

# Learning Run-time Safety Monitors for Machine Learning components

Ozan Vardal<sup>1</sup>, Richard Hawkins<sup>1</sup>, Colin Paterson<sup>1</sup>, Chiara Picardi<sup>1</sup>, Daniel Omeiza<sup>2</sup>, Lars Kunze<sup>2</sup>, and Ibrahim Habli<sup>1</sup>

<sup>1</sup> Department of Computer Science, University of York, York, UK  
{ozan.vardal, richard.hawkins, colin.paterson}@york.ac.uk  
<sup>2</sup> Department of Engineering, University of Oxford, Oxford, UK

**Abstract.** For machine learning components used as part of autonomous systems (AS) in carrying out critical tasks it is crucial that assurance of the models can be maintained in the face of post-deployment changes (such as changes in the operating environment of the system). A critical part of this is to be able to monitor when the performance of the model at runtime (as a result of changes) poses a safety risk to the system. This is a particularly difficult challenge when ground truth is unavailable at runtime. In this paper we introduce a process for creating safety monitors for ML components through the use of degraded datasets and machine learning. The safety monitor that is created is deployed to the AS in parallel to the ML component to provide a prediction of the safety risk associated with the model output. We demonstrate the viability of our approach through some initial experiments using publicly available speed sign datasets.

## 1 Introduction

The use of machine learning (ML) in perception and understanding is essential for many autonomous systems (AS). Where such systems are used for safety related tasks it is critical that the safety of the ML components can be assured prior to deployment. Post deployment, we must be able to demonstrate that the system continues to operate safely throughout operation in complex and dynamic environments. In [18] we introduced transfer assurance, a process for assuring ML components used in AS. Transfer assurance is used when the ML component is required to be updated in response to changes in the AS or, crucially, in response to changes in the environment in which it operates. Figure 1 shows an overview of our transfer assurance process, which is split into three stages containing six activities. The first stage considers the initial development of the ML component for deployment into an operational AS and makes use of the AMLAS assurance process [8]. This stage results in the creation of an ML component along with its safety case and a set of appropriate ML safety monitors. These ML safety monitors are deployed on the AS along with the ML component and are used to identify changes in the system or environment which invalidate the safety case for the ML component. The second and the

third phases deal with analysing and responding to such changes to maintain acceptable safety of the AS. This paper focuses on the first stage of the transfer assurance process, in particular the creation of effective safety monitors which allow us to understand the impact of change on ML models at run-time (activities 2 and 3 in Figure 1).

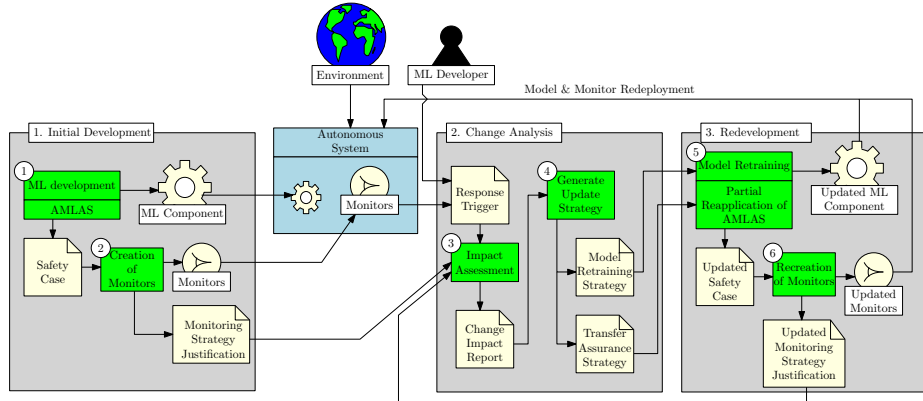


Fig. 1: Three Stage Transfer Assurance Process with activities shown in green and artefacts in yellow.

The challenge for activity 2 is to create monitors capable of detecting when ML component outputs are potentially unsafe due to changes in operating conditions. This is particularly challenging because we are unable to determine ground-truth during operation and because changes in operating conditions are combinatorial in nature, meaning that only by understanding the complex interplay of features can we determine the impact on safety.

We address this challenge by training an ML safety monitor with degraded data to reflect the potential impact of real-world factors. These degraded data sets are presented to the ML component and labelled to reflect their impact on model performance. This labelled data is then used as training data to learn a new model that can be used as a safety monitor during operation of the AS. This approach enables indication of the level of safety risk for the model given current input conditions.

This safety monitor is then deployed to the AS alongside the ML component, taking the same inputs from the operating environment at run-time and providing a prediction of the safety risk associated with the model output (activity 3), potentially triggering a response where the safety risk is assessed to be significant. The response may involve action from the system to ensure the continued safe operation of the AS (for example through fall-back to alternate systems or restricted operation) as well as consideration of the impact on the safety case.

The rest of the paper is organised as follows. In Section 1.1 we start by discussing related work before moving on to describe our proposed methodology

in Section 2. Next, in Section 3, we describe experiments we have undertaken to demonstrate our approach before finally, we discuss future work and provide conclusions in Sections 4 and 5.

### 1.1 Related work

This work fits in the broader research area of monitoring ML models during operation through the detection of *data distribution shift*, which refers to the phenomenon in supervised learning where data in operation change after deployment, differing from the training data used to develop the model [10]. This misalignment between the training and testing distribution can cause models to predict less accurately, despite maintaining confidence that their performance is good [1,19]. This is a significant safety issue, as erroneous predictions that are not appropriately flagged for inspection can lead to harmful outcomes. For example, changes in medical imaging software or hardware that are not anticipated by an ML component could endanger patients through over- or underdiagnosis [5,21].

Detecting samples that are out-of-distribution (OOD) is a problem that has received much research attention in recent years [7, 9, 14, 22]. Hendrycks et al. introduced a solution for detecting OOD samples in neural networks by thresholding softmax probabilities [9]. More recently, Liu et al. proposed use of energy scores, derived from energy-based models, as a more effective method than softmax confidence scores [14]. Another promising approach is use of non-parametric and parametric statistical tests for multiple conditional distribution hypothesis testing [13]. This approach can pinpoint the exact features causing data distribution shifts, addressing shortcomings of traditional OOD detection methods. Other approaches involve use of general adversarial networks (GANs) to detect anomalies [3, 4, 17] in high-dimensional data.

Effective safety monitoring may require more than identifying OOD inputs, as not all OOD data pose safety risks, and recent work has demonstrated that performant OOD detection does not necessarily translate to effective safety monitoring in practical scenarios [6]. Our work differs from OOD detection as we propose an approach to identifying the boundaries at which degradation of input data may make outputs from the ML component pose a hazard for the AS. Our primary contribution is the introduction of a process explicitly designed to create safety monitors for ML components used in complex environments. We demonstrate the effectiveness of our proposed process through application to an image classification problem for sign identification in an autonomous vehicle.

## 2 Methodology

In this section we describe the methodology we have developed for the creation and use of safety monitors that identify when changes to the post-deployment operational environment of the AS have the potential to invalidate the safety case for the ML component.

The seven step methodology, shown in Figure 2, starts with the identification of influencing factors that are likely to impact the performance the ML component and the safety of the AS. These influencing factors, e.g. blur and noise, arise from the combinatorial effects of natural phenomena in the operating environment and the operation of the sensors used in the AS. Once these have been identified, transformations must be specified to simulate the effects of these influencing factors on unperturbed input data (Step 1). For instance, approaches to such transformations could include mathematical functions that transform input data to produce perturbed output, or other data augmentation methods such as GANs [24] used to simulate input data perturbed by influencing factors in a realistic fashion.

The next step is to identify an unperturbed dataset that is representative of input data likely to occur in the operating environment of the AS, and systematically apply the defined transformations to this dataset at different levels of degradation to create a set of degraded datasets (Step 2).

Once a degraded dataset of each combination of influencing factors and perturbation levels has been created, the performance of the ML component on each degraded dataset is recorded (Step 3).

Based on these observations, each dataset is labelled according to its correspondence to a pre-determined performance threshold (Step 4). These thresholds are derived through an assessment of the impact of the performance of the ML component on the safety of the AS, informed by system safety analysis and aligned with levels of safety risk identified in the AS safety case. As an example, based on a consideration of required system behaviour, a dataset may belong to a “safety level 1: Normal operation” if the performance of the ML component on that dataset is above 90%, or ”safety level 3: Emergency stop” if the performance is below 20%.

Labelled datasets are then aggregated and split into training and development data in preparation for a typical machine learning training procedure (Step 5), whereby the safety monitor model is trained to classify samples as belonging to different levels of safety (Step 6). This classification allows the safety monitor to estimate the operational safety of the ML component used in an AS in real-world conditions by identifying the impact of natural phenomena in the operating environment on the images being presented to the ML models.

Finally in (Step 7), the performance of the safety monitor is evaluated using established validation methods. This can include testing on external datasets for unbiased assessment or using k-fold cross-validation techniques to ensure robustness and validity.

In order to illustrate our methodology, we demonstrate how it may be applied to an image classification task in an automotive context. We considered a safety monitor trained to work alongside a road sign classifier. The ML component that the safety monitor was trained for was a pre-trained convolutional model that had been developed for a previous study (see model 3b in [16]).

In the following we explain the rationale underlying each step of our methodology in more detail.

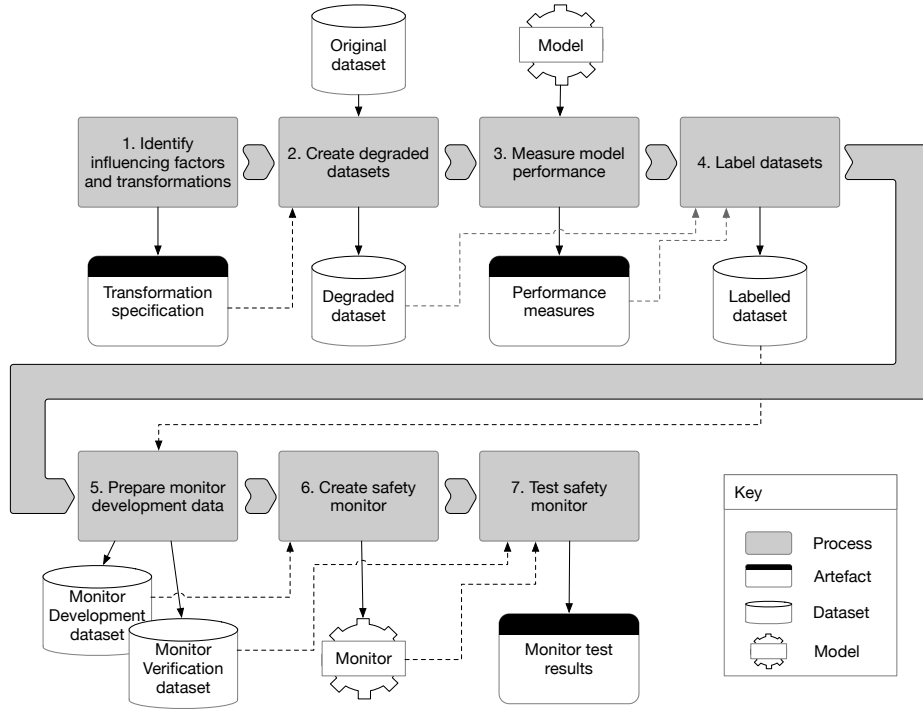


Fig. 2: Overview of process for creation of ML safety monitors

**Step 1: Identifying influencing factors and transformations:** The first stage of the safety monitor training process involves identification of *influencing factors* that affect the ML component during operation in ways that may compromise its safety. These factors can be identified through a consideration of internal factors, i.e. system features, and external factors, i.e. natural phenomena in the operating environment.

For our example application, the operating environment of the ML component is the road network on which the vehicle is operating, and the input data provided to the ML component are images of road signs captured by the vehicle’s camera sensors. We therefore identify two factors that may influence the performance of the ML component: weather conditions such as fog and heavy rain. These factors were chosen as we anticipate that elevated levels in some or all of these factors can lead to misidentification of road signs. This may affect the operating safety of the AS, for example, by resulting in incorrect speed adjustments which place the vehicle in a hazardous state.

Rather than attempting to model all possible real-world phenomena we consider the effect of these, in combination, on the image as presented to the ML model. For example rain, fog and particulates may all be present causing colour shifts, blur and noise in the image. It is not important for the safety monitor

to understand the source of these image perturbations, but rather the combined effect of these on the image and the potential impact on safety.

These transformations may take the form of mathematical functions which take an image sample as input and produce a degraded image that mirrors the effects of the identified influencing factors on the input data. Building on our previous work on test-based verification of neural network image classifiers [16], we define a perturbation encoding function  $g$  as transforming an image  $X$  into a set  $Z$  consisting of a single perturbed image  $X'$ , achieved by the application of a perturbation function on a per-pixel basis:

$$x'_{i,j} = \text{perturbation}(X_{i,j}, \epsilon), \quad (1)$$

where  $x'_{i,j}$  represents the pixel at coordinates  $(i, j)$  in the perturbed image  $X'$ , and  $X_{i,j}$  is the corresponding subset of pixels from the original image  $X$ , and where  $\epsilon \in [0, 1]$  serves as a measure of the intensity of perturbation. For images with color, the image  $X$  is encoded as an array of pixels, with each pixel containing a triplet of values that specify the color's red, green, and blue components.

Having defined the general mathematical form of the transformation in our image classification context, we then consider how it may be adapted for each influencing factor we anticipate to affect the ML component during operation. For example, to simulate the effects of haze in the image, due to fog or water spray, we can apply a *haze transformation* to each image, which may be defined as a color overlay that perturbs the image. Assuming a uniform haze effect, a haze color can be characterized by  $C_f = (r, g, b)$ . This effect is integrated into the image via the equation:

$$x'_{i,j} = (1 - \epsilon)x_{i,j} + \epsilon C^f, \quad (2)$$

A value of  $\epsilon = 0$  signifies no alteration to the image, whereas  $\epsilon = 1$  results in the image being completely filled with the color  $C^f$ . Operations on the pixel values are conducted element-wise.

To simulate blur, which may arise through rain or grease on the camera lens, we can perturb each image using a *blur transformation*. This is modelled using a convolutional filter with a Gaussian kernel, where each pixel  $x'_{i,j}$  in the blurred image is a weighted average of its neighbors:

$$x'_{i,j} = \sum_{k=-k_d}^{k_d} \sum_{l=-k_d}^{k_d} \alpha_{k,l} \cdot x_{i+k,j+l}, \quad (3)$$

with weights  $\alpha_{k,l} \in [0, 1]$  derived from a discretized Gaussian distribution, normalized such that  $\sum_{k=-k_d}^{k_d} \sum_{l=-k_d}^{k_d} \alpha_{k,l} = 1$ . The perturbation intensity  $\epsilon$  scales with the Gaussian standard deviation, dictating the blur level.

**Step 2: Creating degraded datasets:** In the next step, influencing factors and corresponding transformations identified in Step 1 are embedded into datasets used to train the safety monitor. This process has several substages:

1. An initial dataset that is maximally representative of an unperturbed operating environment must first be obtained. For an image classifier, this could be a dataset of images captured in the environment that the ML component is expected to operate in. In the case of our road sign classification model, we use a subset of original images from the German Traffic Sign Recognition Benchmark (GTSRB) dataset.

2. Transformations specified in Step 1 must be suitably encoded so as to perturb any given individual image sample at a given value  $\epsilon$ . Transformations may be applied uniformly across every point in a sample, or may be applied non-uniformly, for example, by applying a blur transformation to the center of an image, and a haze transformation to the edges. Indeed, the application methodology may be considered as an additional influencing factor, with the amount of image, or location, effected as a hyperparameter. For simplicity, in our example, we choose to apply each transformation uniformly to each pixel in the entire image.

3. Ranges of perturbations to apply to each sample must be identified. This specification can take the form of a vector  $E$  of intensity values  $\epsilon$ , where each  $\epsilon$  represents the intensity of a particular transformation to be applied to each sample. It is important to consider upper and lower limits that would be appropriate in light of the anticipated safety domains in the operating environment of the ML component. For example, if the ML component is expected to operate in a domain where fog is omnipresent, then the lower limit of the fog perturbation should be set to a value that is representative of the minimum amount of fog that the ML component is expected to encounter. Similarly, the upper limit of the fog perturbation should be set to a value that is representative of the maximum amount of fog that the ML component is expected to encounter.

4. We also consider the possibility that multiple individual sources of degradation may interact combinatorially to impact the image classification capability of the ML component and aim for coverage in the combined space. We therefore generate a set of all possible combinations of perturbation values that can be applied to each sample. Each combination of perturbation values represents a degraded dataset.

5. Once the above steps have been completed, each degraded dataset is created by sequentially applying each combination of transformations at perturbation values specified in the combinatorial product above to each sample in the initial dataset. For example, if we have two transformations, each with five possible perturbation values, and an initial dataset of 100 samples, then the creation of each degraded dataset would involve applying each of the 25 combinations of perturbation values to each of the 100 samples, resulting in a total of 25 degraded datasets, each comprising 100 samples with identical levels and types of perturbations within each dataset.

At the end of this process, a set of degraded datasets is generated that can represent the full spectrum of anticipated operating environments of the ML component, ranging from mild to severe degradation.

**Step 3: Measuring model performance:** The next step in our methodology involves quantitatively assessing how environmental changes impact the ML component’s performance. This assessment is critical for understanding and anticipating the conditions under which the ML component’s reliability may waver. Once the degraded datasets have been created, they are provided as input to the ML component that a safety monitor is being developed for, and the updated performance of the model on each dataset  $(p_1, p_2, \dots, p_m)$  is observed.

For our example we make use of the image classification model used in automotive scenarios taken from our earlier work [16]. We evaluate model performance across our range of contextually relevant degradations by using each degraded dataset as a test set for the original model and recording its accuracy, thereby establishing a performance baseline under various distortion conditions. Although we only record accuracy here, multiple performance metrics may also be considered.

**Step 4: Labelling datasets:** Our methodology assumes that there exists an acceptably safe state for the system which is related to the current model performance. Indeed this should be argued prior to deployment in the system safety case [8]. Distinguishing between acceptable and unacceptable operating states requires a specification of performance boundaries at which the system may move between states.

It is therefore necessary in our approach to first specify a vector of performance thresholds that define the boundaries of acceptable performance for the model, and then creating labels  $l_1, l_2, \dots, l_m$  for each degraded dataset to represent the expected operating state of the model based on its performance in a degraded environment. Thus,  $l_1$  would indicate very degraded accuracy and a hazardous operating state, whereas  $l_m$  would indicate almost the same accuracy as achieved on the original dataset and an acceptable operating state.

For our example case we assume three modes of operation exist. When the performance of the model (accuracy) is over 70% then the vehicle may operate normally. Between 70% and 40%, the vehicle would have a limited top speed giving the decision making component more time to gather data from the image sensor and hence mitigate the increased uncertainty in the input stream. Finally, when accuracy drops below 40% the vehicle is unable to continue to operate safely and a graceful cessation of autonomous function is initiated, this might be handing off to a human operative or bringing the vehicle to a halt. For such a case the thresholds are defined as [70, 40].

The datasets are then labelled according to these thresholds and the labeling forms the basis for the safety monitor’s training. In this way we ensure that the monitor is attuned to detect and respond to performance dips analogous to those it might encounter during the autonomous system’s operation.

**Step 5: Preparing monitor development data:** The creation of an effective safety monitor relies on the preparation of appropriate training data. This process adopted for safety monitor creation follows a typical machine learning



process [10,20]: labelled datasets are divided into development data and testing data. The development data are used for training and validation of the safety monitor. The test data is then used to evaluate the generalisability of the created safety monitor by checking that the safety monitor generates correct labels for the testing data.

In our application with the traffic sign recognition model, this translates into assembling a dataset that encompasses a range of degraded images, each labeled according to the previously established accuracy categories. In our example, we chose an 80:20 split of training to validation data, as this is common practice in most ML training contexts [12].

**Step 6: Creating the monitor:** The safety monitor is an ML model that is created by following the traditional ML cycle of training and testing [2]. For our case, the safety monitor for the traffic sign recognition model is a convolutional neural network trained to classify images based on the performance labels  $l_1, l_2, \dots, l_m$  specified in steps 4 and 5. Note that this structure was not extensively tuned and is purely to demonstrate the application of our methodology. The safety monitor is then trained to detect the severity of performance degradation in the model based on fluctuations in haze and blur.

**Step 7: Testing the monitor:** The final step involves evaluating the safety monitor’s performance on a set of test data similarly degraded. For our example case, we tested the resultant safety monitor using a 5-fold cross-validation procedure.

### 3 Evaluation

In this section we evaluate the use of our methodology and illustrate the application of each stage as described in the previous section.

We make use of a subset of the German Traffic sign benchmark [23] where each sample represents a speed sign as a 32x32 pixel RGB image. The ML component, for which we are designing a monitor, has a nominal accuracy, prior to image degradation, of 0.98 when assessed with its associated test set. This test set contains 4110 images and is described in Table 1 where the number of samples for each of the classes are reported.

In step 1 of the methodology we applied the blur and haze transformation, as previously described, to the original, unperturbed, dataset. We uniformly varied the intensity of these factors for each image. Initially, five different perturbation levels were applied for each factor:  $E = [0, 0.2, 0.5, 0.8, 1]$ . The number of degraded images may be calculated as

$$N = n_i \times \rho^{n_f} \tag{4}$$

where  $\rho$  is the number perturbation levels in the set  $E$ ,  $n_i$  is the number of initial images to which the degradation will be applied and  $n_f$  is the number of factors

class	0	1	2	3	4	5	6
description	30 km/h	50 km/h	60 km/h	70 km/h	80 km/h	100 km/h	120 km/h
# data samples	720	750	450	660	630	450	450

Table 1: Original test dataset with class descriptions and sample counts.

to be applied. For our case of  $n_i = 4110$ ,  $n_f = 2$  and  $\rho = 5$  we generate 102,750 degraded images.



Fig. 3: Perturbed samples from the GTSRB dataset showing the independent effects of the influencing factors. Numbers under the images indicate the level of the effect applied.

Next we inspected the generated images to assess the range of perturbation for each factor. This is to ensure that we are not generating unrealistic images. We then redefined our perturbation range such that the set of epsilon values for haze, and blur where  $E_{haze} = 0.8 * [0, 0.2, 0.5, 0.8, 1]$  and  $E_{blur} = 0.6 * [0, 0.2, 0.5, 0.8, 1]$ , respectively. Examples of unperturbed and degraded images, together with corresponding perturbation levels, are provided in Figure 3. While this is a simple demonstration of how synthetic transformations can be applied and tuned for a given context, we acknowledge that real-world examples are complex and likely warrant more a larger number of factors and associated mathematical transforms.

Each degraded dataset was then provided as input to the ML component and the accuracy of the model observed (Step 3). For the 25 degraded datasets

generated the model performed with accuracy  $\geq 70\%$  for 15 datasets, with  $70\% \geq \text{accuracy} \geq 40\%$  on 7 datasets, and with accuracy  $\geq 40\%$  on 3 datasets (see Figure 4). Labels were then added to each of the datasets (Step 4). For larger numbers of factors a visualisation of this type would not be possible, however, the approach would still identify regions in the factor space which are mapped to modes of safe operation.

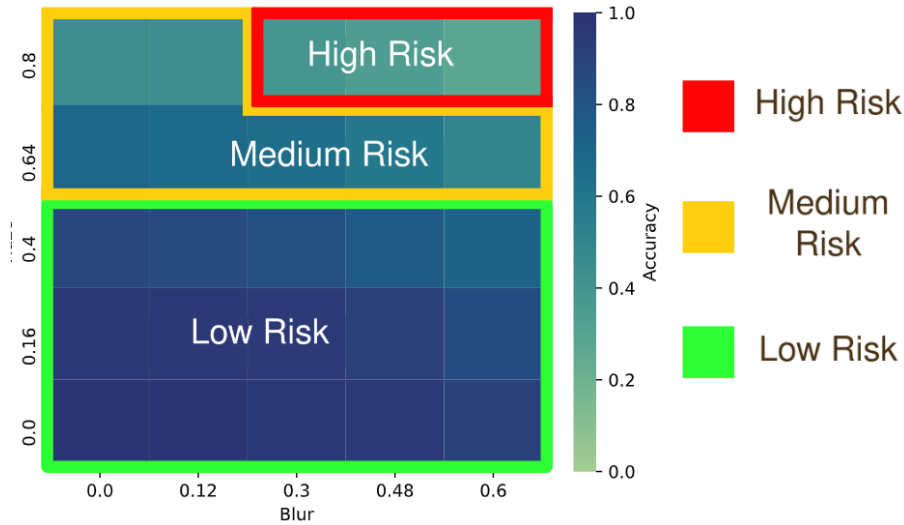


Fig. 4: Heatmap of accuracy of the original road sign classifier for each degraded dataset. Axes represent the  $\epsilon$  (degradation level) of each degradation effect. The region with the green border represents epsilon combinations that result in  $\text{accuracy} \geq 70\%$ , while those with the amber and red border represent epsilon combinations that result in  $40\% \leq \text{accuracy} < 70\%$  and  $\text{accuracy} < 40\%$  respectively.

Labelled degraded datasets were then prepared as training data for a safety monitor tasked with classifying degraded datasets based on their performance label (Step 5). Datasets for each class were concatenated and samples were randomly shuffled within each class. Each concatenated dataset was then split into training and testing sets with an 80:20 split. Once degraded datasets were shuffled, concatenated, and split into training and test sets, all training sets and test sets were concatenated. The prepared data were then saved together with labels for classes in the training (82200 samples) and test set (20550 samples).

In Step 6, we trained the safety monitor by initializing a convolutional neural network and training it to distinguish between degraded samples belonging to different classes. To ensure that the safety monitor’s performance was reliably evaluated we used 5-Fold cross-validation (this number of folds is standard prac-

tice [15]), splitting each data set into 5 subsets, with each subset serving as a test set in one of the folds and as part of the training set in the others.

Finally, we evaluated the performance of the trained safety monitor by providing samples from the test set as input and measuring performance. The final safety monitor model achieved an accuracy of 92%. Figure 5 provides a confusion matrix of the monitor as well as aggregate and individual ROC curves for each class.

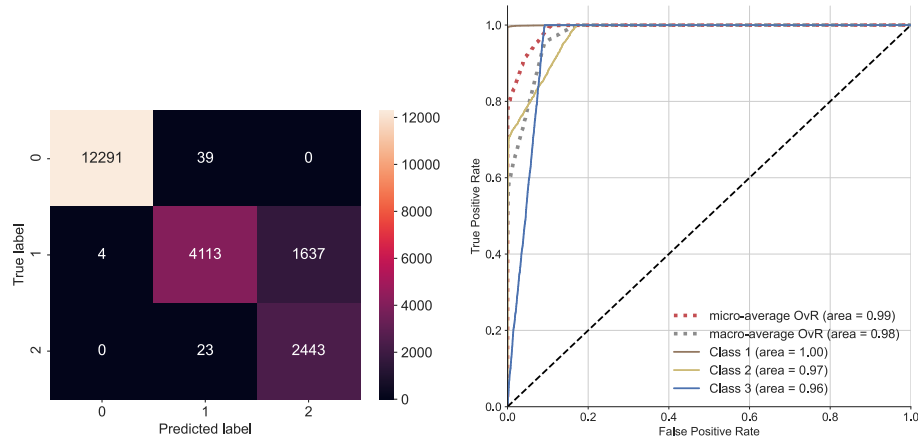


Fig. 5: The left panel shows the performance of the safety monitor across three classes, with each class representing the anticipated operational safety of the ML component in environments with different levels of perturbation. The right panel shows Receiver Operating Characteristic (ROC) curves for each of the three operational safety labels classified by the safety monitor.

## 4 Discussion

In order to illustrate our approach we have used a simple example based on a publicly available image dataset. While results of this case study demonstrate the feasibility of the approach, it is not limited in its application to a specific type of data or domain; our approach could in principle be used to create safety monitors for any ML component used in an AS. Further work will focus on testing the general applicability of the approach, with particular consideration for case studies involving non-image datasets. For example we are keen to test our approach to create safety monitors for ML components used in medical applications where the data takes the form of multi-attribute tabular patient data [11]. One of the key challenges here, compared to the case study presented

in this paper, will be the correct identification of relevant influencing factors and the different types of perturbations that will need to be considered.

The labelling of datasets used to train the safety monitor is crucial to the success of our process. For example if the risk regions defined for our use case in Figure 4 are defined differently, it could result in a higher number of high risk outputs being missed by the monitor, or conversely a higher number of false warnings being provided by the monitor. Although the former case would seem to be of most concern, a high number of false positives could result in operational issues for the AS (if the system takes action to reduce perceived increased operational risk), and ultimately in warnings from the safety monitor being ignored or the monitor being “turned off”. Further work will explore the effect of labeling decisions on the performance of the safety monitor and further guidance developed.

The examples we have used to-date have considered the inputs to the model as single-shot. In reality many of the applications of ML in AS will involve components presented with continuous data streams where a single erroneous output may be insufficient to affect the behaviour of the AS. In such systems it may be possible to create safety monitors that require multiple high-risk inputs before a warning is provided. Such a solution could also help to mitigate system operational issues discussed above.

## 5 Conclusions

In this paper we have introduced a method that begins to address the problem of how to monitor the safety of an ML model at runtime when ground truth is unavailable. Establishing a method to do this is crucial when the model is being used as part of an autonomous system to perform critical tasks. This forms part of a broader process for ensuring that assurance of the ML model is maintained in the face of post-deployment changes.

We have undertaken initial experiments to demonstrate the viability of our approach, with promising results. Further work is required to validate the approach on more realistic examples. In particular we will explore the nature of perturbations that are applied to training data to ensure that they are reflective of the real-world operational context in which the model is being used. This will, for example, explore how to best reflect the non-uniform nature of real-world perturbations in degraded datasets.

## 6 Acknowledgements

This work was supported by the Engineering and Physical Sciences Research Council through the RAILS project (EP/W011344/1) and the Assuring Autonomy International Programme, a partnership between Lloyd’s Register Foundation and the University of York.

## References

1. Amodei, D., Olah, C., Steinhardt, J., Christiano, P.F., Schulman, J., Mané, D.: Concrete problems in AI safety. *CoRR* (2016)
2. Ashmore, R., Calinescu, R., Paterson, C.: Assuring the machine learning lifecycle: Desiderata, methods, and challenges. *ACM Computing Surveys (CSUR)* **54**(5), 1–39 (2021)
3. Carrara, F., Amato, G., Brombin, L., Falchi, F., Gennaro, C.: Combining gans and autoencoders for efficient anomaly detection. In: 2020 25th International Conference on Pattern Recognition (ICPR). pp. 3939–3946 (2021)
4. Deecke, L., Vandermeulen, R., Ruff, L., Mandt, S., Kloft, M.: Image anomaly detection with generative adversarial networks (09 2018)
5. Finlayson, S.G., Subbaswamy, A., Singh, K., Bowers, J., Kupke, A., Zittrain, J., Kohane, I.S., Saria, S.: The clinician and dataset shift in artificial intelligence. *New England Journal of Medicine* **385**(3), 283–286 (2021)
6. Guerin, J., Delmas, K., Ferreira, R., Guiochet, J.: Out-of-distribution detection is not all you need. *Proceedings of the AAAI Conference on Artificial Intelligence* **37**(12), 14829–14837 (Jun 2023)
7. Hashemi, V., Křetínský, J., Rieder, S., Schmidt, J.: Runtime monitoring for out-of-distribution detection in object detection neural networks. In: Chechik, M., Katoen, J.P., Leucker, M. (eds.) *Formal Methods*. pp. 622–634 (2023)
8. Hawkins, R., et al.: Guidance on the assurance of machine learning in autonomous systems (AMLAS). *arXiv preprint arXiv:2102.01564* (2021)
9. Hendrycks, D., Gimpel, K.: A baseline for detecting misclassified and out-of-distribution examples in neural networks. *CoRR* (2017)
10. Huyen, C.: *Designing Machine Learning Systems* (2022)
11. Jia, Y., Kaul, C., Lawton, T., Murray-Smith, R., Habli, I.: Prediction of weaning from mechanical ventilation using convolutional neural networks. *Artificial Intelligence in Medicine* **117**, 102087 (2021)
12. Joseph, V.R.: Optimal ratio for data splitting. *Statistical Analysis and Data Mining: The ASA Data Science Journal* **15**(4), 531–538 (2022)
13. Kulinski, S., Bagchi, S., Inouye, D.I.: Feature shift detection: Localizing which features have shifted via conditional distribution tests. *CoRR* (2021)
14. Liu, W., Wang, X., Owens, J., Li, Y.: Energy-based out-of-distribution detection. In: *Advances in Neural Information Processing Systems*. vol. 33, pp. 21464–21475 (2020)
15. Nti, I., Nyarko-Boateng, O., Aning, J.: Performance of machine learning algorithms with different k values in k-fold cross-validation. *International Journal of Information Technology and Computer Science* **6**, 61–71 (12 2021)
16. Paterson, C., Wu, H., Grese, J., Calinescu, R., Păsăreanu, C.S., Barrett, C.: Deepcert: Verification of contextually relevant robustness for neural network image classifiers. In: *Computer Safety, Reliability, and Security: 40th International Conference, SAFECOMP 2021, York, UK, September 8–10, 2021, Proceedings* 40. pp. 3–17. Springer (2021)
17. Perera, P., Nallapati, R., Xiang, B.: Ocgan: One-class novelty detection using gans with constrained latent representations. pp. 2893–2901 (06 2019)
18. Picardi, C., Hawkins, R.D., Paterson, C., Habli, I.: Transfer assurance for machine learning in autonomous systems. In: *Proceedings of the Workshop on Artificial Intelligence Safety (SafeAI 2023)*. CEUR Workshop Proceedings (2023)

19. Rabanser, S., Günnemann, S., Lipton, Z.C.: Failing loudly: an empirical study of methods for detecting dataset shift (2019)
20. Raschka, S.: Model evaluation, model selection, and algorithm selection in machine learning. arXiv preprint (2018)
21. Roschewitz, M., Khara, G., Yearsley, J., Sharma, N., James, J.J., Ambrózay, É., Heroux, A., Kecskemethy, P., Rijken, T., Glocker, B.: Automatic correction of performance drift under acquisition shift in medical image classification. *Nature Communications* **14**(1), 6608 (Oct 2023)
22. Shafaei, A., Schmidt, M., Little, J.J.: Does your model know the digit 6 is not a cat? A less biased evaluation of "outlier" detectors. *CoRR* (2018)
23. Stallkamp, J., Schlipsing, M., Salmen, J., Igel, C.: The german traffic sign recognition benchmark: a multi-class classification competition. In: *The 2011 international joint conference on neural networks*. pp. 1453–1460. IEEE (2011)
24. Tanaka, F., Aranha, C.: Data augmentation using gans. arxiv (04 2019)