critical nature of learning-enabled systems, it is important to account for the uncertainty associated with predicted safety metric values. As discussed in Section 2, DL-based probabilistic forecasting methods account for such uncertainty by predicting the values of the safety metric probability distribution. Knowing the values at the tail-end of the predicted probability distribution of the safety metric allows us to rely on such values for worst-case metric predictions.

Thus, to address the challenges outlined in Section 3, we propose training a DL-based probabilistic forecaster, given historical execution data of the system and its safety metric, such that it provides forecasts of the safety metric value over the hazard forecast horizon. Note that the weights of the DL model are selected and remain the same for all the different scenarios and time-series data that is used for training, i.e., it is a *global* model. Concretely, the DL-based probabilistic forecaster takes the times series of the safety metric values and learned component outputs, as well as scenario parameters as input, and returns time series predictions of the safety metric probability distribution. Note that the duration of the input time series is equal to the lookback horizon, while the duration of the predicted time series is equal to the forecast horizon.

³Also referred to as *target variable* in time series forecasting literature [11].

⁴We will empirically measure the inference latency as well as the computation resource usage of state-of-the-art DL forecasters in Section 6.

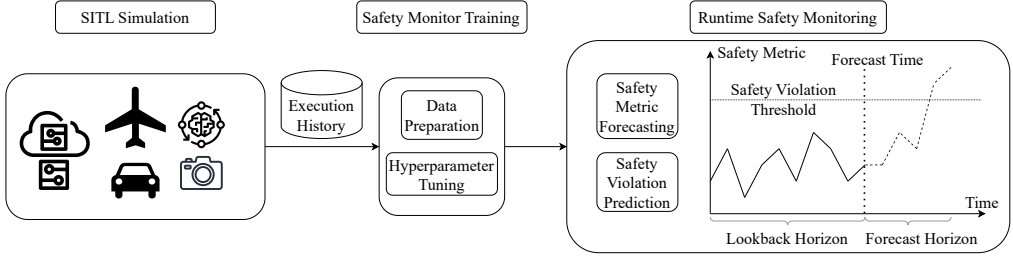


Fig. 1. The overall process for training the temporal safety metric forecaster for safety monitoring.

Safety Violation Prediction. As discussed earlier in Section 3.2, one of the main goals, aside from knowing the value of the safety metric at each timestep, which is crucial for safety-critical decision-making, is to predict safety violations that might occur in the near future, i.e., over the hazard forecast horizon. Therefore, we further describe how a safety metric forecaster can be used to predict safety violations.

Recall that a safety metric forecaster, at timestep t , provides predictions of the safety metric value from timestep $t + 1$ to $t + h$. Further recall that we have assumed that the function measuring the safety metric is defined such that non-negative values imply safety violation, similar to Equation 1. Thus, if the *maximum* of the predicted safety metric values over the hazard forecast horizon is non-negative, we can say that a safety violation has been predicted.

More specifically, given predicted safety metric values and a hazard forecast horizon h , we can define the safety violation function $v(h)$ as detailed in Equation 3.

$$v(h) = \text{sign}(\max\{y_{t+1:t+h}\}) \quad (3)$$

Note that the sign function $\text{sign}(i)$ used in Equation 3 returns $+1$ when $i \geq 0$ and returns -1 otherwise [10, 20]. Based on the definition provided in Equation 3, a safety violation is detected over the hazard forecast horizon, i.e., from timestep $t + 1$ to $t + h$, if $v(h) = 1$.

Figure 1 depicts the overall process for training and deploying the safety monitor. The process starts with data generation using System-in-the-Loop (SITL) simulation where the required data to train the safety monitor, as discussed above, is generated. Then, in the training stage, we preprocess the execution history, tune the hyperparameters of the safety metric prediction model, and train the best model on the complete dataset. Finally, the trained model is deployed during system operation where the future values of the safety metric are predicted, which is in turn used for safety violation prediction.

As surveyed by Benidis et al. [11], various DL models with different architectures have been proposed and applied to time-series forecasting. The number and variety of the proposed models make the problem of selecting the appropriate DL model for safety metric and violation prediction an important challenge, which can only be addressed through empirical investigation. To this end, we have empirically evaluated state-of-the-art DL-based time-series forecasting models in our specific application context.

6 EMPIRICAL EVALUATION

In this section, we report the empirical evaluation of time series-based safety monitors applied to an ACT and an ADS system. We aim to answer the following research questions:

RQ₁ (Safety Metric Prediction Accuracy) How do different forecasting models score and compare in terms of safety metric prediction accuracy?

RQ₂ (Safety Violation Prediction Accuracy) How do different forecasting models perform and compare in terms of safety violation prediction accuracy?

RQ₃ (Accuracy Sensitivity Analysis) What is the impact of varying lookback and hazard horizon window sizes on safety metrics and safety violation prediction accuracy?

RQ₄ (Resource Overhead Sensitivity Analysis) How do different forecasting models compare in terms of the memory and time overhead of making predictions?

RQ₁ and RQ₂ are motivated by the wide variety of potentially applicable forecasting models, which raises the need for experimental evaluation to determine which one scores best in terms of safety metric forecast and safety violation prediction accuracy, respectively. RQ₂ is particularly relevant in scenarios where a safety monitor does not have a particularly high accuracy in predicting safety metric values, yet its predictions sufficiently contribute to accurate safety violation predictions.

Note that hazard forecast horizon and lookback window sizes are design choices for system developers. Increasing the hazard forecast horizon is expected to decrease safety metric prediction accuracy, as indicated by previous studies [17]. Conversely, increasing the lookback window size is expected to enhance accuracy. However, we also expect such changes to impact models' runtime performance, as they impact the number of model parameters, influencing factors like latency and memory overhead. Given that the forecasters are destined for deployment in resource-constrained safety-critical systems, understanding the consequences of altering window configurations—specified by hazard forecast horizon window size and lookback to forecast window size ratio parameters—on prediction accuracy and runtime performance is crucial for system developers in practice. Consequently, our evaluation further explores the effects of different window configurations on prediction accuracy (RQ₃) and runtime performance (RQ₄).

6.1 Evaluation Subjects

We evaluate the DL-based forecasting models by applying them to two case studies related to an autonomous centerline tracking system for autonomous taxiing (ACT), and an autonomous driving systems (ADS) focused on lane keeping. In this section, we provide for each case study, an overview of the subject system, explain the details of the evaluation dataset and the simulation workflow used to generate it.

6.1.1 ACT Case Study.

Subject System and Simulation Platform. We used an open-source ACT system [41], similar to previous studies that had ACT-related case studies [8, 18, 65]. Note that the ACT system is crucial for safe taxiing operation of autonomous aviation systems. As reported by the U.S. Department of Transportation, National Transportation Safety Board (NTSB) [61], as well as major commercial aircraft manufacturers [1, 87], fatalities, loss of aircraft, and other substantial damages have occurred during the taxi stage of flight.

As illustrated in Figure 2a, the ACT system consists of a camera, a learned component (i.e., a DNN estimator that outputs cross-track error *cte* and heading error *he* estimates given an image input), and a proportional controller that generates control commands steering the aircraft.

Previous studies that used ACT as a case study used the TaxiNet model [8, 18, 65], a DNN developed by Boeing for ACT applications, as their learned component. Since we were not granted access to TaxiNet, we relied on the open source version called TinyTaxiNet [41],⁵ which has a lower number of deep layers and input vector size. We used X-Plane 11 [93], a high-fidelity flight simulator—which is used for training pilots [93] in all flight phases including taxiing—to control various aircraft and environmental parameters. Based on the simulator's controllable elements,

⁵To the best of our knowledge, this is the only open-source DNN model trained for ACT.

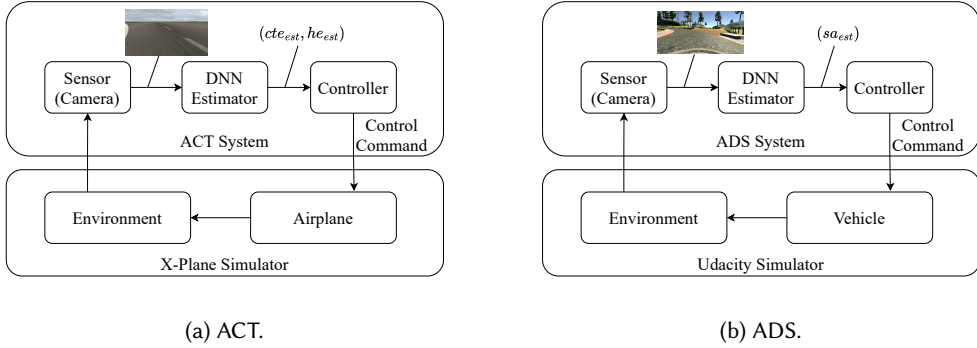


Fig. 2. System-in-the-loop simulation of the ACT system and the X-Plane simulator (a), and the ADS system and the Udacity simulator (b).

in line with previous studies [41], we considered the following four scenario elements: period of the day, cloud cover, starting cross-track error, and starting heading error. A detailed explanation of the scenario elements and their value ranges, and the TinyTaxiNet model can be found in the supporting material (see Section 6.9).

Considering the capability of X-Plane in controlling environmental parameters and the major functionality of our target learned component (i.e., ACT), we focus on the following two safety requirements: (1) “The aircraft should have a distance ($|cte_{act}|$) no more than cte_{thr} from the centerline, while taxiing on the runway.” (2) “The aircraft should have a heading angle ($|he_{act}|$) no more than he_{thr} from the centerline, while taxiing on the runway.” To compute the safety metric related to the above requirements, we measure the actual distance of the center of the aircraft from the centerline ($|cte_{act}|$), as well as the actual angle between the longitudinal axis of the aircraft and the centerline ($|he_{act}|$), at every time step. Given our definition of the safety metric provided in Section 3.1 (Equation 1), negative values of the computed safety metric ($|cte_{act}| < cte_{thr}$) imply that the cte safety requirement is *not* violated, whereas zero and positive values ($|cte_{act}| \geq cte_{thr}$) indicate a safety violation. Note that the same applies to the he safety requirement. Based on the size of the aircraft and its taxiing speed, we set cte_{thr} to 5 m and he_{thr} to 5° , in line with other ACT-related studies [5, 7, 65].

Evaluation Dataset. Given the input and output of the safety metric forecaster, discussed in Section 5, the dataset used to train and evaluate the DL-based forecasters must contain time series values of the learned component’s output and the safety metric, given a specific scenario. Thus, we generated a dataset based on the autonomous taxiing aircraft case study (Section 6.1.1, inspired by the simulation setup proposed by [41]).

Concretely, we used the Latin Hypercube Sampling (LHS) method to generate 1,996 unique scenarios to ensure coverage of the scenario parameter search space [54]. LHS is a sampling method that is used to generate near-random samples from a multidimensional distribution, while ensuring uniform coverage of the simulation input space (i.e., *scenario space*) by stratifying each input dimension [54, 76]. Given that the high-fidelity simulations that we require to generate the dataset are computationally expensive, LHS allows us to obtain diverse scenarios with a limited number (1996) of simulations. We executed each scenario using the ACT simulation stack until the aircraft reaches to its destination on the runway (taking an average duration of 200 s), whereby time series data of TinyTaxiNet outputs (cte_{est} , he_{est}) are recorded, as well as the safety metrics (computed

based on Equation 1, with $cte_{thr} = 5$ m and $he_{thr} = 5^\circ$, as discussed in Section 6.1.1). Thanks to our sampling strategy, 52% and 21.9% of the scenarios recorded in the dataset for the cte and he safety requirements, respectively, include safety violations in their test set (the division of the dataset is discussed in Section 6.3.1). We therefore generate a large number of diverse cte and he safety violations in the test sets. In many simulations, the ACT system misidentifies an imaginary line, which is a parallel offset of the runway centerline by more than 5 m, as the actual centerline and continues taxiing along it. This can explain why the dataset includes more cte safety violations than he safety violations.

Finally, following best practice, similar to the guidelines provided by previous studies on training and evaluating DL-based time series forecasting models [48, 75, 91], we normalized the time series data, i.e., cte_{est} , he_{est} and the safety metrics values using Z-score normalization, thereby reducing model bias caused by differences in time series magnitudes among various parameters and scenarios [50]. Further details regarding our generated dataset can be found in the supporting material (see Section 6.9).

6.1.2 ADS Case Study. For the ADS case study, we relied on the artifacts available in the paper by Stocco et al. [83], where the authors tested a lane-keeping ADS in a driving simulator [83]. The dataset includes not only the time series of the learned component outputs and the scenario parameters but also the time series of the safety metric, which is crucial for our safety monitoring method. Over the rest of this section, we provide an overview of the ADS subject system and the simulation platform used to generate the dataset. We then discuss the details of the raw dataset, as well as our preprocessing steps leading to the final dataset used by our study.

Subject System and Simulation Platform. The ADS case study involves a lane keeping ADS widely used in previous studies [32, 36, 73, 83, 85, 85]. Note that the lane keeping ADS is critical for safe operation of AVs. As described in the original paper, in the U.S., run-off-road crashes are one the most important types of road accidents in terms of frequency and cost.

As depicted in Figure 2b, the ADS consists of a camera, a learned component (i.e., a DNN which estimates the required *steering angle* sa_{est} to keep the vehicle within the lane given an image input), and a controller that issues the control commands to the vehicle. The learning component is based on the NVIDIA Dave-2 model architecture [14], a DNN-based steering angle estimator that is trained with a set of images collected while a human driver is driving a vehicle. The simulator used to evaluate the ADS case study was the Udacity simulator for self-driving cars [88], a driving simulator which has been widely used in the ADS testing literature [32, 36, 36, 83–85]. Udacity provides close-loop tracks to simulate an ADS driving under various scenario conditions. Based on the controllable elements of the Udacity simulator, the original paper considered the following two scenario elements: weather conditions (i.e., clear, fog, rain, and snow), and weather intensity.⁶ Since the authors record the time series of cte_{act} , we use the recorded value at the beginning of each episode as our third scenario element, i.e., the starting cross-track error, which indicates the initial position of the vehicle. A detailed explanation of the value ranges for the scenario parameters can be found in the supporting material (Section 6.9).

Considering the main functionality of the target learned component, i.e., lane keeping, the following safety requirement is considered: “*The vehicle should have a distance ($|cte_{act}|$) no more than cte_{thr} from the centerline, while driving on track*”. To calculate the safety metric, it is necessary to obtain measurements of cte_{act} , which are provided in the original dataset. Similar to the cte safety requirement for the ACT case study, the safety requirement is violated when $|cte_{act}| \geq cte_{thr}$

⁶Note that we are only mentioning the parts of the study conducted by Stocco et al. [83] that are relevant to our study. We refer the reader to the study itself for comprehensive details on all the evaluations conducted by the authors.

and not violated otherwise. Given that the total width of the track set in the simulator, the size and the speed of the vehicle, we set cte_{thr} to 5 m, which is reasonable as it provides the vehicle less than 0.5 m of side clearance from the edge of the track.

Evaluation Dataset. The dataset generated by Stocco et al. [83] contains a time series of the Dave-2 outputs (estimated steering angles), a time series of cte_{act} (which we used to compute the time series of the safety metric, as described above), and the scenario parameters.

Concretely, the authors executed the ADS in the Udacity simulator, i.e., let it drive for one lap around the track under various weather conditions with intensity increments of 10%. Therefore, to cover a diverse range of the scenario space, the authors recorded 31 one-lap simulations ($1 \times \text{clear} + 10 \times \text{fog} + 10 \times \text{rain} + 10 \times \text{snow} = 31$) which include the time series of cte_{act} measurements [83]. Note that in some scenarios, the ADS drives the car out of the track, in which case the car is reset on the next waypoint on the track. This reset during the execution of the scenario leads to a discontinuity during the execution of the learned component output and safety metric time series. Therefore, we divide the executions at the points where the vehicle has gone out of the track completely into separate episodes, to handle the discontinuity in the time series data. However, this has led to having diverse execution lengths, e.g., from 5 s to more than 100 s. To make the size of the episodes more uniform, we discarded the very short episodes, i.e., less than 15 s, as they do not contain the minimum number of timesteps required to train and test the DL-forecasting models, and divided the larger episodes into shorter chunks. The resulting dataset contains 175 episodes with an average duration of 17.5 s. Note that the size of the ADS dataset is a fraction of (approximately 0.8%) the size of the dataset we generated for the ACT case study. As we will see, this will have an impact on our results and conclusions. Despite the ADS dataset size, we observe that 33.7% of the episodes include safety violations in their test set. Therefore, the ADS dataset includes a considerable number of diverse safety violations thanks to the sampling strategy used to search the scenario space, as discussed above.

Finally, similar to the ACT case study (Section 6.1.1), we normalized the time series data, i.e., estimated steering angle and the safety metric values using Z-score normalization. We have included more details about the raw and the preprocessed datasets in our supporting material (Section 6.9).

6.2 Models Under Evaluation

In this section, we outline the chosen forecasting models for evaluation and discuss the hyperparameter tuning process applied to optimize the selected models.

As discussed in Section 5, we are interested in global univariate probabilistic forecasting models that take in input both dynamic time series and static scenario data and provide probabilistic forecasts of the safety metric. We selected the models for evaluation from the GluonTS library [3], a widely used probabilistic DL-based forecasting Python library, containing the implementations of many state-of-the-art models. From the list of available models in the library, we selected four models, three of which satisfy all the requirements of the safety metric forecasting problem (i.e., a global univariate probabilistic forecasting model capable of processing both static and dynamic inputs), whereas the fourth model acts as a competitive baseline, even though it does not fully satisfy all requirements.

- MQCNN: a sequence-to-sequence model which is the CNN-based variant of Multi Quantile Recurrent Forecaster, using a CNN encoder instead of an RNN [91].
- Temporal Fusion Transformer (TFT): a sequence-to-sequence transformer-based model [48].
- Seq2Seq: a vanilla sequence-to-sequence model with a CNN encoder and an MLP decoder [3].
- DeepAR: an iterative model which utilizes both RNNs and autoregressive techniques to iteratively capture temporal dependencies [75]. Due to DeepAR's architecture, it only takes

as input the static scenario parameters and time series of the target variable, i.e., safety metric. Despite DeepAR architecture's lack of ability to process the time series of the learned component output, we have included it in our evaluation as a competitive baseline, due to its high performance in the forecasting benchmarks [57] and wide use in industry [11].

Hyperparameter Tuning. We fine-tuned the hyperparameters of each considered model before answering the research questions, considering the values (fixed or range) retrieved from the original publications. We have tuned the hyperparameters for each safety requirement separately, as they lead to different datasets for the models to be trained and tested on. The relevant hyperparameters and their value ranges for each model are as follows:

- *MQCNN* [91]. *Number of layers in the MLP decoder (or dim)* was selected in the range {2, 4, 8}. *Number of neurons in the hidden layer* was selected in the range {20, 40, 80}. *Number of channels per layer of the CNN encoder* was chosen in the range {20, 40, 80}.
- *TFT* [48]. We set the *dropout rate* to values ranging from 0.1 to 0.3 in steps of 0.1. We took the values of *number of attention heads* and *state size* in the ranges {1, 4} and {80, 160, 320}, respectively. We kept the loss function the same as the original paper, i.e., *quantile (pinball) loss* [48].
- *Seq2Seq*. *Number of layers* was chosen from {2, 4, 8}. *Number of neurons per layer* was selected from {10, 20, 40}.
- *DeepAR* [75]. The *number of RNN layers* was set to 3 as in the original study [75]. The *number of RNN nodes per layer* was selected in the range {40, 100}. We selected the type of RNN nodes in each layer as being one among {LSTM, GRU}. We selected the *dropout rate* in the range {0.1, 0.2, 0.3}. The loss function *negative log likelihood* was used, in line with the original study [75].
- *DL Training*. For all the deep learning models above, we selected the hyper-parameters related to training as follows. We chose *learning rate* in the range {0.0001, 0.001, 0.01}. We selected *max gradient norm* in the range {0.01, 1.0, 100.0}. We evaluated *batch size* in the range {64, 128, 256}. The Adam optimizer [44] was used for training all the models, as per the original paper implementations or that of the GluonTS library.

For other hyperparameters, we relied on the suggested values used in the original studies or the default value in their implementation (additional details on the parameter settings are available in the supporting material in Section 6.9).

We trained each model configuration (defined by a combination of hyperparameters) on the training set (i.e., 70% of the dataset) and evaluated it on the validation dataset (i.e., 10% of the dataset), 5 times, to account for randomness, for example, due to random seeds. Similar to RQ_1 , we compared the models based on their q-Risk (Equation 4) values at quantiles considered in RQ_1 . The details of the evaluation metric, i.e., q-Risk and the quantiles under consideration, are presented in Section 6.3.1. Table 1 summarizes the hyperparameters for each model, their possible values, and the selected hyperparameters, for *cte* (ACT_{cte}) and *he* (ACT_{he}) safety requirements of the ACT case study and the *cte* (ADS_{cte}) safety requirement of the ADS case study, respectively.

Evaluation Hardware. To train each configuration of the models under investigation on the evaluation dataset and evaluate them, we used the following compute resources: 1x NVIDIA V100 GPU with 32GB HBM2 memory, 16 cores of Intel Silver 4216 Cascade Lake 2.1GHz CPU, and 128GB of RAM.

Table 1. Hyperparameters of the models under evaluation, for *cte* and *he* safety requirements of the ACT case study, and the *cte* requirement of the ADS case study.

Hyperparameter	Value Range	ACT _{cte}	ACT _{he}	ADS _{cte}
Seq2Seq				
Batch Size	64, 128, 256	128	64	64
Learning Rate	1e-4, 1e-3, 1e-2	1e-3	1e-4	1e-4
Gradient Clipping Value	1e-2, 1.0, 1e+2	1.0	1e-2	1e+2
Number of MLP Decoder Layers	1, 2, 4	2	2	2
Number of Neurons per MLP Layer	20, 80	80	20	80
DeepAR				
Batch Size	64, 128, 256	64	64	64
Learning Rate	1e-4, 1e-3, 1e-2	1e-2	1e-3	1e-3
Gradient Clipping Value	1e-2, 1.0, 1e+2	1e-2	1.0	1.0
RNN Node Type	LSTM, GRU	GRU	GRU	LSTM
Number of RNN Nodes	40, 100	40	40	100
Dropout Rate	0.1, 0.2, 0.3	0.1	0.1	0.1
TFT				
Batch Size	64, 128, 256	256	256	128
Learning Rate	1e-4, 1e-3, 1e-2	1e-3	1e-3	1e-2
Gradient Clipping Value	1e-2, 1.0, 1e+2	1e-2	1.0	1e+2
State Size	40, 80, 160	160	160	160
Number of Attention Heads	1, 4	4	4	4
Dropout Rate	0.1, 0.2, 0.3	0.1	0.1	0.1
MQCNN				
Batch Size	64, 128, 256	256	64	128
Learning Rate	1e-4, 1e-3, 1e-2	1e-3	1e-4	1e-4
Gradient Clipping Value	1e-2, 1.0, 1e+2	1.0	1e-2	1.0
Number of MLP Decoder Layers	1, 2, 4	2	2	2
Number of Neurons per MLP Layer	20, 80	20	20	80
Number of Channels	20, 40	20	20	20

6.3 RQ₁: Safety Metric Forecast Accuracy

In this section, first we provide the details of our evaluation methodology to answer RQ₁ (Section 6.3.1). Then, we present the results for the autonomous taxiing (ACT) case study (Section 6.3.2), followed by the results of the autonomous driving (ADS) case study (Section 6.3.3). Finally, we draw conclusions from both case studies (Section 6.3.4) and present our answer to RQ₁.

6.3.1 Methodology. To answer RQ₁, we divide each dataset (detailed in Section 6.1.1 and Section 6.1.2) into training, validation, and test datasets, which correspond to 70%, 10% and 20% of the dataset, respectively. The training set is used to train the time series forecasting models, whereas the validation dataset is used for hyper-parameter tuning (Section 6.2). Finally, we generated predictions using the trained models on the test dataset, which is disjoint from the training and validation datasets. To avoid *look-ahead bias* [35], we used time-based splitting [11, 49, 67], such that all the samples in the test dataset occur after the validation dataset, whose samples occur after the training dataset. Our evaluation method, also referred to as *rolling-horizon out-of-sample testing* [86], provides an evaluation of the model forecasting accuracy on multiple rolling-window samples in each time series that are not seen by the model during training, and aggregates them

over all time series in the dataset (see Equation 4)⁷. Note that, the evaluation of the model is conducted over multiple subsequent samples in the test set of each time series, thus providing an accuracy measure of the forecasting model over time.

The most accurate model resulting from the hyper-parameter tuning phase, in terms of the loss function, was selected and retrained on the union of training and validation datasets. The hyper-parameters for each optimized model selected for evaluation, are listed in Table 1. The retrained model was evaluated against the test dataset and its corresponding evaluation metric was computed. To account for randomness (in the training process), we repeated the above process, i.e., training the best model on the joint training and validation set and evaluating it on the test set, 30 times and reported descriptive statistics of the evaluation metric. To evaluate the statistical significance of the difference in accuracy metrics of different DL-based safety metric forecasters, we used the Mann-Whitney U test [58]. To measure the effect size of the differences, we measured Vargha and Delaney's \hat{A}_{AB} , where $0 \leq \hat{A}_{AB} \leq 1$ [89]. Generally, the value of \hat{A}_{AB} indicates a small, medium, and large difference (effect size) between populations A and B when it is higher than 0.56, 0.64, and 0.71, respectively.

We investigated RQ_1 while considering, as the window size for the hazard forecast horizon, the minimum reaction time required for a human to take over control of the system in case of a hazard. Considering that each time step in the ACT dataset corresponds to one second and the minimum reaction time for a human with vehicles traveling at 30 mi/h is 3 s [83], we set the minimum hazard forecast horizon to 3 timesteps. Furthermore, the lookback to forecast horizon ratio was set to 3 times, since it has been frequently considered in previous studies [48, 75].

Similarly, for the ADS dataset, we selected the hazard forecast horizon of 3 timesteps equaling to 3 s, which is in line with the minimum reaction time suggested by the original study [83]. However, due to the small size of the dataset, as described in Section 6.1.2, we could only select the lookback to a forecast horizon ratio of 1, i.e., 3 s. Note that larger lookbacks to forecast horizon ratios significantly increase the total window size and reduce the number of samples available to train and test the model.

Evaluation Metric. As discussed in Section 5, a probabilistic forecast is better suited than a point forecast for critical applications, such as predicting a safety violation, as it provides forecast intervals with attached probabilities. Similar to previous studies in other application domains, where the performance of probabilistic forecasting models has been reported [48, 75], we report the q -Risk metric at multiple quantiles. Equation 4 provides the definition of q -Risk. Intuitively, q -Risk measures the quantile loss [91] (Equation 5) across the entire hazard forecast horizon, normalized over the length of the horizon and over all samples in the test set. Thus, it allows us to compare the safety metric prediction accuracy of models under evaluation at each prediction quantile. Formally, q -Risk is defined as follows:

$$q\text{-Risk} = \frac{2 \sum_{y_t \in \tilde{\Omega}} \sum_{\tau=1}^{\tau_{max}} QL(y_t, \hat{y}(q, t - \tau, \tau), q)}{\sum_{y_t \in \tilde{\Omega}} \sum_{\tau=1}^{\tau_{max}} |y_t|}, \quad (4)$$

where $\tilde{\Omega}$ is the test set, q is the quantile, $\tau = 1, \dots, \tau_{max}$ is the time step counter of the hazard forecast horizon⁸ and QL is the quantile loss function, which is defined in Equation 5.

$$QL(y, \hat{y}, q) = q(y - \hat{y})_+ + (1 - q)(\hat{y} - y)_+, \quad (5)$$

⁷This method is also in line with the evaluation method used by the reference studies of the models evaluated in paper [11, 48, 75, 91]

⁸ $\tau = 1$ and $\tau = \tau_{max}$ correspond to $t + 1$ and $t + h$ in Equation 2, respectively.

Table 2. q-Risk values for different models and quantiles, for the *cte* and *he* safety requirements of the ACT case study and the *cte* safety requirement of the ADS case study, respectively.

Model	Average q-Risk $\pm 0.5 \times CI_{0.95}$						
	$q = 0.005$	$q = 0.025$	$q = 0.05$	$q = 0.5$	$q = 0.95$	$q = 0.975$	$q = 0.995$
<i>ACT_{cte}</i>							
Seq2Seq	0.038 \pm 0.0076	0.046 \pm 0.0056	0.052 \pm 0.0028	0.040 \pm 0.0013	0.026 \pm 0.0014	0.024 \pm 0.0019	0.020 \pm 0.0026
DeepAR	0.004 \pm 0.0006	0.012 \pm 0.0012	0.020 \pm 0.0019	0.062 \pm 0.0035	0.016 \pm 0.0008	0.009 \pm 0.0005	0.003 \pm 0.0002
MQCNN	0.012 \pm 0.0015	0.015 \pm 0.0011	0.018 \pm 0.0016	0.030 \pm 0.0010	0.026 \pm 0.0025	0.023 \pm 0.0022	0.017 \pm 0.0026
TFT	0.001 \pm 0.0001	0.003 \pm 0.0001	0.004 \pm 0.0002	0.012 \pm 0.0002	0.005 \pm 0.0001	0.003 \pm 0.0001	0.001 \pm 0.0001
<i>ACT_{he}</i>							
Seq2Seq	0.208 \pm 0.0420	0.300 \pm 0.0363	0.420 \pm 0.0370	0.841 \pm 0.0204	0.541 \pm 0.0390	0.458 \pm 0.0314	0.330 \pm 0.0351
DeepAR	0.041 \pm 0.0022	0.118 \pm 0.0040	0.199 \pm 0.0056	0.732 \pm 0.0119	0.296 \pm 0.0092	0.196 \pm 0.0077	0.086 \pm 0.0053
MQCNN	0.092 \pm 0.0197	0.171 \pm 0.0161	0.260 \pm 0.0151	0.705 \pm 0.0170	0.418 \pm 0.0372	0.316 \pm 0.0215	0.180 \pm 0.0273
TFT	0.041 \pm 0.0026	0.096 \pm 0.0037	0.140 \pm 0.0041	0.379 \pm 0.0025	0.158 \pm 0.0036	0.112 \pm 0.0034	0.049 \pm 0.0025
<i>ADS_{cte}</i>							
Seq2Seq	0.019 \pm 0.0025	0.056 \pm 0.0064	0.074 \pm 0.0065	0.222 \pm 0.0033	0.137 \pm 0.0073	0.097 \pm 0.0089	0.045 \pm 0.0127
DeepAR	0.007 \pm 0.0004	0.023 \pm 0.0008	0.042 \pm 0.0011	0.204 \pm 0.0031	0.124 \pm 0.0049	0.090 \pm 0.0045	0.048 \pm 0.0041
MQCNN	0.015 \pm 0.0021	0.052 \pm 0.0076	0.070 \pm 0.0080	0.214 \pm 0.0031	0.119 \pm 0.0068	0.077 \pm 0.0060	0.028 \pm 0.0068
TFT	0.023 \pm 0.0040	0.045 \pm 0.0023	0.069 \pm 0.0035	0.226 \pm 0.0066	0.098 \pm 0.0049	0.065 \pm 0.0039	0.023 \pm 0.0006

where $(.)_+ = \max(0, .)$.

Given the safety-critical nature of the decisions that need to be made based on the predicted safety metric values, we reported the q-risk at quantiles that correspond to tail-end values of the prediction interval, namely 90% ($q=0.05, 0.95$), 95% ($q=0.025, 0.975$), 99% ($q=0.005, 0.995$), as well as the median ($q=0.5$) of the prediction distribution.⁹ Recall that, according to the safety metric function defined in Section 3.1 (Equation 1¹⁰), negative values of the safety metric imply no violation of the safety requirement while non-negative safety metric values imply a violation. Further, the higher the negative values, the closer the system is to a safety violation. Moreover, note that safety metric values predicted at higher prediction quantiles ($q > 0.5$) provide the upper bounds of the predicted safety metric value, which are more conservative estimates based on the definition of the safety metric function. Therefore, given their safety-critical application, we expect the predictions that our proposed safety monitors generate at quantiles $q > 0.5$ to be more useful for predicting safety violations, than predictions for other quantiles.

6.3.2 ACT Case Study Results. Table 2 reports the achieved q-Risk values for Seq2Seq, DeepAR, MQCNN, and TFT over 30 repetitions at different quantiles, where the best value for each quantile is written in **bold**.

Overall, we observe that TFT consistently outperforms the other models at all reported quantiles, except in the case of the *he* safety requirement when $q < 0.025$, where TFT and DeepAR both have the lowest q-Risk value. Furthermore, for the *cte* safety requirement, DeepAR is the second best at very high and low ends of the quantile spectrum (specifically, when $q < 0.025$ or $q \geq 0.95$), while yielding the worst accuracy at the median of the forecast probability distribution ($q = 0.5$). Similarly, for the *he* safety requirement, DeepAR is the second best at the ends of the quantile spectrum (except when $q < 0.025$, as mentioned above), while being second to last at the median ($q = 0.5$). We suspect that the fact that DeepAR is an iterative forecasting model, as opposed to the other three

⁹Note that the median of the predicted probability distribution often corresponds to the single value predicted by *point* forecasting methods [11].

¹⁰As mentioned in Section 6.1.1, substituting *cte* with *he* in Equation 3, provides us with the safety metric function definition for the *he* safety requirement.

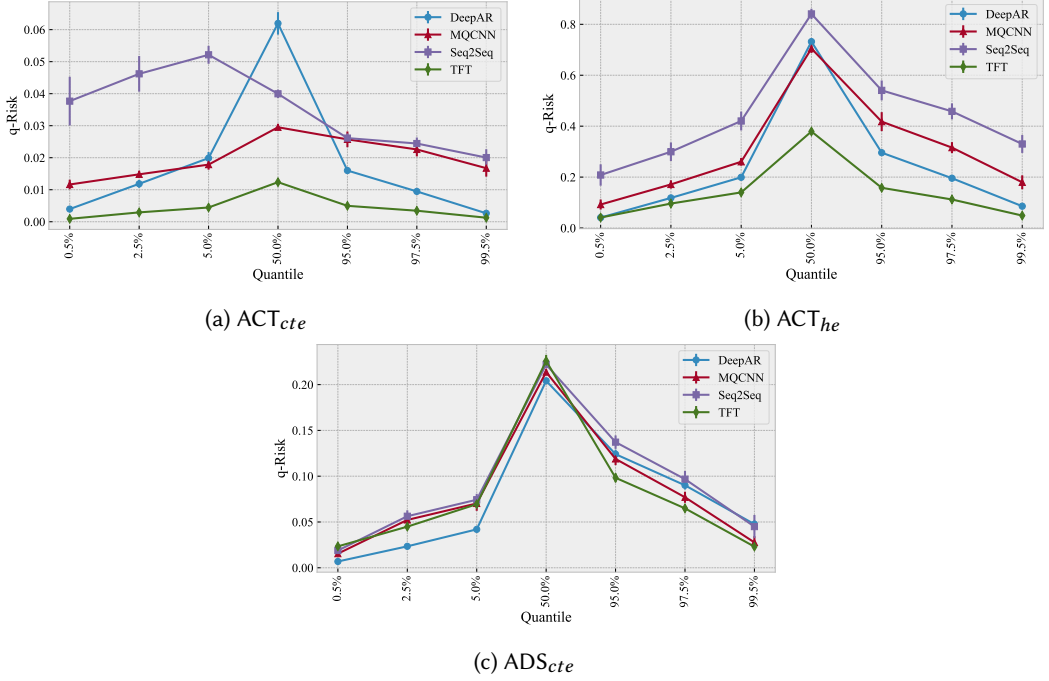


Fig. 3. Average q-Risk values and their corresponding 95% confidence interval ($CI_{0.95}$) for all models over all reported quantiles, for the *cte* and *he* safety requirements of the autonomous taxiing (ACT) case study and the *cte* safety requirement of the autonomous driving (ADS) case study, respectively. Note that the x-axis is *not* drawn to scale, in favor of a more readable presentation.

models which are sequence-to-sequence forecasting models, could explain the large variability in q-Risk values over quantiles. Since iterative forecasting models, such as DeepAR, only predict the target value for the next timestep and use the predicted value to predict the timestep after that (as explained in Section 2), they are prone to accumulating forecasting errors from previous forecast timesteps. DeepAR's error accumulation is more extreme when predicting at quantiles closer to the median ($q = 0.5$), where other models also have higher q-Risk values than for other quantiles.

Figure 3a, which depicts the q-Risk average values and their 95% confidence interval for different models at all measured quantiles, illustrates the large change in performance (average q-Risk values) of DeepAR.

Our visual observations are supported by the statistical comparison results reported in Table 3. Columns *A* and *B* indicate the DL-forecasting models being compared. Columns *p* and \hat{A}_{AB} indicate statistical significance and effect size (as described in Section 6.3.1), respectively, when comparing *A* and *B* in terms of q-Risk at different quantiles q . Given a significance level of $\alpha = 0.01$, for the *cte* safety requirement, we observe that the differences between the best model (TFT) in all quantiles, and the other models are significant. Furthermore, for all quantiles, \hat{A}_{AB} is greater than 0.71 when $B = \text{TFT}$, indicating that the difference between TFT and other models is large. For the *he* safety requirement, TFT and DeepAR are equally the best models, when $q < 0.025$. In this case, we observe that TFT is significantly better than the second-best model, i.e., MQCNN, with a large difference, though that is not the case for DeepAR.

Table 3. Statistical comparison of q-Risk values for different DL-based forecasters at different quantiles q.

Comparison		q-Risk													
A	B	q = 0.005		q = 0.025		q = 0.05		q = 0.5		q = 0.95		q = 0.975		q = 0.995	
		p	\hat{A}_{AB}	p	\hat{A}_{AB}	p	\hat{A}_{AB}	p	\hat{A}_{AB}	p	\hat{A}_{AB}	p	\hat{A}_{AB}	p	\hat{A}_{AB}
ACT _{cte}															
Seq2Seq	DeepAR	8.15×10^{-11}	0.99	6.72×10^{-10}	0.96	3.34×10^{-11}	1.00	3.02×10^{-11}	0.00	1.33×10^{-10}	0.98	3.34×10^{-11}	1.00	3.02×10^{-11}	1.00
Seq2Seq	MQCNN	2.00×10^{-5}	0.82	1.56×10^{-8}	0.93	3.02×10^{-11}	1.00	7.39×10^{-11}	0.99	4.04×10^{-1}	0.56	9.33×10^{-2}	0.63	4.51×10^{-2}	0.65
Seq2Seq	TFT	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
DeepAR	MQCNN	2.87×10^{-10}	0.03	1.41×10^{-4}	0.21	7.48×10^{-2}	0.63	3.02×10^{-11}	1.00	3.20×10^{-9}	0.05	3.02×10^{-11}	0.00	3.02×10^{-11}	0.00
DeepAR	TFT	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
MQCNN	TFT	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
ACT _{he}															
Seq2Seq	DeepAR	8.48×10^{-9}	0.93	4.50×10^{-11}	1.00	3.02×10^{-11}	1.00	4.62×10^{-10}	0.97	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
Seq2Seq	MQCNN	2.60×10^{-5}	0.82	1.60×10^{-7}	0.89	1.29×10^{-9}	0.96	9.92×10^{-11}	0.99	4.64×10^{-5}	0.81	1.56×10^{-8}	0.93	1.87×10^{-7}	0.89
Seq2Seq	TFT	8.48×10^{-9}	0.93	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
DeepAR	MQCNN	8.15×10^{-5}	0.20	8.20×10^{-7}	0.13	8.48×10^{-9}	0.07	1.70×10^{-2}	0.68	3.08×10^{-8}	0.08	6.72×10^{-10}	0.04	2.44×10^{-9}	0.05
DeepAR	TFT	8.77×10^{-1}	0.49	4.18×10^{-9}	0.94	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.69×10^{-11}	1.00
MQCNN	TFT	8.66×10^{-5}	0.80	3.16×10^{-10}	0.97	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
ADS _{cte}															
Seq2Seq	DeepAR	4.98×10^{-11}	0.99	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	2.39×10^{-8}	0.92	3.67×10^{-3}	0.72	5.40×10^{-1}	0.55	3.39×10^{-2}	0.34
Seq2Seq	MQCNN	4.36×10^{-2}	0.65	2.12×10^{-1}	0.59	9.05×10^{-2}	0.63	1.25×10^{-4}	0.79	5.56×10^{-4}	0.76	1.17×10^{-3}	0.74	2.50×10^{-3}	0.73
Seq2Seq	TFT	2.06×10^{-1}	0.40	3.02×10^{-2}	0.67	5.79×10^{-1}	0.54	8.53×10^{-1}	0.49	2.44×10^{-9}	0.95	3.96×10^{-8}	0.91	1.58×10^{-1}	0.61
DeepAR	MQCNN	8.99×10^{-11}	0.01	3.02×10^{-11}	0.00	3.02×10^{-11}	0.00	3.02×10^{-4}	0.22	6.35×10^{-2}	0.64	5.56×10^{-4}	0.76	1.75×10^{-5}	0.82
DeepAR	TFT	3.02×10^{-11}	0.00	3.02×10^{-11}	0.00	3.02×10^{-11}	0.00	2.38×10^{-7}	0.11	9.26×10^{-9}	0.93	7.12×10^{-9}	0.94	3.34×10^{-11}	1.00
MQCNN	TFT	1.77×10^{-3}	0.26	6.00×10^{-1}	0.54	1.54×10^{-1}	0.39	1.24×10^{-3}	0.26	2.88×10^{-6}	0.85	1.17×10^{-3}	0.74	7.73×10^{-2}	0.37

6.3.3 ADS Case Study Results. The q-Risk values achieved by Seq2Seq, DeepAR, MQCNN, and TFT over 30 repetitions at different quantiles are reported in Table 2. Note that the best value for each quantile is written in **bold**.

Overall, we observe that at each quantile, the difference between the most accurate and least accurate models are less than what is observed for the ACT case study results at similar quantiles (compare Figure 3c vs Figure 3a and Figure 3b). We believe that this is due to the sample size of the ADS dataset which is significantly lower than the size of the ACT dataset, as discussed in Section 6.1.1, where a model like TFT is expected to suffer the most, as it is a transformer-based model, which has been shown to require significantly more training data than other DL-based models such as CNNs in vision tasks [21]. Nonetheless, we observe that TFT achieves the lowest q-Risk values (most accurate predictions) when $q > 0.5$. Whereas, for $q \leq 0.5$, DeepAR outperforms other models.

Our statistical test results (Table 3), confirm that TFT significantly outperforms other models when $0.5 < q < 0.995$, with a large effect size as the corresponding \hat{A}_{AB} values are greater than 0.71. However, at $q = 0.995$, using the Mann-Whitney U-test, when $A \in \{\text{Seq2Seq, MQCNN}\}$ and $B \in \{\text{TFT}\}$, indicate that their difference is not statistically significant (with a confidence level of 95%), as p-values are larger than 0.05. Therefore, at $q = 0.995$, we conclude that sequence-to-sequence models, i.e., TFT, MQCNN and Seq2Seq, all equally yield the best safety metric prediction accuracy. Finally, we observe that for $q \leq 0.5$, DeepAR consistently outperforms other models with a large effect size.

6.3.4 Discussion. Given the results of the ACT case study (Section 6.3.2), we observed that, for probabilistic prediction of both *cte* and *he* safety metrics, at all measured quantiles q , with a practical window configuration, i.e., hazard forecast horizon of 3 s (which correlates to the minimum reaction time, as discussed in Section 6.3.1) and lookback to forecast horizon ratio of three (which is similar to the ratio used by the literature in multiple forecasting problems [48, 75]), TFT yields significantly more accurate quantile forecasts than Seq2Seq, DeepAR, and MQCNN. Our observations for the *he* safety requirement is the same as *cte*, for $q \geq 0.025$. Whereas, for $q < 0.025$, TFT and DeepAR both yield the highest time series prediction accuracy. Therefore, we can conclude that for the ACT

dataset, where the dataset size is large and contains numerous safety violations, TFT is the best model or one of the best models to be used for probabilistic forecasting of the safety metric values at all quantiles, given a practical window configuration.

For the ADS case study results (Section 6.3.3), we observed that for probabilistic safety metric prediction, again with a practical window configuration, i.e., hazard forecast horizon of 3 s and a lookback to forecast horizon ratio of one, as discussed in Section 6.3.1, TFT yields significantly more accurate predictions when $0.5 < q < 0.995$. At $q = 0.995$, all sequence-to-sequence models, i.e., TFT, MQCNN and Seq2Seq, yield the lowest q -Risk value. Finally, DeepAR yields significantly more accurate predictions when $q \leq 0.5$. Therefore, for the ADS case study, where the size of the dataset is a fraction of the ACT dataset, as discussed in Section 6.1.2, TFT is the best or one of the best models to be used for safety metric forecasting when $q > 0.5$. Though DeepAR is the most accurate model when $q \leq 0.5$, as discussed in Section 6.3.1, given the definition of the safety metric (Equation 1), forecasts for quantiles $q > 0.5$ are more important for safety violation prediction.

For the ACT case study, where the size of the dataset is large, given a practical window configuration, i.e., hazard forecast horizon of 3 s and lookback to forecast horizon ratio of 3, TFT is more suitable than Seq2Seq, DeepAR and MQCNN, for probabilistic safety metric forecasting over all reported quantiles, for both *cte* and *he* safety requirements.

For the ADS case study, where the size of the dataset is small, given a practical window configuration of $h = 3$ s and $cm = 1$, DeepAR is significantly more accurate than TFT, MQCNN, and Seq2Seq, for $q \leq 0.5$. However, TFT is the most accurate model, among the evaluated models, when $q > 0.5$, which is a more important quantile range in a safety monitoring context, as discussed in Section 6.3.1.

Therefore, TFT is the most accurate model, when $q > 0.5$, for both case studies.

6.4 RQ₂: Safety Violation Prediction Accuracy

In this section, similar to Section 6.3, first we provide the details of our evaluation methodology to answer RQ₂ (Section 6.4.1). Then, we present the results for the ACT and ADS case studies (Section 6.4.2 and Section 6.4.3, respectively). Finally, we draw conclusions from the results of both case studies (Section 6.4.4) and present our answer to RQ₂.

6.4.1 Methodology. To answer RQ₂, we reuse the models trained to answer RQ₁ with the aim of predicting safety violations. To do so, we applied the safety violation function (Equation 3) to the safety metric values predicted by the forecasting models, as reported in RQ₁. A non-negative safety function value implies a safety violation. We compared the predicted safety violations with the true safety violation value of the test samples used in RQ₁. True safety violation values are calculated by applying the safety violation function to the true safety metric values of the test samples.

Evaluation Metric. To report the safety violation prediction accuracy of the models, we report *Precision* ($Pr = TP / (TP + FP)$) and *Recall* ($Re = TP / (TP + FN)$), where *true positives* (TP), *false positive* (FP), and *false negatives* (FN) are the correct, false, and missed safety violation predictions, respectively. Note that Precision measures the fraction of correct warnings that a safety monitor raises, whereas Recall measures the fraction of safety violations that a safety monitor can successfully predict [83]. We further compute and report the F_β score [9] as a weighted balance between precision and recall, with $\beta = 3.0$ ($F_\beta = \frac{10 \cdot Precision \times Recall}{9 \cdot Precision + Recall}$), granting higher importance to Recall as compared to Precision, as false negatives have severe consequences for safety-critical systems [25, 83]. We recall that false negatives, in the context of safety-critical systems, are safety violations that were not predicted by the safety monitor and thus could lead to system hazards. In contrast, false positives, although an

important consideration for the design of the safety monitor, lead to inconvenience or inefficiencies for the users, which are less harmful than safety violations. For the ACT system specifically, false positives could lead to the disengagement of the autonomous taxiing operation or emergency stops, which could lead to delayed flight operations and schedules.

Recall that the predictions for quantiles $q > 0.5$ are more important than lower quantiles, as they provide more conservative estimates of the safety metric, as discussed in Section 6.3.1.

6.4.2 ACT Case Study Results. Figure 4a and Figure 4b illustrate the False Negative (FN) and False Positive (FP) values for all the models at the reported quantiles, respectively.¹¹ Overall, we observe that with an increase in the prediction quantile, the FN value decreases while the FP value increases. This is expected based on the definitions of the safety metric and the safety violation function (Equation 1 and Equation 3, respectively). Recall that, as mentioned in Section 6.4.1, the predictions at quantiles $q > 0.5$ provide more conservative estimates of the safety metric values. Therefore, the higher the prediction quantile, the higher the probability of the safety monitor correctly predicting safety violations (lower FN) and raising false alarms (higher FP).

As mentioned in Section 6.4.1, the priority in safety-critical applications is having the lowest FN value possible (since FNs lead to system hazards), while having a reasonably low FP value (since FPs lead to inefficiencies) is the second priority. Therefore, for the ACT and ADS case studies targeted in this paper, hereafter we focus on the results for prediction quantiles $q \geq 0.5$, for which the FN and FP values are reported in Table 4. We have provided the values at other quantiles in the supporting material (Section 6.9). Table 4 reports the averages and their 95% confidence intervals for the evaluation metrics, namely Pr, Re, and F_3 .

Overall, Precision for all models drops as q increases, whereas Recall increases, for both *cte* and *he* safety requirements. This general trend is in line with the FN and FP trends above. Since estimates become more conservative, more safety violations are correctly predicted (Recall \uparrow) while the proportion of false alarms increases (Precision \downarrow). For the *he* safety requirement, since the proportion of time steps including a safety violation is lower than that of the *cte* safety requirement (Section 6.1.1), we expect and observe that the precision scores of the models recorded for *he* are lower than the scores recorded for *cte*. In the case of the *cte* safety requirement, note that TFT and DeepAR reach a Recall value of 1.0 at a $q = 0.995$, thus indicating all of the safety violations are correctly predicted, which is also confirmed by the low corresponding FN values; Seq2Seq and MQCNN, on the other hand, yield the lowest Recall scores (high FN values). However, we observe that DeepAR has the lowest Precision among the models for $q \geq 0.5$, indicating a higher fraction of false alarms, which is also confirmed by the fact that the FP value for DeepAR is an order of magnitude greater than that of other models. In contrast, TFT has a precision above 0.92, for $q \geq 0.5$, which makes it a more suitable choice for safety monitoring than DeepAR. In the case of the *he* safety requirement, we observe that TFT consistently yields the highest Precision and Recall scores for $q \geq 0.5$, whereas the Recall of other models is 20 – 40% lower. The superiority of TFT over other models is further confirmed by the reported F_3 scores, where TFT consistently has the highest F_3 scores for $q \geq 0.5$, for both the *cte* and *he* safety requirements. Note that the highest F_3 scores for each quantile are highlighted in bold. Our visual observations are supported by the statistical comparison results reported in Table 6. Columns *A* and *B* indicate the DL-forecasting models being compared. Given a significance level of $\alpha = 0.01$, we observe that the differences between TFT and the other models are significant for all quantiles. Furthermore, for all quantiles, \hat{A}_{AB} is greater than 0.71 when $B = \text{TFT}$, indicating that the difference between TFT and other models is large.

¹¹Note that the x-axis, i.e., quantile q , is not drawn up to scale for better readability.

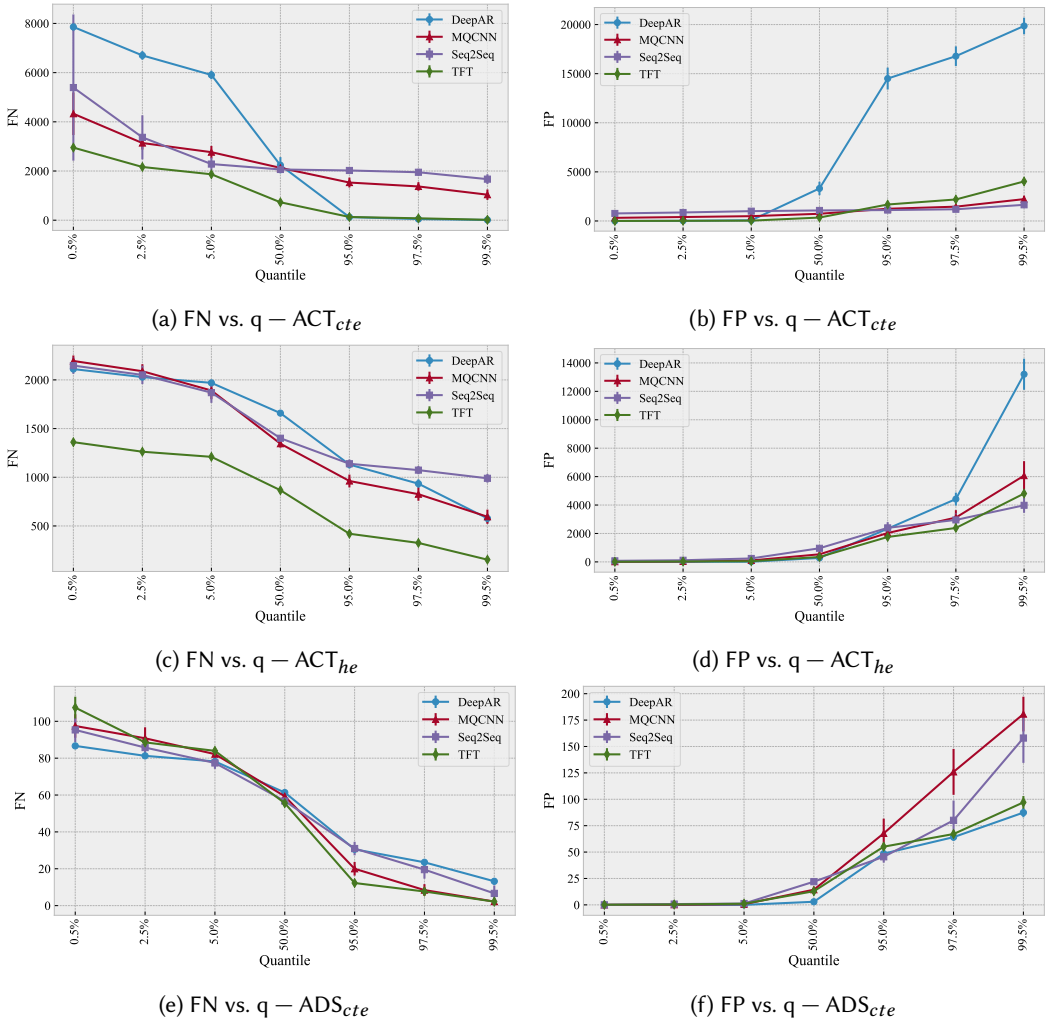


Fig. 4. False Negatives (FN) and False Positives (FP) at different prediction quantiles (q), for cte and he safety requirements of the ACT case study and the cte safety requirement of the ADS case study, respectively. Note that the x-axis is *not* drawn to scale, in favor of a more readable presentation.

6.4.3 ADS Case Study Results. Similar to the ACT case study results (Section 6.4.2), Figure 4e and Figure 4f, indicate that FN decreases and FP increases with increase in prediction quantile. As explained in Section 6.4.2, we will focus on the results for prediction quantiles $q \geq 0.5$. We report the averages and 95% confidence intervals of the FN, FP, Precision, Recall and F_3 score values in Table 5.

Overall, in line FP and FN trends, Precision for all models drops when q increases, while Recall increases, as discussed in Section 6.4.2. Note that TFT and MQCNN yield the highest Recall score of 0.985 at $q = 0.995$, while for the same quantile, MQCNN yields the lowest Precision score. We further observe for $q \geq 0.5$, TFT is the second best performing model in terms of Precision score while yielding the highest Recall score. This results in TFT achieving the highest F_3 score over all

Table 4. Various safety violation prediction accuracy metric values for different models and quantiles, for *cte* and *he* (rows highlighted in gray) safety requirements.

Model	Average Metric Value $\pm 0.5 \times CI_{0.95}$			
	$q = 0.5$	$q = 0.95$	$q = 0.975$	$q = 0.995$
FN				
Seq2Seq	2062.9 \pm 111.5	2022.0 \pm 114.5	1949.7 \pm 142.3	1669.6 \pm 198.9
	1400.9 \pm 32.9	1138.8 \pm 43.8	1073.4 \pm 42.0	989.4 \pm 44.7
DeepAR	2236.5 \pm 329.5	124.5 \pm 73.9	46.0 \pm 31.1	10.3 \pm 7.6
	1658.8 \pm 20.3	1130.5 \pm 40.0	934.1 \pm 45.2	577.7 \pm 46.2
MQCNN	2128.0 \pm 199.9	1531.2 \pm 201.6	1372.3 \pm 174.0	1037.6 \pm 211.2
	1344.7 \pm 40.3	961.7 \pm 65.9	826.2 \pm 67.1	594.6 \pm 73.2
TFT	729.2 \pm 27.2	131.4 \pm 17.6	78.1 \pm 11.6	16.8 \pm 2.5
	866.9 \pm 20.4	420.0 \pm 24.3	325.7 \pm 22.3	153.9 \pm 15.6
FP				
Seq2Seq	1073.4 \pm 129.0	1108.9 \pm 150.6	1191.4 \pm 175.6	1642.5 \pm 291.9
	958.8 \pm 165.1	2395.2 \pm 402.2	2953.1 \pm 460.7	3982.8 \pm 528.9
DeepAR	3311.0 \pm 684.6	14501.4 \pm 1113.0	16781.9 \pm 998.2	19858.7 \pm 842.3
	276.4 \pm 24.8	2337.9 \pm 240.0	4420.1 \pm 455.4	13203.1 \pm 1093.1
MQCNN	732.8 \pm 110.2	1243.0 \pm 160.6	1457.3 \pm 186.4	2224.9 \pm 336.1
	533.7 \pm 76.3	2032.3 \pm 363.5	3118.9 \pm 522.3	6069.3 \pm 1018.0
TFT	348.4 \pm 13.3	1677.3 \pm 61.3	2190.1 \pm 78.9	4026.0 \pm 138.0
	344.6 \pm 19.1	1752.8 \pm 73.0	2385.5 \pm 102.9	4808.9 \pm 221.0
Precision				
Seq2Seq	0.978 \pm 0.0026	0.977 \pm 0.0030	0.975 \pm 0.0035	0.967 \pm 0.0056
	0.496 \pm 0.0304	0.343 \pm 0.0254	0.309 \pm 0.0237	0.261 \pm 0.0233
DeepAR	0.936 \pm 0.0120	0.773 \pm 0.0138	0.746 \pm 0.0116	0.712 \pm 0.0089
	0.693 \pm 0.0156	0.338 \pm 0.0186	0.242 \pm 0.0162	0.118 \pm 0.0082
MQCNN	0.985 \pm 0.0022	0.975 \pm 0.0031	0.970 \pm 0.0036	0.956 \pm 0.0062
	0.646 \pm 0.0261	0.422 \pm 0.0330	0.343 \pm 0.0304	0.249 \pm 0.0335
TFT	0.993 \pm 0.0003	0.967 \pm 0.0012	0.957 \pm 0.0015	0.924 \pm 0.0024
	0.804 \pm 0.0081	0.515 \pm 0.0085	0.451 \pm 0.0087	0.308 \pm 0.0086
Recall				
Seq2Seq	0.958 \pm 0.0023	0.959 \pm 0.0023	0.960 \pm 0.0029	0.966 \pm 0.0041
	0.383 \pm 0.0145	0.499 \pm 0.0193	0.527 \pm 0.0185	0.564 \pm 0.0197
DeepAR	0.954 \pm 0.0067	0.997 \pm 0.0015	0.999 \pm 0.0006	1.000 \pm 0.0002
	0.270 \pm 0.0089	0.502 \pm 0.0176	0.589 \pm 0.0199	0.746 \pm 0.0203
MQCNN	0.956 \pm 0.0041	0.969 \pm 0.0041	0.972 \pm 0.0036	0.979 \pm 0.0043
	0.408 \pm 0.0178	0.577 \pm 0.0290	0.636 \pm 0.0295	0.738 \pm 0.0322
TFT	0.985 \pm 0.0006	0.997 \pm 0.0004	0.998 \pm 0.0002	1.000 \pm 0.0001
	0.618 \pm 0.0090	0.815 \pm 0.0107	0.857 \pm 0.0098	0.932 \pm 0.0069
F ₃				
Seq2Seq	0.960 \pm 0.0019	0.960 \pm 0.0020	0.962 \pm 0.0024	0.966 \pm 0.0032
	0.390 \pm 0.0125	0.472 \pm 0.0122	0.486 \pm 0.0107	0.498 \pm 0.0103
DeepAR	0.952 \pm 0.0049	0.969 \pm 0.0012	0.966 \pm 0.0015	0.961 \pm 0.0015
	0.287 \pm 0.0090	0.476 \pm 0.0135	0.510 \pm 0.0129	0.481 \pm 0.0146
MQCNN	0.959 \pm 0.0036	0.969 \pm 0.0035	0.972 \pm 0.0030	0.976 \pm 0.0034
	0.423 \pm 0.0170	0.549 \pm 0.0212	0.576 \pm 0.0182	0.595 \pm 0.0159
TFT	0.986 \pm 0.0005	0.994 \pm 0.0002	0.994 \pm 0.0001	0.992 \pm 0.0003
	0.633 \pm 0.0084	0.770 \pm 0.0079	0.785 \pm 0.0064	0.774 \pm 0.0048

quantiles $q \geq 0.5$. Moreover, we observe that DeepAR yields the lowest Recall and F₃ score over all quantiles $q \geq 0.5$.

This case study thus confirms the superiority of TFT over Seq2Seq, MQCNN, and DeepAR, when $q > 0.5$, in terms of F₃ score, based on the results of our statistical analysis, reported in Table 6.

Table 5. Various safety violation prediction accuracy metric values for different models and quantiles, for the ADS case study.

Model	Average Metric Value $\pm 0.5 \times CI_{0.95}$			
	$q = 0.5$	$q = 0.95$	$q = 0.975$	$q = 0.995$
FN				
Seq2Seq	56.8 \pm 1.4	31.0 \pm 3.7	19.5 \pm 4.9	6.6 \pm 4.3
DeepAR	61.4 \pm 0.7	30.6 \pm 1.2	23.5 \pm 1.5	13.2 \pm 1.8
MQCNN	59.3 \pm 2.0	20.0 \pm 3.7	8.4 \pm 3.2	2.2 \pm 1.8
TFT	55.6 \pm 1.8	12.3 \pm 1.6	7.7 \pm 1.2	2.2 \pm 0.3
FP				
Seq2Seq	22.0 \pm 1.9	45.5 \pm 5.5	80.1 \pm 18.7	158.0 \pm 23.5
DeepAR	2.9 \pm 0.3	48.5 \pm 3.1	64.1 \pm 3.2	87.4 \pm 4.4
MQCNN	14.4 \pm 2.5	67.7 \pm 14.0	126.0 \pm 21.7	180.7 \pm 16.4
TFT	12.9 \pm 1.5	55.1 \pm 3.1	67.1 \pm 3.0	97.1 \pm 6.0
Precision				
Seq2Seq	0.810 \pm 0.0128	0.728 \pm 0.0153	0.649 \pm 0.0380	0.505 \pm 0.0449
DeepAR	0.968 \pm 0.0035	0.711 \pm 0.0120	0.663 \pm 0.0101	0.610 \pm 0.0113
MQCNN	0.867 \pm 0.0179	0.676 \pm 0.0311	0.555 \pm 0.0405	0.460 \pm 0.0282
TFT	0.880 \pm 0.0105	0.714 \pm 0.0101	0.679 \pm 0.0085	0.605 \pm 0.0142
Recall				
Seq2Seq	0.619 \pm 0.0092	0.792 \pm 0.0246	0.869 \pm 0.0329	0.955 \pm 0.0289
DeepAR	0.588 \pm 0.0047	0.795 \pm 0.0079	0.842 \pm 0.0097	0.911 \pm 0.0121
MQCNN	0.602 \pm 0.0131	0.866 \pm 0.0248	0.943 \pm 0.0216	0.985 \pm 0.0122
TFT	0.627 \pm 0.0121	0.918 \pm 0.0109	0.949 \pm 0.0080	0.985 \pm 0.0022
F_3				
Seq2Seq	0.634 \pm 0.0084	0.784 \pm 0.0197	0.833 \pm 0.0215	0.866 \pm 0.0144
DeepAR	0.612 \pm 0.0046	0.785 \pm 0.0062	0.820 \pm 0.0079	0.868 \pm 0.0096
MQCNN	0.620 \pm 0.0119	0.838 \pm 0.0169	0.874 \pm 0.0099	0.880 \pm 0.0028
TFT	0.645 \pm 0.0113	0.892 \pm 0.0086	0.912 \pm 0.0057	0.927 \pm 0.0031

We further observe that, when $q > 0.5$, TFT is significantly better than other models with a high effect size. However, for $q = 0.5$, we observe that TFT is not significantly better than Seq2Seq, though it is still outperforming DeepAR and MQCNN significantly. Moreover, the effect size of the difference between the F_3 scores of TFT and MQCNN when $q = 0.5$, i.e., \hat{A}_{AB} , when $A = \text{TFT}$ and $B = \text{MQCNN}$, is $0.64 < \hat{A}_{AB} = 0.69 < 0.71$, indicating a medium effect size.

6.4.4 Discussion. Based on the results of the ACT case study (Section 6.4.2), where the dataset is large and includes numerous safety violations, we conclude that, for probabilistic prediction of the safety violation, at all measured quantiles q , with again a hazard forecast horizon of 3 s and a lookback to forecast horizon ratio of 3, as discussed in Section 6.3.2, TFT is significantly more accurate than Seq2Seq, DeepAR, and MQCNN, for both *cte* and *he* safety requirements.

Based on the ADS case study results (Section 6.4.3), where the size of the dataset is much smaller than the ACT dataset, we observe that, given the hazard forecast horizon of 3 s (equal to the minimum reaction time) and a lookback to forecast horizon ratio of 1, as discussed in Section 6.3.1, TFT is significantly more accurate than Seq2Seq, DeepAR, and MQCNN with a large effect size, when $q > 0.5$. At $q = 0.5$, TFT and Seq2Seq both yield the most accurate predictions, while significantly outperforming DeepAR and MQCNN with a large and medium effect size, respectively. Therefore, we conclude that TFT is the most suitable model, among all evaluated models, for safety metric forecasting, for all reported quantiles $q \geq 0.5$. Predictions for $q \leq 0.5$ are in any case not sufficiently accurate regardless of the model employed and therefore comparisons for such q values are not of practical utility.

Table 6. Statistical comparison of F_3 score values for different DL-based forecasters at quantiles $q \geq 0.5$, for *cte* and *he* safety requirements of the ACT case study and the *cte* safety requirement of the ADS case study, respectively.

Comparison		F_3 score							
A	B	$q = 0.5$		$q = 0.95$		$q = 0.975$		$q = 0.995$	
		p	\hat{A}_{AB}	p	\hat{A}_{AB}	p	\hat{A}_{AB}	p	\hat{A}_{AB}
TFT	Seq2Seq	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
TFT	DeepAR	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
TFT	MQCNN	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
Seq2Seq	DeepAR	8.99×10^{-2}	0.68	1.16×10^{-7}	0.10	1.60×10^{-3}	0.26	3.64×10^{-2}	0.66
Seq2Seq	MQCNN	7.62×10^{-1}	0.48	2.28×10^{-5}	0.18	1.25×10^{-5}	0.17	3.37×10^{-5}	0.19
DeepAR	MQCNN	3.64×10^{-2}	0.34	1.37×10^{-1}	0.39	4.71×10^{-4}	0.24	7.69×10^{-8}	0.10
ACT_{he}									
TFT	Seq2Seq	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
TFT	DeepAR	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
TFT	MQCNN	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00	3.02×10^{-11}	1.00
Seq2Seq	DeepAR	8.99×10^{-11}	0.99	5.59×10^{-1}	0.46	5.26×10^{-4}	0.24	3.64×10^{-2}	0.66
Seq2Seq	MQCNN	3.34×10^{-3}	0.28	3.01×10^{-7}	0.11	1.55×10^{-9}	0.05	1.96×10^{-10}	0.02
DeepAR	MQCNN	6.70×10^{-11}	0.01	8.20×10^{-7}	0.13	6.53×10^{-7}	0.13	7.39×10^{-11}	0.01
ADS_{cte}									
TFT	Seq2Seq	1.41×10^{-1}	0.61	2.92×10^{-9}	0.95	2.19×10^{-8}	0.92	1.61×10^{-11}	1.00
TFT	DeepAR	2.58×10^{-5}	0.82	3.01×10^{-11}	1.00	3.01×10^{-11}	1.00	4.96×10^{-11}	0.99
TFT	MQCNN	9.88×10^{-3}	0.69	4.98×10^{-7}	0.88	1.98×10^{-8}	0.92	7.85×10^{-12}	1.00
Seq2Seq	DeepAR	2.12×10^{-4}	0.78	5.89×10^{-1}	0.46	2.77×10^{-1}	0.58	1.52×10^{-1}	0.61
Seq2Seq	MQCNN	6.25×10^{-2}	0.64	9.79×10^{-5}	0.21	3.25×10^{-2}	0.34	6.87×10^{-1}	0.47
DeepAR	MQCNN	2.97×10^{-1}	0.42	1.49×10^{-6}	0.14	1.74×10^{-8}	0.08	6.60×10^{-3}	0.30

For the ACT case study, where the size of the dataset is large, given a hazard forecast horizon of 3 s and a lookback to a forecast horizon ratio of 3 for safety violation prediction (minimum reaction time), TFT is significantly more accurate, with a large effect size, than Seq2Seq, DeepAR, and MQCNN for all quantiles $q \geq 0.5$, for both the *cte* and *he* safety requirements.

For the ADS case study, where the dataset size is much smaller, given a practical window configuration, i.e., $h = 3$ s and $cm = 1$, TFT is the most suitable model, among all evaluated models, for probabilistic safety metric forecasting, for all quantiles $q \geq 0.5$. Predictions for $q \leq 0.5$ are poor for all models.

Therefore, when $q > 0.5$, TFT is consistently the best model for both case studies.

6.5 RQ₃: Prediction Accuracy Sensitivity Analysis

In Section 6.5.1, we provide the details of our evaluation methodology as it relates to answering RQ₃. As discussed in Section 6.3.1, due to the low number of samples available in the dataset of the ADS case study, we were only able to study the effect of varying window sizes on prediction accuracy for the two safety requirements of the ACT case study, for which we report the results in Section 6.5.2. Finally, we present our answer to RQ₃, based on the results presented in Section 6.5.2.

6.5.1 Methodology. We have answered RQ₁ and RQ₂ based on the minimum reaction time (3 s) for hazard forecast horizon and a commonly used lookback horizon that is three times longer (12 s). However, the size of the hazard forecast and lookback horizons are design choices of the system developers. To investigate the impact of window configuration, for each forecasting model, in terms

of safety metric and violation accuracy, inference latency, and computations resource usage, we also assessed the tuned models with different hazard forecast horizons and lookback-to-forecast horizon values.

Particularly, we selected the hazard forecast horizon from $\{3, 12\}$ and the ratio of lookback to forecast horizon window size, also known as *context multiplier* (cm), from $\{1, 3, 9\}$. Thus, we studied the following forecast horizon window size and context multiplier combinations (h, cm): (3, 1), (3, 3), (3, 9), (12, 1), (12, 3), and (12, 9). Note that the smallest total window size¹² is $3 + 1 \times 3 = 6$ s, whereas the largest total window size is $12 + 9 \times 12 = 120$ s. In our preliminary experiments, we observed that increasing the forecast horizon increases the likelihood that samples include safety violations. We further observed that for forecast horizons longer than 12 s, the distribution of test samples becomes highly imbalanced, leading to biased comparisons of safety violation prediction accuracy metrics between short ($h = 3$ s) and very long ($h > 12$ s) forecast horizons. Furthermore, we would not have been able to study the effect of varying cm on very large forecast horizons since their total window size on higher cm values would become longer than the maximum training sample length. Thus, for the evaluation of RQ₃, we did not include forecast horizons larger than 12 s. The results for a representative instance of our preliminary experiments with long forecast horizons, i.e., a forecast horizon of 36 s and context multiplier of 1, are included in our supporting material (see Section 6.9).

6.5.2 Results. As discussed in Section 6.4.2, we are interested in the predictions at quantiles that primarily detect as many hazards as possible while having a low number of false alarms. Due to the safety-critical nature of our problem, we only focus here on predictions at $q = 0.995$ since it is the most conservative measured prediction quantile (Section 6.4.2). Recall that when comparing models, a lower q-Risk value implies a more accurate model at predicting safety metric values.

From Figure 5a, we observe that for the *cte* safety requirement, given a fixed cm value, increasing the forecast horizon h leads to higher q-Risk values and thus lower accuracy. Further, we can also see that, in contrast, for a fixed forecast horizon, increasing cm does not significantly improve q-Risk. We observe that TFT outperforms all models (i.e., has the lowest q-Risk value) at the forecast horizon of 3, across all cm values. At longer forecast horizons ($h = 12$), DeepAR outperforms all other models. Moreover, note that due to the iterative forecasting architecture of DeepAR, it is prone to accumulating forecasting errors. Thus, its confidence interval increases significantly with the increase in forecast horizon. For instance, at $cm = 3$, the confidence interval of DeepAR includes the q-Risk values for both TFT and MQCNN.

For the *he* safety requirement, similar to *cte*, we observe that given a fixed cm value, increasing h leads to higher q-Risk values, except for MQCNN at $cm = 9$, where q-Risk decreases when increasing the forecast horizon.

In terms of precision for the *cte* safety requirement (Figure 6a), we observe that when increasing forecast horizon, given a constant cm value, the precision of models with a sequence-to-sequence architecture (TFT, MQCNN, and Seq2Seq) drops while DeepAR's precision increases, most particularly at $cm = 1$. A similar trend is observed for the *he* safety requirement (Figure 6b) with the difference that the models reach a lower precision score than similar models evaluated on the *cte* safety requirement data. This can be explained by the lower proportion of samples including safety violations for *he*, when compared to *cte* (as discussed in Section 6.3.2). Further note that, for both *cte* and *he*, given a constant forecast horizon, increasing cm can slightly improve the precision of the model, most particularly for DeepAR.

In terms of recall for the *cte* safety requirement (Figure 7a), we observe that all models yield a lower recall when increasing forecast horizon while increasing cm does not significantly improve

¹²Total window size = $h + (cm \times h)$

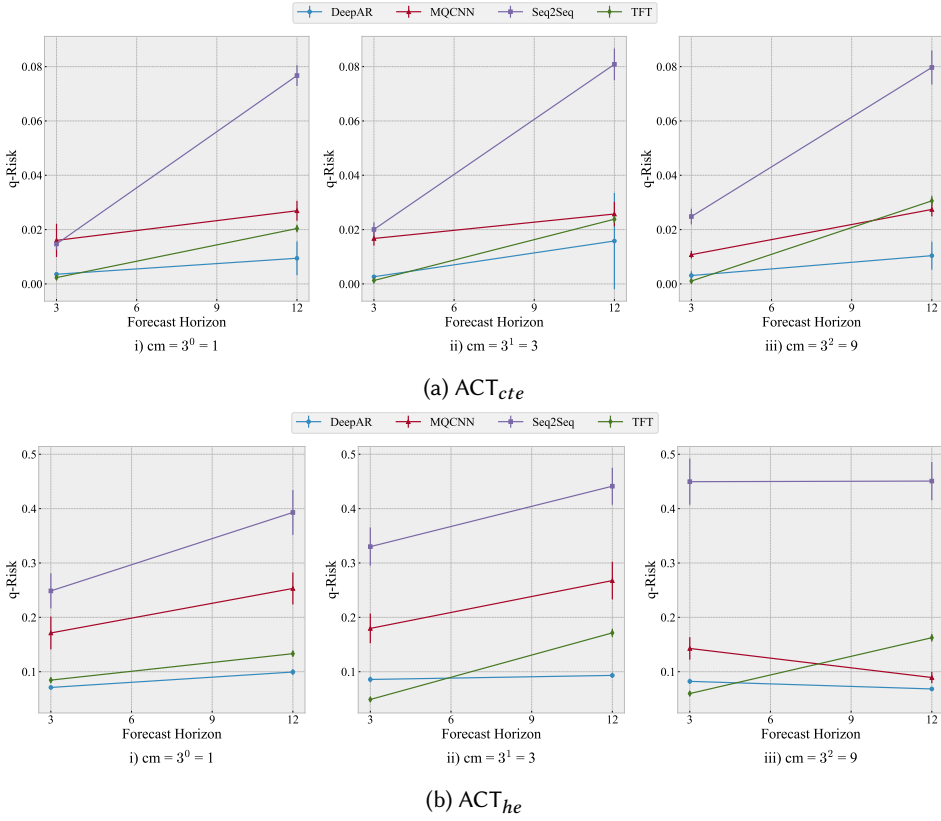


Fig. 5. Safety metric prediction accuracy metrics for various window configurations of DeepAR, MQCNN, Seq2Seq and TFT, for the cte and he safety requirements, respectively.

recall given the same forecast horizon. Note that TFT and DeepAR similarly reach the highest recall at $h = 3$. Whereas, at $h = 12$, TFT experiences a larger drop in recall and is outperformed by DeepAR. Concerning the effect of cm , we observe that MQCNN experiences the most improvement when cm changes from 3 to 9. At the same time, we observe that the recall score of DeepAR and Seq2Seq drops when cm increases given a fixed forecast horizon. Finally, the recall score of TFT is not significantly affected by changes in cm values. Similarly, for the he safety requirement, TFT reaches the highest recall score at $h = 3$. However, instead of all models experiencing a drop in recall score with increases in forecast horizon, the recall score for DeepAR increases in all cm . MQCNN also experiences a recall score increase with an increase in h at $cm = 9$. To conclude, the effect of cm is similar to that observed for the cte safety requirement, with the exception of the differences mentioned above.

In terms of overall safety violation prediction accuracy, i.e., F_3 (Figure 8), for both cte and he safety requirements, we observe that TFT yields the highest F_3 score when $h = 3$ for all values of cm , except for he when $cm = 1$. However, we observe that F_3 is much lower than the scores reached for other cm values, suggesting that the small lookback horizon does not contain sufficient information for the models to generate accurate forecasts. Note that for the same window configuration, DeepAR yields the lowest F_3 score, which is mainly due to its very low precision. Nevertheless, given the increase in precision and a smaller decrease in recall for cte , as well as an increase for he , the F_3

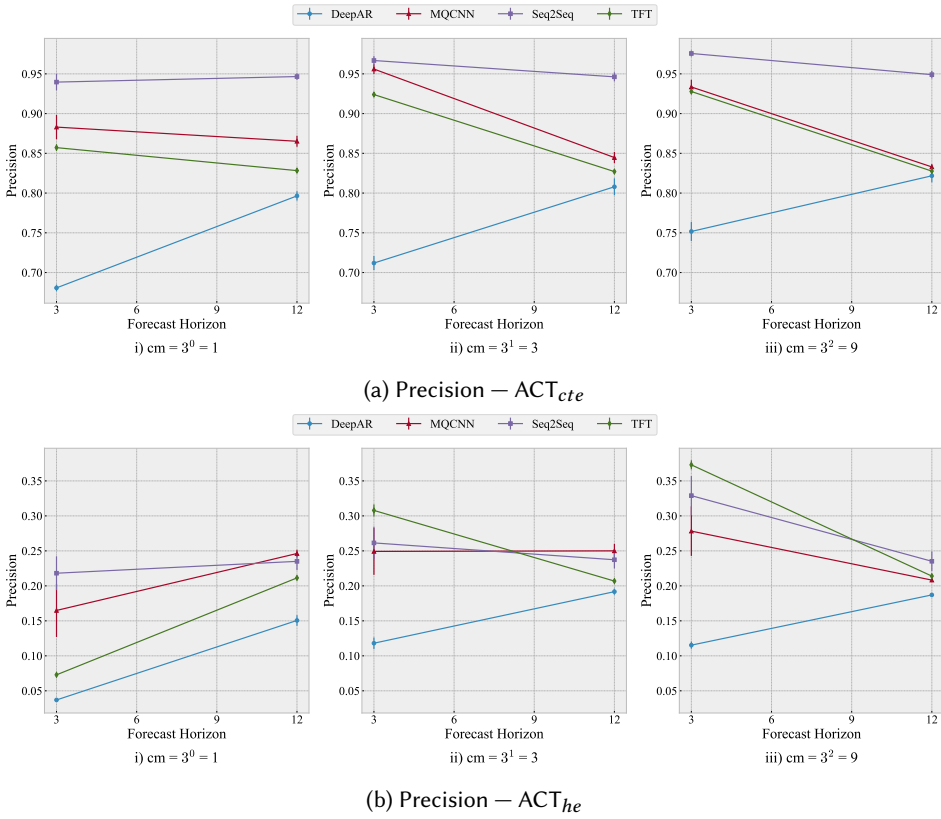


Fig. 6. Precision score measurements for DeepAR, MQCNN, Seq2Seq, and TFT, in various window configurations, for the *cte* (a) and *he* (b) safety requirements.

score of DeepAR for configurations with $h = 12$ at any cm value rises to become the best model. However, note that the highest F_3 score at $h = 12$ is still lower than the highest F_3 score reached by TFT at $h = 3$, suggesting that increasing the forecast horizon h leads to lower safety violation prediction accuracy in general.

Given a hazard forecast horizon of 3 s, for both *cte* and *he* safety requirements of the ACT case study, TFT yields the most accurate safety metric and safety violation predictions, with an improving accuracy when the lookback horizon length increases ($cm \times h$). For the *he* safety requirement at $cm = 1$, the lookback horizon does not contain sufficient information for any model to generate accurate forecasts. We further conclude that for prediction horizons longer than the minimum reaction time, i.e., $h = 12$ s), regardless of the cm value, DeepAR yields the most accurate predictions.

6.6 RQ₄: Latency and Memory Overhead

In section Section 6.6.1, we provide the details of our evaluation methodology for answering RQ₄. As discussed in the beginning of Section 6.5, we were only able to evaluate the effect of varying window sizes, on runtime performance, for the two safety requirements of the ACT case study (we

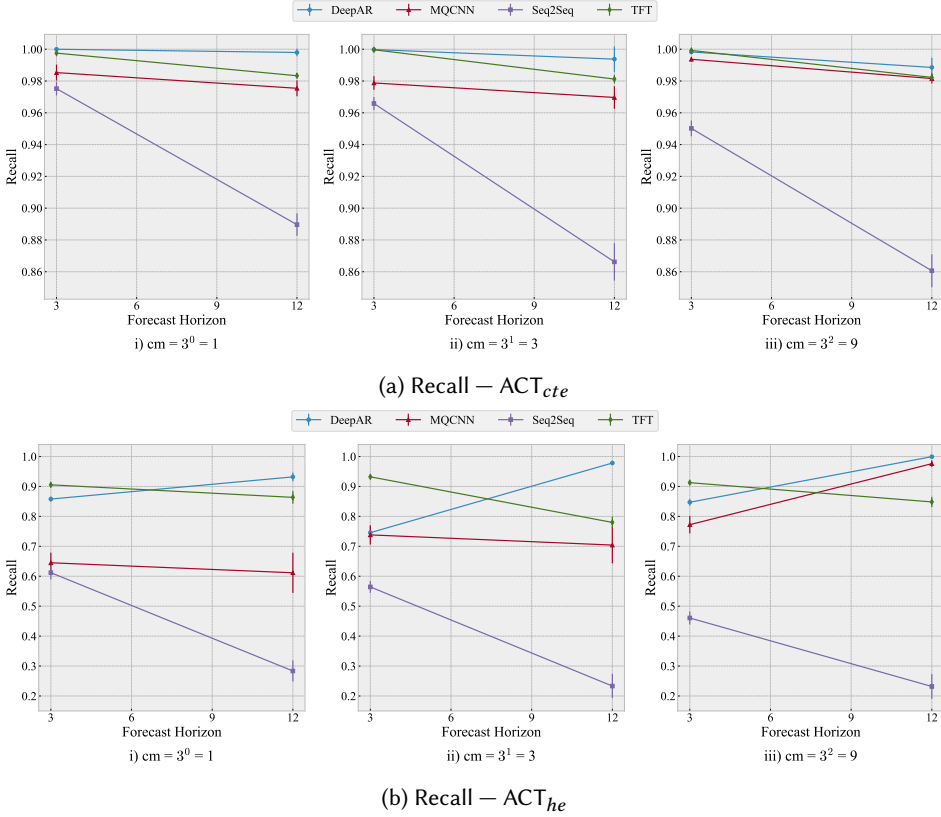


Fig. 7. Recall score measurements for DeepAR, MQCNN, Seq2Seq, and TFT, in various window configurations, for the *cte* (a) and *he* (b) safety requirements.

present the results in Section 6.6.2). We also measured the runtime performance of the ADS case study, when $h = 3$ s and $cm = 1$, which is also discussed in Section 6.6.2. We finally provide our answer to RQ₄ based on the results discussed in Section 6.6.2.

6.6.1 Methodology. To answer RQ₄, we queried each of the models trained in RQ₃ over the whole test set for both *cte* and *he* safety requirements of the ACT case study, and collected the average latency (in milliseconds *ms*) and peak GPU memory¹³ usage during inference (in megabytes, *MB*). Furthermore, we measure the model size in terms of GPU memory usage (in megabytes, *MB*), by measuring the difference in GPU memory usage before and after a model is loaded to the GPU. Note that the base ML libraries are preloaded in the GPU and their GPU memory usage is not reported. For the ADS case study, we applied the above process to the models that were trained in RQ₁ ($h = 3$ s and $cm = 1$).

6.6.2 Results. The RQ₄ evaluation results for the *cte* safety requirement of the ACT case study are similar to those for its *he* safety requirement. This is to be expected, as the size of the input and output vectors for the same model do not change between the two safety requirements. The only difference between them is due to differences in the hyperparameter values selected to best

¹³Recall that, as mentioned in Section 6.2, we use a GPU to train the models and generate predictions.

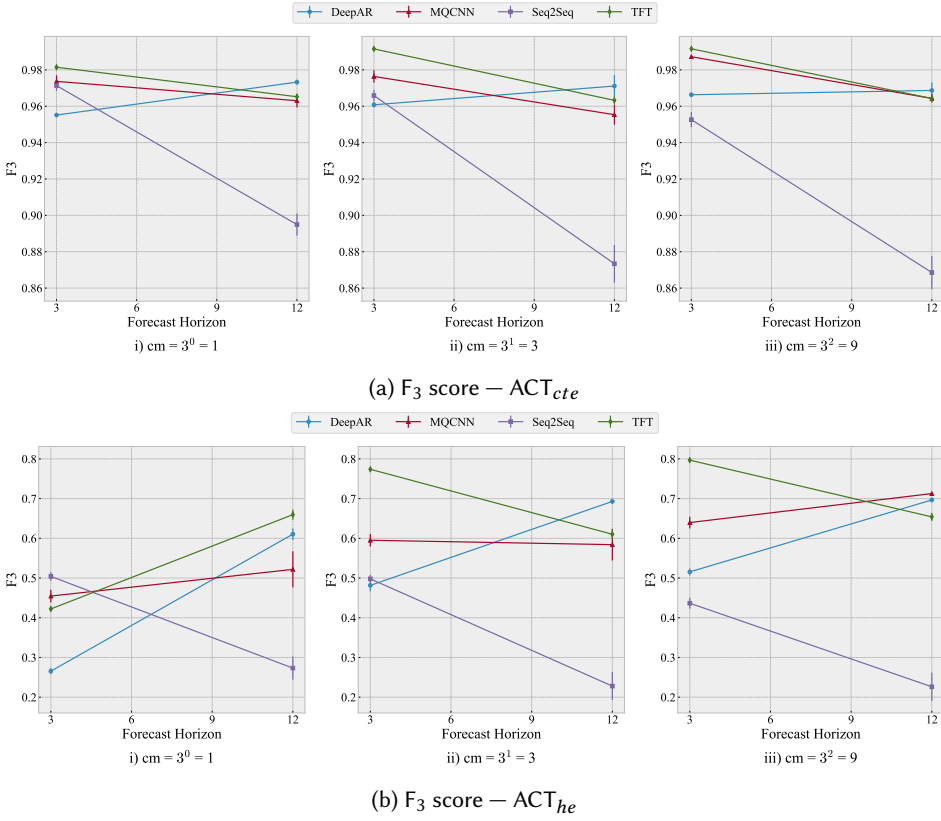


Fig. 8. F_3 score measurements for DeepAR, MQCNN, Seq2Seq, and TFT, in various window configurations, for the cte (a) and he (b) safety requirements.

address each safety requirement, which in turn can change the number of learnable parameters in the model. However, as observed, such change has not substantially changed the results. Therefore, although we present the results for the cte safety requirement of the ACT case study (Figure 9) in the paper, our discussion of the results and conclusions similarly hold for the he safety requirement. Moreover, we have observed that the ADS case study results, are not substantially different from the ACT case study results for the same window configuration, i.e., when $h = 3$ s and $cm = 1$. Thus, the discussion of the results and conclusions for the cte safety requirement of the ACT case study, given a similar window configuration, i.e., when $h = 3$ s and $cm = 1$, similarly holds for the ADS case study. The figures and results for the ADS case study, and the he safety requirement of the ACT case study, are available in our replication package (Section 6.9).

As shown in Figure 9a, all models in all window configurations have low GPU memory usage. TFT, which has the highest usage of all, only consumes around 175 MB. We further observe that an increase in context multiplier does not lead to a significant increase in model memory usage, while an increase in forecast horizon leads to slight increases in MQCNN and Seq2Seq model memory usage.

From Figure 9b, we observe that an increase in forecast horizon leads to increased peak memory usage during inference. When increasing cm , the peak memory usage of the models does not increase significantly for the same forecast horizon, except for TFT where the rate of increase in

peak memory usage with increasing forecast horizon grows at higher cm values. Nevertheless, the largest peak memory usage during inference, which we observe for TFT when $h = 12$ s and $cm = 9$, is 700 MB, only consuming 17.5% of the available memory of an NVIDIA Jetson Nano GPU, the *least* powerful embedded GPU made by NVIDIA. Therefore, we conclude that all the evaluated models, for all h and cm combinations, yield practical GPU memory usage in terms of model size and peak inference memory usage.

Finally, Figure 9c suggests that the average inference latency for sequence-to-sequence forecasting models (MQCNN, Seq2Seq, and TFT) does not significantly change when increasing the forecast horizon or context multiplier. However, as expected, DeepAR's inference latency linearly increases with forecast horizon. Furthermore, when increasing cm , DeepAR's prediction latency increases for the same forecast horizon. Thus, for longer prediction horizons or higher context multipliers, DeepAR is prone to having a high inference latency which could render its use prohibitive in the context of safety-critical systems. In general, for a safety monitor to be effective it should be able to predict safety violations and raise an alarm before the planning or control modules update their command. In our ACT example, the system does not include a planner and the controller directly generates control commands based on the perception system outputs. Thus, in our ACT case study, the use of a safety monitor is only meaningful if its average inference latency is less than the controller cycle time, i.e., the time period between two generated control commands. For practical reference, the design requirement for the maximum cycle time of planning and control modules at the Indy Autonomous Challenge [34], is 10 ms [39].¹⁴ In autonomous aviation, Paredes-Vallés et al. [64] and Navardi et al. [62] report an average vision cycle latency of 12 ms and 13 ms, respectively on an NVIDIA Jetson GPU. In another example, for a real-time vision-based drone system developed by Farrukh and West [22], the authors empirically measure and report that the average latency for the vision pipeline, where the safety monitor should have an average 16 ms or less to be useful. Note that DeepAR with a forecast horizon of 12, reaches the maximum latency of 10 ms at $cm = 3$ and largely exceeds it at $cm = 9$. In contrast, sequence-to-sequence models maintain a constant average inference latency of approximately 2 ms for all h and cm values. This is expected since sequence-to-sequence models generate their forecast for all forecast horizon timesteps at once.

For both *cte* and *he* safety requirements in the ACT case study, all models in all configurations yield practical model size and peak inference memory usage. Furthermore, although MQCNN, Seq2Seq, and TFT exhibit a constant and very low inference latency, DeepAR's inference latency significantly increases with the forecast horizon and context multiplier, which can render DeepAR impractical for longer forecast horizons at $cm > 1$.

For the ADS case study, all models, in a practical window configuration of $h = 3$ s and $cm = 1$, yield practical model size, peak inference memory usage and inference latency.

6.7 Discussion

Safety Monitoring via Safety Metric Forecasting. Overall, the results of our study suggest that safety metric forecasting, given learned component outputs and scenarios, is effective for safety monitoring. Indeed, the models, when evaluated on a dataset with a balanced distribution of safety violations, i.e., the *cte* safety requirement of the ACT case study, have yielded F_3 scores above

¹⁴The Indy Autonomous Challenge is an international challenge where teams from universities develop autonomous racing vehicles with the ultimate goal of improving the safety and performance of autonomous driving technology [34].

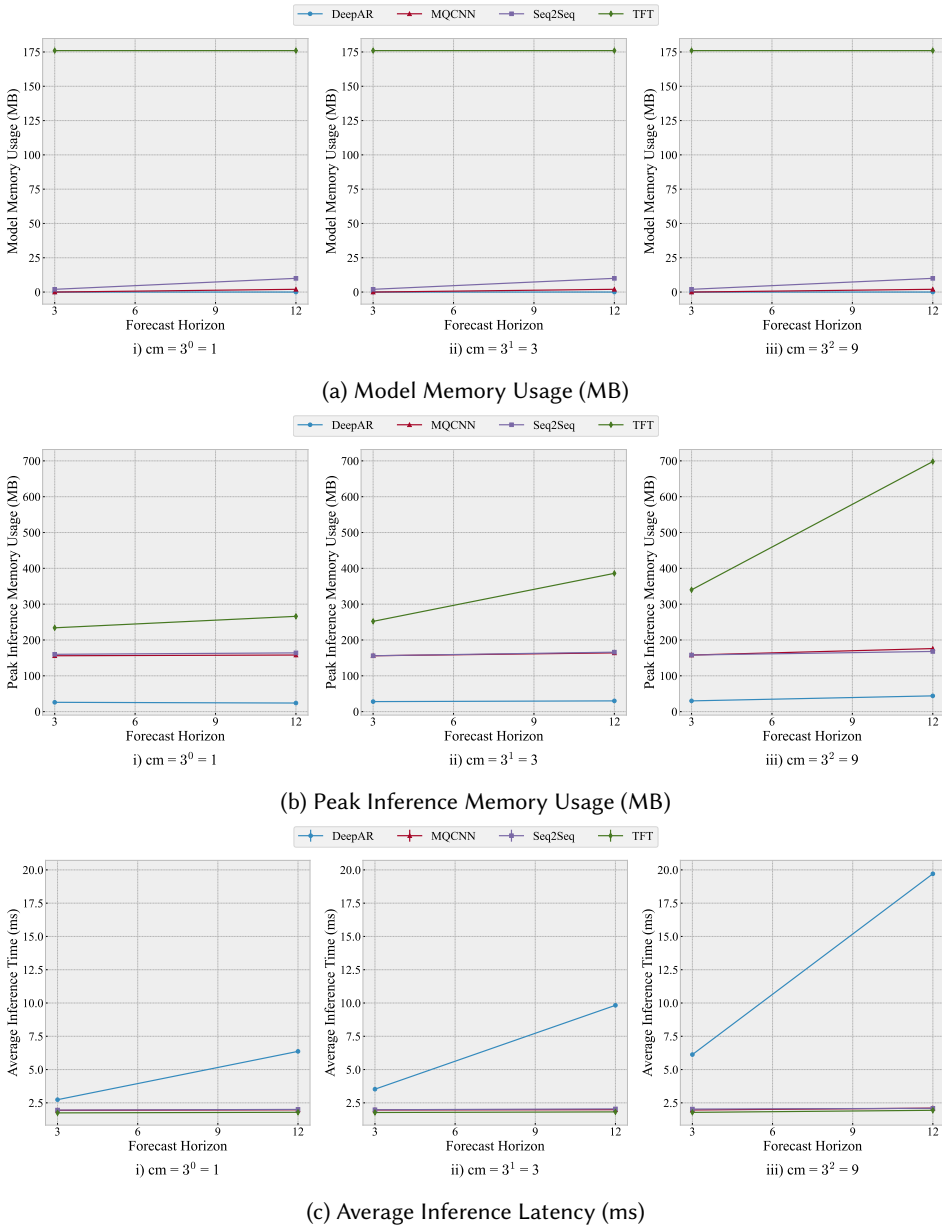


Fig. 9. The plot for a) model memory usage, b) peak inference memory usage, and c) average inference latency for DeepAR, MQCNN, Seq2Seq, and TFT models at different window configurations, for the *cte* safety requirement of the ACT case study, which is similar to the results for the *he* safety requirement. The above results at forecast horizon of 3 s and $cm = 1$ are also similar to the results for the ADS case study.

95% for all cm and h combinations, except for Seq2Seq at $h = 12$ s.¹⁵ However, the fact that all the evaluated models have high F_3 scores, for the *cte* safety requirement in the ACT case study, does *not* mean that all of them are equally accurate in predicting safety violations. This is highlighted in Table 4, where 1-2% differences in F_3 scores, for the *cte* safety requirement, translate into two orders of magnitude increase in the number of false negatives (FN). Furthermore, for both the *he* safety requirement of the ACT case study, where the distribution of safety violations in the dataset contains less safety violations compared to *cte*, and the ADS case study, where the size of the dataset is substantially smaller than that of the ACT case study, we observe that at least one model yields an F_3 score of 77% or more, for $q > 0.5$. This suggests that training the DL-forecasting models to achieve high safety violation prediction accuracy, collecting a large number of diverse safety violations in the dataset is crucial.

Moreover, we observe that the ranking of evaluated models in terms of safety metric accuracy (q-Risk) matches the ranking of models in terms of recall scores, suggesting that q-Risk values could be indicators of the recall scores for safety violation prediction (compare Figure 5 with Figure 7). However, we observe that the same order does not hold for precision and F_3 scores. For instance, MQCNN and Seq2Seq have higher (q-Risk) values than DeepAR for $h = 3$ s (Figure 5), whereas their precision and F_3 scores (Figure 6 and Figure 8, respectively) are significantly higher. Thus, we conclude that although q-Risk values, which are readily available after training the models and testing them on the test dataset, might be an indicator of the recall score for safety violation prediction, they are not indicators of precision and overall accuracy (F_3) scores for safety violation prediction. Thus, in practice, one needs to compute precision and F_3 scores before choosing a model for runtime deployment, and not only rely on q-Risk scores.

Our results further illustrate that the use of DL-based probabilistic forecasting methods, especially those with sequence-to-sequence architecture, leads to low inference latency while consuming feasible computing resources in terms of model size and peak memory usage during inference.

Furthermore, the results confirm the superiority of probabilistic forecasting over point forecasting for use in safety monitoring. This conclusion is drawn based on our empirical results and the fact that point forecast predictions correspond to the median ($q = 0.5$) value of the probability distribution predicted by probabilistic forecasting methods [11]. Our empirical results show that using values from the tail-end of the forecast probability distribution ($q \geq 0.95$ in our case) leads to more accurate safety metric and safety violation predictions than predictions for $q = 0.5$.

Window Configurations. We explored the effect of different combinations of varying hazard forecast horizons and context multipliers (*window configurations*), on prediction accuracy (RQ_3) and runtime performance (RQ_4), on the ACT case study only, due to the limited size of the dataset used for the ADS case study (as discussed in Section 6.3.1). Given all the results discussed in Section 6.5.2 and Section 6.6.2, we conclude that for both *cte* and *he* safety requirements in the ACT safety monitoring problem, TFT is the best model to be used for predicting imminent safety violations, i.e., $h = 3$ s, for all cm values. We further suggest that high cm values ($cm = 9$) be used as they improve overall safety violation accuracy. Nevertheless, as the peak inference memory of TFT increases with the increase in cm , a lower cm might also be considered depending on the available GPU memory onboard the learning-enabled autonomous system. The results for $h = 12$ s further highlight that, although DeepAR has superior prediction accuracy than other models, it is not as accurate as TFT for $h = 3$ s. Furthermore, the high average inference latency of DeepAR at $h = 12$ s prohibits it from being used as a safety monitor of the learned component, as discussed in Section 6.6.2. However, if

¹⁵Recall that by safety monitoring, we refer to runtime monitoring of learned components and the system operational context to predict a system safety requirement violation, which is, different from predicting when a learned component might mispredict.

model optimizations and specialized inference hardware can reduce the DeepAR’s inference latency to an acceptable range, it can be considered for predicting longer horizon safety violations given its good prediction accuracy. Nevertheless, considering both *accuracy* and *inference latency*, using TFT on shorter forecast horizons than 12 s, is a better option.

Challenging Scenarios. Although the trained safety metric forecasters yield high overall accuracy in predicting safety violations, it is important to identify the scenarios during which the safety monitor is more likely to mispredict safety violations. Characterizing such scenarios will allow the developers to generate more relevant execution data which can be used to train the safety monitors further and increase their safety violation prediction accuracy. Moreover, knowing the scenarios under which the safety monitor is expected to yield lower safety violation prediction accuracy would allow a system to be vigilant during the run-time of such scenarios and intervene in the automated operation of the system, if necessary.

One potential method relies on fitting a regression tree [51] to the safety violation prediction results (such as the ones provided in Section 6.4.2 and Section 6.4.3). Concretely, a regression tree is fitted to a dataset whose features are the scenario parameters, and target variable is the F_3 score that the safety monitor yields for the corresponding scenario. A notable benefit of a regression tree is that it allows the extraction of explainable rules that specify the part of the scenario space where the safety monitor yields a lower accuracy.

As an example, to explore the feasibility of explaining variation in safety violation prediction accuracy, we have fitted a regression tree to the safety violation prediction accuracy results of the safety monitor based on the TFT forecasting model, at the prediction quantile $q = 0.995$, for the *cte* safety requirement of the ACT case study (Section 6.4.2). Therefore, the features of the regression tree are the ACT scenario parameters, i.e., *time of day*, *cloud cover*, and *starting cte and he of the aircraft*¹⁶, while the target variable is the F_3 score of the TFT model at $q = 0.995$. Using grid search, we fitted regression trees by exploring combinations of values for its hyperparameters, i.e., *maximum depth* and *minimum number of samples per leaf node*. We computed the average *mean squared error* (MSE) for each model using 10-fold cross validation [23] and selected the most accurate tree, i.e., the one with the lowest average MSE over all ten folds ($MSE = 3 \times 10^{-4}$). The computed measure of determination (R^2) for the most accurate model is 0.68, indicating that most of the variance in F_3 is explained by the tree. We have provided the dataset used to train the regression tree, its preprocessing details, the model selection and cross validation script, as well as the selected regression tree (with detailed accuracy metrics) in our replication package (see Section 6.9).

Based on the most accurate regression tree, we observe that the following three rules characterize part of the scenario space where the safety monitor yields its lowest F_3 score:

- $\text{time_of_day} \in \{\text{afternoon}\} \wedge \text{cloud_cover} \in \{\text{moderate, high}\} \implies F_3 = 0.952$
- $\text{time_of_day} \in \{\text{afternoon}\} \wedge \text{cloud_cover} \in \{\text{none, low}\} \implies F_3 = 0.973$
- $\text{time_of_day} \in \{\text{morning}\} \wedge \text{he_start} \in [-10^\circ, -7.5^\circ] \implies F_3 = 0.984$

Given the above rules, we observe that time of day, cloud cover conditions and starting *he* values are the most important features explaining variations in safety violation prediction accuracy across scenarios. We further conclude, based on the above rules, that the lowest accuracy scores are observed during the afternoon, when the sky is moderately or highly cloudy ($\text{cloud_cover} \in \{\text{moderate, high}\}$). As mentioned earlier, knowing the scenario subspace where prediction is challenging, specified by the above rules, can help the developer understand where

¹⁶The detailed description and value ranges for the scenario parameters are provided in our replication package (Section 6.9).

more scenarios can be generated to re-train the safety monitor and potentially increase its prediction accuracy for low-accuracy scenarios. Moreover, the user can take the uncertainty of the safety monitor predictions into account at runtime, e.g., by being vigilant during scenarios where safety violation prediction accuracy is low and intervening when necessary.

6.8 Threats to Validity

In this section, we discuss potential threats to the validity of our study, namely internal, external, conclusion, and construct validity [79, 92, 100].

Internal Validity. Internal validity is concerned with the accuracy of the cause-and-effect relationships established by the experiments. Due to the limitations of the GluonTS library at the time of our evaluation, we had to use models for evaluation that were implemented in different ML frameworks. Concretely, DeepAR, MQCNN, and Seq2Seq models were implemented in MXNet while the TFT model was implemented in PyTorch, as the MXNet implementation of the TFT model was faulty and the MXNet models for MQCNN and Seq2Seq were not available. The use of different ML frameworks could impact the internal validity of the results. However, we conducted preliminary experiments on a model implemented in both frameworks¹⁷ to compare the impact of differences in ML framework on safety metric forecasting accuracy and runtime performance (memory and time overhead). We found that the differences in accuracy and runtime performance metrics were less than the standard deviation of the measurements and negligible.

External Validity. External validity is concerned with the generalizability of our results. One notable factor to consider is that in this study, we relied only on a specific ACT system (TinyTaxiNet) and simulation platform (X-Plane), for the ACT case study, and relied on a specific ADS system (Dave-2) and a simulation platform (Udacity simulator). However, X-Plane is a widely used high-fidelity simulator, and TinyTaxiNet was the best open-source ACT available at the time of our evaluation. Through preliminary experiments, we confirmed the superior *cte* estimation accuracy of TinyTaxiNet against two other pre-trained models that were available on the NASA ULI X-Plane Simulator project repository on GitHub [41]. Regarding the ADS case study, Dave-2 is a widely used lane keeping ADS, and Udacity is a popular simulator used for closed-track simulation of ADS. Nonetheless, further studies involving other learning-enabled autonomous systems in aviation and autonomous driving, as well as other domains, such as autonomous agriculture, and manufacturing, are required. We should however keep in mind that experiments such as the ones reported here entail substantial computations and extensive calendar time, i.e., 7500+ hours of GPU computation which were performed over 42 calendar days, thanks to having access to multiple GPUs on the Digital Research Alliance of Canada compute clusters. Another relevant factor is that due to X-Plane's capabilities, we were only able to set the weather statically, i.e., without sudden changes that rarely happen. The same static weather has been used in the ADS simulation used by our study [83]. Nevertheless, both X-Plane and Udacity simulator are widely used high-fidelity simulators, as mentioned above. Moreover, given that the maximum duration of a scenario execution, in both the ACT and ADS cases studies, are less than 4 min and 2 min, respectively, assuming a static weather over the execution is not unreasonable. An additional factor potentially impacting our proposed method's generalizability, is the fact that we assume that the monitored safety metrics are directly measurable (e.g., cte_{act} can be measured directly using a GPS) or can be estimated during the system operation (e.g., Time-To-Collision or TTC in the case of autonomous driving is estimated based the relative distance and velocity between the ADS and the object in front of it [59], which can be measured using a front radar on the ADS). However,

¹⁷At the time of our evaluation, only similar DeepAR implementations were available in both frameworks.

this is not a restrictive assumption as safety requirements, similar to any type of requirement, should have already been defined by the system developers and safety engineers, such that they are *measurable* [72]. Therefore, the safety requirements are expected to rely on metrics that can be measured or estimated to assess their satisfaction or violation. Another factor that could impact the generalizability of our results relates to the fact that we have not evaluated the performance results of our models on embedded hardware similar to the one that might be used during the operation of real ACT and ADS systems. Nevertheless, our analysis revealed that the memory demand by the models is quite low and in line with what is reported in resource-constrained environments. Regarding the average inference latency measurements, model optimizations such as model quantization [74], can potentially improve the average inference latency of models, especially DeepAR such that its latency falls below runtime constraints. However, we have performed our latency measurements using cloud-based NVIDIA V100 GPUs, in an inference setting, e.g., disabling gradient calculations [68, 69], where the inference latency is expected to be comparable to or lower than that of embedded GPUs.

Conclusion Validity. Conclusion validity relates to the conclusions that can be drawn from the collected data and their statistical significance. We followed the widely accepted rule-of-thumb of 30 repetitions for the experiments and we report every statistical value with its confidence interval.

Construct Validity. Construct validity is concerned with the degree to which the measured variables in the study represent the underlying concept being studied. As discussed in Section 6.3.1, q-Risk is a widely used metric to measure the accuracy of time series predictions (forecast safety metric values in our case) for a specific prediction quantile [11, 48, 49, 75], since it provides a summation of quantile loss (QL) over the forecast horizon for all predictions, normalized over all samples in the test set. As discussed in Section 6.4.1, we have used precision and recall which are widely used as metrics to capture the accuracy of the models in terms of missed safety violations and false prediction, respectively. Furthermore, similar to [83], we have used F_3 score as an aggregate metric to compare the overall safety violation prediction accuracy of safety monitors while capturing the relative importance of false negatives and false positives. As discussed in Section 6.6.1, the performance overhead introduced by the safety monitor can be refined to space and time overheads. Since models are loaded in the GPU, GPU memory usage is a metric that successfully captures the space overhead of the models as model size and peak memory usage during inference. Whereas, average prediction latency effectively captures the time overhead of the model at runtime.

6.9 Data Availability

The evaluated DL-based probabilistic forecasting models have been implemented in Python. We made the aforementioned implementations, the instructions to set up the ACT case study, the detailed description of the scenario parameters used to generate data, the generated ACT dataset, the raw and preprocessed ADS dataset, and the detailed evaluation results, for both the ACT and ADS case studies, will be made available online, once the paper is accepted.

7 CONCLUSION AND FUTURE WORK

In this paper, we proposed a method for safety monitoring of learned components in autonomous systems via probabilistic safety metric forecasting. We address the practical challenges of lacking access to internal information of the learned component and the system having limited operational resources, by using state-of-the-art DL-based probabilistic time series forecasters. They rely on scenarios and learned component output values to provide predictions of the safety metric probability distribution with acceptable inference latency and memory usage. We apply these forecasters to widely used case studies in autonomous aviation and autonomous driving, namely ACT and lane

keeping ADS, respectively, where we run extensive experiments to evaluate the safety metric and violation prediction accuracy, inference latency, and computation resource usage of state-of-the-art models, with a varying lookback and hazard forecast horizons while comparing them against a very competitive baseline (DeepAR). Our evaluation results suggest that probabilistic forecasting of safety metrics, given learned component outputs and scenarios, is effective for safety monitoring. Moreover, the evaluation results show that using Temporal Fusion Transformer (TFT) for predicting imminent safety violations ($h = 3$ s), for all lookback horizons, leads to the most accurate predictions with acceptable inference latency while consuming reasonable computational resources.

As part of future work, we plan to apply our proposed safety monitoring method to other learning-enabled autonomous systems in various domains such as automated driving, agriculture, and manufacturing. Furthermore, we plan to investigate whether using search-based methods that identify the hazard boundary of a learned component, e.g., MLCSHE [77], reduces the size of the dataset required to train an accurate safety monitor. Finally, in the future, we plan to further investigate the impact of re-training the safety monitor, using additional scenario generated according the regression tree analysis results (Section 6.7), on its safety violation prediction accuracy.

ACKNOWLEDGMENTS

We thank Corina Păsăreanu for her feedback in the early stages of the work and her pointer to the TinyTaxiNet model. We would also like to thank Nathan Aschbacher and Frederic Risacher for their constructive feedback on the work. This work was partially supported by funding from the Natural Sciences and Engineering Research Council of Canada (NSERC), through the Canada Research Chairs and discovery programs, the Science Foundation Ireland grant 13/RC/2094-2, Ontario Graduate Scholarship, Mitacs Accelerate Program, and Auxon Corporation. This research was enabled in part by support provided by British Columbia Digital Research Infrastructure (<https://www.bc.net>), Compute Ontario (<https://www.computeontario.ca>), and the Digital Research Alliance of Canada (<https://alliancecan.ca>). Andrea Stocco was supported by the Bavarian Ministry of Economic Affairs, Regional Development, and Energy.

REFERENCES

- [1] Airbus 2021. *A Statistical Analysis of Commercial Aviation Accidents 1958 – 2021*. Airbus. Retrieved October 01, 2024 from <https://skybrary.aero/sites/default/files/bookshelf/34487.pdf#:~:text=Fourth-generation%20commercial%20jet%20aircraft%20flew%2054%20of%20flights%20in%202021>.
- [2] Arden Albee, Steven Battel, Richard Brace, Garry Burdick, John Casani, Jeffrey Lavell, Charles Leising, Duncan MacPherson, Peter Burr, and Duane Dipprey. 2000. *JPL D-18709: Report on the Loss of the Mars Polar Lander and Deep Space 2 Missions*. Technical Report. Jet Propulsion Laboratory, Pasadena, CA.
- [3] Alexander Alexandrov, Konstantinos Benidis, Michael Bohlke-Schneider, Valentin Flunkert, Jan Gasthaus, Tim Januschowski, Danielle C. Maddix, Syama Rangapuram, David Salinas, Jasper Schulz, Lorenzo Stella, Ali Caner Türkmen, and Yuyang Wang. 2020. GluonTS: Probabilistic and Neural Time Series Modeling in Python. *Journal of Machine Learning Research* 21, 116 (2020), 1–6. <http://jmlr.org/papers/v21/19-820.html>
- [4] Hugo Araujo, Mohammad Reza Mousavi, and Mahsa Varshosaz. 2023. Testing, Validation, and Verification of Robotic and Autonomous Systems: A Systematic Review. *ACM Trans. Softw. Eng. Methodol.* 32, 2, Article 51 (March 2023), 61 pages. <https://doi.org/10.1145/3542945>
- [5] Erfan Asaadi, Steven Beland, Alexander Chen, Ewen Denney, Dragos Margineantu, Matthew Moser, Ganesh Pai, James Paunicka, Douglas Stuart, and Huafeng Yu. 2020. Assured Integration of Machine Learning-based Autonomy on Aviation Platforms. In *2020 AIAA/IEEE 39th Digital Avionics Systems Conference (DASC)*. IEEE, Institute of Electrical and Electronics Engineers (IEEE), San Antonio, TX, USA, 1–10.
- [6] Erfan Asaadi, Ewen Denney, Jonathan Menzies, Ganesh J. Pai, and Dimo Petroff. 2020. Dynamic Assurance Cases: A Pathway to Trusted Autonomy. *Computer* 53, 12 (2020), 35–46. <https://doi.org/10.1109/MC.2020.3022030>
- [7] Erfan Asaadi, Ewen Denney, and Ganesh Pai. 2019. Towards Quantification of Assurance for Learning-Enabled Components. In *2019 15th European Dependable Computing Conference (EDCC)*. IEEE, New York, NY, US, 55–62. <https://doi.org/10.1109/EDCC.2019.00021>

- [8] Erfan Asaadi, Ewen Denney, and Ganesh Pai. 2020. Quantifying Assurance in Learning-Enabled Systems. In *Computer Safety, Reliability, and Security*, António Casimiro, Frank Ortmeier, Friedemann Bitsch, and Pedro Ferreira (Eds.). Springer International Publishing, Cham, 270–286.
- [9] Ricardo Baeza-Yates and Berthier Ribeiro-Neto. 1999. *Modern Information Retrieval*. ACM, New York, NY, USA.
- [10] Tony Bellotti, Ilia Nouretdinov, Meng Yang, and Alexander Gammerman. 2014. Chapter 6 - Feature Selection. In *Conformal Prediction for Reliable Machine Learning*, Vineeth N. Balasubramanian, Shen-Shyang Ho, and Vladimir Vovk (Eds.). Morgan Kaufmann, Boston, 115–130. <https://doi.org/10.1016/B978-0-12-398537-8.00006-7>
- [11] Konstantinos Benidis, Syama Sundar Rangapuram, Valentin Flunkert, Yuyang Wang, Danielle Maddix, Caner Turkmen, Jan Gasthaus, Michael Bohlke-Schneider, David Salinas, Lorenzo Stella, François-Xavier Aubet, Laurent Callot, and Tim Januschowski. 2022. Deep Learning for Time Series Forecasting: Tutorial and Literature Survey. *ACM Comput. Surv.* 55, 6, Article 121 (dec 2022), 36 pages. <https://doi.org/10.1145/3533382>
- [12] Jennifer Black and Philip Koopman. 2009. System Safety as an Emergent Property in Composite Systems. In *2009 IEEE/IFIP International Conference on Dependable Systems & Networks*. IEEE, IEEE, New York, NY, USA, 369–378.
- [13] Daniel Bogdoll, Maximilian Nitsche, and J. Marius Zöllner. 2022. Anomaly Detection in Autonomous Driving: A Survey. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*. IEEE, New York, NY, USA, 4488–4499.
- [14] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. 2016. End to End Learning for Self-Driving Cars. arXiv:1604.07316 [cs.CV]
- [15] Markus Borg, Jens Henriksson, Kasper Socha, Olof Lennartsson, Elias Sonnsjö Lönegren, Thanh Bui, Piotr Tomaszewski, Sankar Raman Sathyamoorthy, Sebastian Brink, and Mahshid Helali Moghadam. 2023. Ergo, SMIRK is safe: A Safety Case for a Machine Learning Component in a Pedestrian Automatic Emergency Brake System. *Software Quality Journal* 31, 2 (2023), 335–403.
- [16] George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. 2015. *Time series analysis: Forecasting and Control*. John Wiley & Sons, Hoboken, NJ, USA. <https://doi.org/10.1002/9781118619193>
- [17] Cristian Challu, Kin G. Olivares, Boris N. Oreshkin, Federico Garza Ramirez, Max Mergenthaler Canseco, and Artur Dubrawski. 2023. NHITS: Neural Hierarchical Interpolation for Time Series Forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence* 37, 6 (Jun. 2023), 6989–6997. <https://doi.org/10.1609/aaai.v37i6.25854>
- [18] Darren Cofer, Isaac Amundson, Ramachandra Sattigeri, Arjun Passi, Christopher Boggs, Eric Smith, Limei Gilham, Taejoon Byun, and Sanjai Rayadurgam. 2020. Run-Time Assurance for Learning-Enabled Systems. In *NASA Formal Methods*, Ritchie Lee, Susmit Jha, Anastasia Mavridou, and Dimitra Giannakopoulou (Eds.). Springer International Publishing, Cham, 361–368.
- [19] On-Road Automated Driving (ORAD) Committee. 2021. *Taxonomy and Definitions for Terms Related to Driving Automation Systems for On-Road Motor Vehicles*. SAE International. https://doi.org/10.4271/J3016_202104
- [20] Nello Cristianini and John Shawe-Taylor. 2000. *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*. Cambridge University Press, Cambridge, UK.
- [21] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. arXiv:2010.11929 [cs.CV] <https://arxiv.org/abs/2010.11929>
- [22] Anam Farrukh and Richard West. 2023. FlyOS: rethinking integrated modular avionics for autonomous multicopters. *Real-Time Systems* 59, 2 (2023), 256–301.
- [23] Tadayoshi Fushiki. 2011. Estimation of Prediction Error by using K-fold Cross-Validation. *Statistics and Computing* 21 (2011), 137–146.
- [24] Yarin Gal and Zoubin Ghahramani. 2016. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In *Proceedings of The 33rd International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 48)*, Maria Florina Balcan and Kilian Q. Weinberger (Eds.). PMLR, New York, NY, USA, 1050–1059. <https://proceedings.mlr.press/v48/gal16.html>
- [25] Ruben Grewal, Paolo Tonella, and Andrea Stocco. 2024. Predicting Safety Misbehaviours in Autonomous Driving Systems using Uncertainty Quantification. In *Proceedings of 17th IEEE International Conference on Software Testing, Verification and Validation (ICST '24)*. IEEE, New York, NY, USA, 12 pages.
- [26] Stephen Haben, Marcus Voss, and William Holderbaum. 2023. *Time Series Forecasting: Core Concepts and Definitions*. Springer International Publishing, Cham, 55–66. https://doi.org/10.1007/978-3-031-27852-5_5
- [27] Franz Hell, Gereon Hinz, Feng Liu, Sakshi Goyal, Ke Pei, Tetiana Lytvynenko, Alois Knoll, and Chen Yiqiang. 2021. Monitoring Perception Reliability in Autonomous Driving: Distributional Shift Detection for Estimating the Impact of Input Data on Prediction Accuracy. In *Proceedings of the 5th ACM Computer Science in Cars Symposium (Ingolstadt, Germany) (CSCS '21)*. Association for Computing Machinery, New York, NY, USA, Article 8, 9 pages.

<https://doi.org/10.1145/3488904.3493382>

- [28] Dan Hendrycks and Kevin Gimpel. 2017. A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, online, 12 pages. <https://openreview.net/forum?id=Hkg4TI9xl>
- [29] Jens Henriksson, Christian Berger, Markus Borg, Lars Tornberg, Cristofer Englund, Sankar Raman Sathyamoorthy, and Stig Ursing. 2019. Towards Structured Evaluation of Deep Neural Network Supervisors. In *2019 IEEE International Conference On Artificial Intelligence Testing (AITest)*. IEEE, New York, NY, USA, 27–34. <https://doi.org/10.1109/aitest.2019.00-12>
- [30] Julia Hoffman and Dev Metha. 2023. *Artificial intelligence: in-depth market analysis market insights report*. Technical Report. Statista.
- [31] Xiaowei Huang, Daniel Kroening, Wenjie Ruan, James Sharp, Youcheng Sun, Emese Thamo, Min Wu, and Xinpeng Yi. 2020. A Survey of Safety and Trustworthiness of Deep Neural Networks: Verification, Testing, Adversarial Attack and Defence, and Interpretability. *Computer Science Review* 37 (2020), 100270.
- [32] Nargiz Humbatova, Gunel Jahangirova, and Paolo Tonella. 2021. DeepCrime: mutation testing of deep learning systems based on real faults. In *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis (Virtual, Denmark) (ISSTA 2021)*. Association for Computing Machinery, New York, NY, USA, 67–78. <https://doi.org/10.1145/3460319.3464825>
- [33] Manzoor Hussain, Nazakat Ali, and Jang-Eui Hong. 2022. DeepGuard: A Framework for Safeguarding Autonomous Driving Systems from Inconsistent Behaviour. *Automated Software Engineering* 29, 1 (2022), 1.
- [34] Indy Autonomous Challenge 2024. *Indy Autonomous Challenge*. Indy Autonomous Challenge. Retrieved March 22, 2024 from <https://www.indyautonomouschallenge.com/>
- [35] Anand Iyer and Aditya Prakash. 2019. *Controlling Biases*. John Wiley & Sons, Ltd, NY, NY USA, Chapter 10, 77–82. <https://doi.org/10.1002/9781119571278.ch10> arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/9781119571278.ch10>
- [36] Gunel Jahangirova, Andrea Stocco, and Paolo Tonella. 2021. Quality Metrics and Oracles for Autonomous Vehicles Testing. In *2021 14th IEEE Conference on Software Testing, Verification and Validation (ICST)*. IEEE, NY, NY USA, 194–204. <https://doi.org/10.1109/ICST49551.2021.00030>
- [37] Tim Januschowski, Jan Gasthaus, Yuyang Wang, David Salinas, Valentin Flunkert, Michael Bohlke-Schneider, and Laurent Callot. 2020. Criteria for Classifying Forecasting Methods. *International Journal of Forecasting* 36, 1 (2020), 167–177. <https://doi.org/10.1016/j.ijforecast.2019.05.008> M4 Competition.
- [38] Kyle D. Julian, Mykel J. Kochenderfer, and Michael P. Owen. 2018. Deep Neural Network Compression for Aircraft Collision Avoidance Systems. arXiv:1810.04240
- [39] Chanyoung Jung, Andrea Finazzi, Hyunki Seong, Daegyoo Lee, Seungwook Lee, Bosung Kim, Gyuri Gang, Seungil Han, and David Hyunuchul Shim. 2023. An Autonomous System for Head-to-Head Race: Design, Implementation and Analysis; Team KAIST at the Indy Autonomous Challenge. arXiv:2303.09463 [cs.RO]
- [40] Ismet Burak Kadron, Divya Gopinath, Corina S. Păsăreanu, and Huafeng Yu. 2022. Case Study: Analysis of Autonomous Center Line Tracking Neural Networks. In *Software Verification*, Roderick Bloem, Rayna Dimitrova, Chuchu Fan, and Natasha Sharygina (Eds.). Springer International Publishing, Cham, 104–121.
- [41] Sydney M. Katz, Anthony Corso, Sandeep Chinchali, Amine Elhafi, Apoorva Sharma, Mykel J. Kochenderfer, and Marco Pavone. 2021. *NASA ULI X-Plane Simulator*. Stanford ASL. Retrieved May 7, 2024 from https://github.com/StanfordASL/NASA_ULI_Xplane_Simulator
- [42] Sydney M Katz, Anthony L Corso, Christopher A Strong, and Mykel J Kochenderfer. 2022. Verification of image-based neural network controllers using generative models. *Journal of Aerospace Information Systems* 19, 9 (2022), 574–584.
- [43] Alex Kendall and Yarin Gal. 2017. What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision?. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., New York, NY, USA. https://proceedings.neurips.cc/paper_files/paper/2017/file/2650d6089a6d640c5e85b2b88265dc2b-Paper.pdf
- [44] Diederik P. Kingma and Jimmy Ba. 2017. Adam: A Method for Stochastic Optimization. arXiv:1412.6980 [cs.LG]
- [45] Stephan Kolassa. 2016. Sometimes It’s Better to Be Simple than Correct. *Foresight: The International Journal of Applied Forecasting* 40 (2016), 20 – 26. <https://search.ebscohost.com/login.aspx?direct=true&db=bth&AN=114335722&site=ehost-live>
- [46] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2017. Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., New York, NY, USA. https://proceedings.neurips.cc/paper_files/paper/2017/file/9ef2ed4b7fd2c810847ffa5fa85bce38-Paper.pdf
- [47] Nancy G. Leveson. 2012. *Engineering a Safer World*. The MIT Press, Boston, MA, USA. 608 pages. <https://doi.org/10.7551/mitpress/8179.001.0001>

- [48] Bryan Lim, Sercan Ö. Arik, Nicolas Loeff, and Tomas Pfister. 2021. Temporal Fusion Transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting* 37, 4 (2021), 1748–1764. <https://doi.org/10.1016/j.ijforecast.2021.03.012>
- [49] Bryan Lim and Stefan Zohren. 2021. Time-series Forecasting with Deep Learning: A Survey. *Philosophical Transactions of the Royal Society A* 379, 2194 (2021), 20200209.
- [50] Felipe Tomazelli Lima and Vinicius M.A. Souza. 2023. A Large Comparison of Normalization Methods on Time Series. *Big Data Research* 34 (2023), 100407. <https://doi.org/10.1016/j.bdr.2023.100407>
- [51] Wei-Yin Loh. 2011. Classification and regression trees. *WIREs Data Mining and Knowledge Discovery* 1, 1 (2011), 14–23. <https://doi.org/10.1002/widm.8> arXiv:<https://wires.onlinelibrary.wiley.com/doi/pdf/10.1002/widm.8>
- [52] Guannan Lou, Yao Deng, Xi Zheng, Mengshi Zhang, and Tianyi Zhang. 2022. Testing of autonomous driving systems: where are we and where should we go?. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (Singapore, Singapore) (ESEC/FSE 2022)*. Association for Computing Machinery, New York, NY, USA, 31–43. <https://doi.org/10.1145/3540250.3549111>
- [53] Yuan Luo, Ya Xiao, Long Cheng, Guojun Peng, and Danfeng (Daphne) Yao. 2021. Deep Learning-based Anomaly Detection in Cyber-physical Systems: Progress and Opportunities. *ACM Comput. Surv.* 54, 5, Article 106 (may 2021), 36 pages. <https://doi.org/10.1145/3453155>
- [54] R. J. Beckman M. D. Mckay and W. J. Conover. 2000. A Comparison of Three Methods for Selecting Values of Input Variables in the Analysis of Output From a Computer Code. *Technometrics* 42, 1 (2000), 55–61. <https://doi.org/10.1080/00401706.2000.10485979> arXiv:<https://www.tandfonline.com/doi/pdf/10.1080/00401706.2000.10485979>
- [55] David J. C. MacKay. 1992. A Practical Bayesian Framework for Backpropagation Networks. *Neural Computation* 4, 3 (05 1992), 448–472. <https://doi.org/10.1162/neco.1992.4.3.448> arXiv:<https://direct.mit.edu/neco/article-pdf/4/3/448/812348/neco.1992.4.3.448.pdf>
- [56] Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. 2020. The M4 Competition: 100,000 Time Series and 61 Forecasting Methods. *International Journal of Forecasting* 36, 1 (2020), 54–74. <https://doi.org/10.1016/j.ijforecast.2019.04.014> M4 Competition.
- [57] Spyros Makridakis, Evangelos Spiliotis, Assimakopoulos Vassilios, Artemios-Anargyros Semenoglou, Gary Mulder, and Konstantinos Nikolopoulos. 2023. Statistical, Machine Learning and Deep Learning Forecasting Methods: Comparisons and Ways Forward. *Journal of the Operational Research Society* 74, 3 (2023), 840–859. <https://doi.org/10.1080/01605682.2022.2118629> arXiv:<https://doi.org/10.1080/01605682.2022.2118629>
- [58] H. B. Mann and D. R. Whitney. 1947. On a Test of Whether one of Two Random Variables is Stochastically Larger than the Other. *The Annals of Mathematical Statistics* 18, 1 (1947), 50–60. <http://www.jstor.org/stable/2236101>
- [59] Michiel M. Minderhoud and Piet H.L. Bovy. 2001. Extended time-to-collision measures for road traffic safety assessment. *Accident Analysis & Prevention* 33, 1 (2001), 89–97. [https://doi.org/10.1016/S0001-4575\(00\)00019-1](https://doi.org/10.1016/S0001-4575(00)00019-1)
- [60] Sina Mohseni, Haotao Wang, Chaowei Xiao, Zhiding Yu, Zhangyang Wang, and Jay Yadawa. 2022. Taxonomy of Machine Learning Safety: A Survey and Primer. *ACM Comput. Surv.* 55, 8, Article 157 (dec 2022), 38 pages. <https://doi.org/10.1145/3551385>
- [61] National Transporations Safety Board 2024. *General Aviation Accident Dashboard: 2012-2021*. National Transporations Safety Board. Retrieved October 01, 2024 from <https://www.ntsb.gov/safety/data/Pages/GeneralAviationDashboard.aspx#:~:text=Data%20Spreadsheet%20General%20Aviation%20Accidents,%20Findings,%20and%20Safety%20Recommendations:%202012-2021>
- [62] Mozghan Navardi, Aidin Shiri, Edward Humes, Nicholas R. Waytowich, and Tinoosh Mohsenin. 2022. An Optimization Framework for Efficient Vision-Based Autonomous Drone Navigation. In *2022 IEEE 4th International Conference on Artificial Intelligence Circuits and Systems (AICAS)*. IEEE, NY, NY USA, 304–307. <https://doi.org/10.1109/AICAS54282.2022.9869975>
- [63] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D. Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. 2019. Can You Trust Your Model's Uncertainty? Evaluating Predictive Uncertainty under Dataset Shift. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc., New York, NY, USA. https://proceedings.neurips.cc/paper_files/paper/2019/file/8558cb408c1d76621371888657d2eb1d-Paper.pdf
- [64] F. Paredes-Vallés, J. J. Hagenaaers, J. Dupeyroux, S. Stroobants, Y. Xu, and G. C. H. E. de Croon. 2024. Fully neuromorphic vision and control for autonomous drone flight. *Science Robotics* 9, 90 (2024), eadi0591. <https://doi.org/10.1126/scirobotics.adi0591> arXiv:<https://www.science.org/doi/pdf/10.1126/scirobotics.adi0591>
- [65] Corina S. Păsăreanu, Ravi Mangal, Divya Gopinath, Sinem Getir Yaman, Calum Imrie, Radu Calinescu, and Huafeng Yu. 2023. Closed-Loop Analysis of Vision-Based Autonomous Systems: A Case Study. In *Computer Aided Verification*, Constantin Enea and Akash Lal (Eds.). Springer Nature Switzerland, Cham, 289–303.
- [66] Corina S. Păsăreanu, Ravi Mangal, Divya Gopinath, Sinem Getir Yaman, Calum Imrie, Radu Calinescu, and Huafeng Yu. 2023. Closed-Loop Analysis of Vision-Based Autonomous Systems: A Case Study. In *Computer Aided Verification*,

Constantin Enea and Akash Lal (Eds.). Springer Nature Switzerland, Cham, 289–303.

- [67] Marco Peixeiro. 2022. *Time Series Forecasting in Python*. Simon and Schuster, New York City, NY, USA.
- [68] predictMode 2024. *MXNet Predict Mode*. Retrieved September 18, 2024 from https://mxnet.apache.org/versions/1.6/api/python/docs/api/autograd/index.html#mxnet.autograd.predict_mode
- [69] Pyt-InferMode 2024. *PyTorch Inference Mode*. Retrieved September 17, 2024 from https://pytorch.org/docs/stable/generated/torch.autograd.grad_mode.inference_mode.html
- [70] Quazi Marufur Rahman, Peter Corke, and Feras Dayoub. 2021. Run-Time Monitoring of Machine Learning for Robotic Perception: A Survey of Emerging Trends. *IEEE Access* 9 (2021), 20067–20075. <https://doi.org/10.1109/ACCESS.2021.3055015>
- [71] Elizabeth M Renieris, David Kiron, and Steven Mills. 2023. Building robust RAI programs as third-party AI tools proliferate. *MIT Sloan Management Review* (2023).
- [72] 2018. ISO/IEC/IEEE International Standard - Systems and software engineering – Life cycle processes – Requirements engineering. *ISO/IEC/IEEE 29148:2018(E)* (2018), 1–104. <https://doi.org/10.1109/IEEESTD.2018.8559686>
- [73] Vincenzo Riccio, Gunel Jahangirova, Andrea Stocco, Nargiz Humbatova, Michael Weiss, and Paolo Tonella. 2020. Testing Machine Learning based Systems: A Systematic Mapping. *Empirical Software Engineering* 25 (2020), 5193–5254.
- [74] Wolfgang Roth, Günther Schindler, Bernhard Klein, Robert Peharz, Sebastian Tschiatzschek, Holger Fröning, Franz Pernkopf, and Zoubin Ghahramani. 2024. Resource-Efficient Neural Networks for Embedded Systems. *Journal of Machine Learning Research* 25, 50 (2024), 1–51. <http://jmlr.org/papers/v25/18-566.html>
- [75] David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. 2020. DeepAR: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting* 36, 3 (2020), 1181–1191. <https://doi.org/10.1016/j.ijforecast.2019.07.001>
- [76] S.M. Sanchez. 2005. Work smarter, not harder: guidelines for designing simulation experiments. In *Proceedings of the Winter Simulation Conference, 2005*. 14 pp.–. <https://doi.org/10.1109/WSC.2005.1574241>
- [77] Sepehr Sharifi, Donghwan Shin, Lionel C. Briand, and Nathan Aschbacher. 2023. Identifying the Hazard Boundary of ML-Enabled Autonomous Systems Using Cooperative Coevolutionary Search. *IEEE Transactions on Software Engineering* 49, 12 (2023), 5120–5138. <https://doi.org/10.1109/TSE.2023.3327575>
- [78] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. arXiv:1409.1556 [cs.CV]
- [79] Dag I. K. Sjøberg and Gunnar Rye Bergersen. 2023. Construct Validity in Software Engineering. *IEEE Transactions on Software Engineering* 49, 3 (2023), 1374–1396. <https://doi.org/10.1109/TSE.2022.3176725>
- [80] Mark A. Skoog, Loyd R. Hook, and Wes Ryan. 2020. Leveraging ASTM Industry Standard F3269-17 for Providing Safe Operations of a Highly Autonomous Aircraft. In *2020 IEEE Aerospace Conference*. Institute of Electrical and Electronics Engineers (IEEE), Big Sky, Montana, USA, 1–7. <https://doi.org/10.1109/AERO47225.2020.9172434>
- [81] Mohsen Soori, Behrooz Arezoo, and Roza Dastres. 2023. Artificial intelligence, machine learning and deep learning in advanced robotics, a review. *Cognitive Robotics* 3 (2023), 54–70. <https://doi.org/10.1016/j.cogr.2023.04.001>
- [82] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *The journal of machine learning research* 15, 1 (2014), 1929–1958.
- [83] Andrea Stocco, Paulo J. Nunes, Marcelo D’Amorim, and Paolo Tonella. 2023. ThirdEye: Attention Maps for Safe Autonomous Driving Systems. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering* (Rochester, MI, USA) (ASE ’22). Association for Computing Machinery, New York, NY, USA, Article 102, 12 pages. <https://doi.org/10.1145/3551349.3556968>
- [84] Andrea Stocco and Paolo Tonella. 2022. Confidence-driven Weighted Retraining for Predicting Safety-critical Failures in Autonomous Driving Systems. *Journal of Software: Evolution and Process* 34, 10 (2022), e2386. <https://doi.org/10.1002/smr.2386> arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/smr.2386>
- [85] Andrea Stocco, Michael Weiss, Marco Calzana, and Paolo Tonella. 2020. Misbehaviour Prediction for Autonomous Driving Systems. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering* (Seoul, South Korea) (ICSE ’20). Association for Computing Machinery, New York, NY, USA, 359–371. <https://doi.org/10.1145/3377811.3380353>
- [86] Leonard J. Tashman. 2000. Out-of-sample Tests of Forecasting Accuracy: An Analysis and Review. *International Journal of Forecasting* 16, 4 (2000), 437–450. [https://doi.org/10.1016/S0169-2070\(00\)00065-0](https://doi.org/10.1016/S0169-2070(00)00065-0) The M3- Competition.
- [87] The Boeing Company 2022. *Statistical Summary of Commercial Jet Airplane Accidents: Worldwide Operations | 1959 – 2022*. The Boeing Company. Retrieved October 01, 2024 from https://www.faa.gov/sites/faa.gov/files/2023-10/statsum_summary_2022.pdf#:~:text=In%20this%2054th%20edition%20of%20the%20Statistical%20Summary%20of%20Commercial
- [88] Udacity 2022. *Udacity’s Self-Driving Car Simulator*. Udacity. Retrieved September 24, 2024 from <https://github.com/udacity/self-driving-car-sim>

- [89] András Vargha and Harold D. Delaney. 2000. A Critique and Improvement of the CL Common Language Effect Size Statistics of McGraw and Wong. *Journal of Educational and Behavioral Statistics* 25, 2 (2000), 101–132. <https://doi.org/10.3102/10769986025002101> arXiv:<https://doi.org/10.3102/10769986025002101>
- [90] Huiyan Wang, Jingwei Xu, Chang Xu, Xiaoxing Ma, and Jian Lu. 2020. Dissector: Input Validation for Deep Learning Applications by Crossing-layer Dissection. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering* (Seoul, South Korea) (*ICSE '20*). Association for Computing Machinery, New York, NY, USA, 727–738. <https://doi.org/10.1145/3377811.3380379>
- [91] Ruofeng Wen, Kari Torkkola, Balakrishnan Narayanaswamy, and Dhruv Madeka. 2018. A Multi-Horizon Quantile Recurrent Forecaster. arXiv:1711.11053 [stat.ML]
- [92] Hyrum K. Wright, Miryung Kim, and Dewayne E. Perry. 2010. Validity Concerns in Software Engineering Research. In *Proceedings of the FSE/SDP Workshop on Future of Software Engineering Research* (Santa Fe, New Mexico, USA) (*FoSER '10*). Association for Computing Machinery, New York, NY, USA, 411–414. <https://doi.org/10.1145/1882362.1882446>
- [93] X-Plane Core Team. 2024. *X-Plane 11 Flight Simulator*. Laminar Research, Columbia, South Carolina. <https://www.x-plane.com/product/desktop/>
- [94] Yan Xiao, Ivan Beschastnikh, David S. Rosenblum, Changsheng Sun, Sebastian Elbaum, Yun Lin, and Jin Song Dong. 2021. Self-Checking Deep Neural Networks in Deployment. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, New York, NY, US, 372–384. <https://doi.org/10.1109/ICSE43902.2021.00044>
- [95] Jie M. Zhang, Mark Harman, Lei Ma, and Yang Liu. 2022. Machine Learning Testing: Survey, Landscapes and Horizons. *IEEE Transactions on Software Engineering* 48, 1 (2022), 1–36. <https://doi.org/10.1109/TSE.2019.2962027>
- [96] Mengshi Zhang, Yuqun Zhang, Lingming Zhang, Cong Liu, and Sarfraz Khurshid. 2018. 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* (Montpellier, France) (*ASE 2018*). ACM, New York, NY, USA, 132–142. <https://doi.org/10.1145/3238147.3238187>
- [97] Mengshi Zhang, Yuqun Zhang, Lingming Zhang, Cong Liu, and Sarfraz Khurshid. 2018. 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* (Montpellier, France) (*ASE '18*). Association for Computing Machinery, New York, NY, USA, 132–142. <https://doi.org/10.1145/3238147.3238187>
- [98] Qianqian Zhang, Haifeng Wang, Hongya Lu, Daehan Won, and Sang Won Yoon. 2018. Medical image synthesis with generative adversarial networks for tissue recognition. In *2018 IEEE International Conference on Healthcare Informatics (ICHI)*. IEEE, IEEE, New York, NY, US, 199–207.
- [99] Juan Zhao, QiPing Feng, Patrick Wu, Roxana A. Lupu, Russell A. Wilke, Quinn S. Wells, Joshua C. Denny, and Wei-Qi Wei. 2019. Learning from Longitudinal Data in Electronic Health Record and Genetic Data to Improve Cardiovascular Event Prediction. *Scientific Reports* 9, 1 (2019), 717. <https://doi.org/10.1038/s41598-018-36745-x>
- [100] Xin Zhou, Yuqin Jin, He Zhang, Shanshan Li, and Xin Huang. 2016. A Map of Threats to Validity of Systematic Literature Reviews in Software Engineering. In *2016 23rd Asia-Pacific Software Engineering Conference (APSEC)*. IEEE, New York, NY, USA, 153–160. <https://doi.org/10.1109/APSEC.2016.031>
- [101] Amirhossein Zolfagharian, Manel Abdellatif, and Lionel C. Briand. 2023. SMARLA: A Safety Monitoring Approach for Deep Reinforcement Learning Agents. arXiv:2308.02594