



This is a repository copy of *In-flight novelty detection with convolutional neural networks*.

White Rose Research Online URL for this paper:  
<https://eprints.whiterose.ac.uk/181593/>

Version: Submitted Version

---

**Article:**

Hartwell, A., Montana, F., Jacobs, W. et al. (3 more authors) (Submitted: 2021) In-flight novelty detection with convolutional neural networks. arXiv. (Submitted)

---

© 2021 The Authors. Preprint available under the terms of the CC-BY-NC-ND licence (<https://creativecommons.org/licenses/by-nc-nd/4.0/>).

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.



[eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk)  
<https://eprints.whiterose.ac.uk/>

# In-flight Novelty Detection with Convolutional Neural Networks

Adam Hartwell<sup>1</sup>, Felipe Montana<sup>1</sup>, Will Jacobs<sup>1</sup>, Visakan Kadirkamanathan<sup>1</sup>, Andrew R Mills<sup>1</sup>, and Tom Clark<sup>2</sup>

<sup>1</sup>Department of Automatic Control and Systems Engineering, University of Sheffield, UK  
<sup>2</sup>Rolls-Royce Plc, UK

December 8, 2021

## 1 Abstract

Gas turbine engines are complex machines that typically generate a vast amount of data, and require careful monitoring to allow for cost-effective preventative maintenance. In aerospace applications, returning all measured data to ground is prohibitively expensive, often causing useful, high value, data to be discarded. The ability to detect, prioritise, and return useful data in real-time is therefore vital. This paper proposes that system output measurements, described by a convolutional neural network model of normality, are prioritised in real-time for the attention of preventative maintenance decision makers.

Due to the complexity of gas turbine engine time-varying behaviours, deriving accurate physical models is difficult, and often leads to models with low prediction accuracy and incompatibility with real-time execution. Data-driven modelling is a desirable alternative producing high accuracy, asset specific models without the need for derivation from first principles.

We present a data-driven system for online detection and prioritisation of anomalous data. Biased data assessment deriving from novel operating conditions is avoided by uncertainty management integrated into the deep neural predictive model. Testing is performed on real and synthetic data, showing sensitivity to both real and synthetic faults. The system is capable of running in real-time on low-power embedded hardware and is currently in deployment on

the Rolls-Royce Pearl 15 engine flight trials.

## 2 Introduction

Aerospace Gas turbine engines (GTE)s are complex machines that must be monitored and maintained in order to ensure reliable operation over long time periods [1]. A key concept to maximise reliability is condition based, rather than schedule based, maintenance [2]. This requires an accurate assessment of engine condition, which can be challenging to obtain.

The condition of an engine is typically assessed relative to similar assets in a fleet, as well as functional performance thresholds, and requires expert engineering knowledge of the maintenance decision makers. The finite human resource to assess engine performance is traditionally supported by small data snapshots and failure-mode specific feature engineering. Larger data packets are collected from in-service on an ad-hoc basis, but this is logistically costly and often limited to cases where operational disruption has been observed. New methodologies are needed to support more efficient in-service operations.

Data, obtained from a large variety of different sensors placed on the engine, is the main resource available to assess engine condition. However, the large number of sensors, in-flight bandwidth limitations, and on-board storage limitations mean that it is not possible to return all data to ground for analysis in most current systems [3]. It is, therefore, necessary

to prioritise which data is returned and, in the context of condition based maintenance, prioritise data that is anomalous as this may indicate issues that require maintenance [4]. Furthermore, computational and data storage resources are extremely limited on the current generation of civil aerospace GTEs due to cost and weight restrictions and the harsh operating environment. Any data prioritisation system must function in real-time with limited computational resources lest data become unavailable before it can be saved. 46

Modern detection schemes make use of real-time engine models that run on-board in parallel to the real system, recently referred to as a Digital Twin [5, 6]. Model-based schemes have the advantage of not requiring fault data to validate their performance, and assessments of nominal model accuracy determine the detectability of an asset's condition divergence from its twin. Such models are traditionally derived from physical principals by domain experts. Physical models have advantages, particularly in terms of explicability and their ability to extrapolate (when designed correctly). However, they are costly to develop, difficult to tune to an individual asset, typically bespoke, and are computationally expensive to run [7, 8]. An alternative paradigm is data-driven modelling, where a flexible model structure is fitted to the system by training on historical data. Data-driven models are often poor in terms of explicability and their ability to extrapolate beyond the training set - but can be more accurate in prediction for individual assets with a smaller computational demands relative to a physical model.

There is, therefore, a need for an anomaly detection and prioritisation system that is capable of running on-board under strict computational resource limitations. To address this need we present a data-driven system that is capable of real-time anomaly detection and prioritisation on embedded hardware. The system can be used on subsets of the GTE system, or on the system as a whole, facilitating the prioritisation of high value data for condition monitoring. This system is currently in deployment on Pearl 15 engines.

The developed system is based on the Digital Twin paradigm, that is, accurate simulation of an individ-

ual 'as-operated' gas turbine. Simulation of the whole gas turbine is infeasible given the computational constraints of embedded devices, however by selecting interesting input-output subsets it is possible to generate many data-driven models to form a Digital Twin. Each model may then be run in parallel to give wide coverage of the gas turbine, allowing for generalised anomaly detection without excessive computational overhead.

Each model is trained using nominal data to make a prediction about a single signal given its input subset. The model also predicts a confidence on its prediction based on the model input. Model inputs that are not well described by the training data lead to a low confidence - thus allowing the distinction between model errors that arise from the unusual data of interest, and those from previously unseen operating conditions [9]. The system assumes that a sufficient amount of nominal data that covers expected operating conditions is available at training time.

The difference between model prediction and measurement, termed the residual, has been widely researched as a fault detection technique. The difficulty in obtaining an accurate model for the plant led to significant early work [10] on the development of diagnostic methods with plant observers (using data measurements to adapt state estimates), and with parity equation approaches decoupling known disturbances from faults into model subspaces [11, 12]. Such methods rely on analytical forms of the plant, faults and expected disturbances. The difficulties in suitably modeling gas turbine behaviour (e.g. its non-linear nature and inter-unit variance) have led to these approaches being augmented with data-driven 'tuning' elements [13, 14] at the cost of increased complexity. Machine learning offers a more generic, lower cost, and faster to market approach, as well as the potential to accurately capture modelled physics behaviours embedded in the data.

Other researchers have explored full or partial gas turbine anomaly detection using the full gamut of machine learning techniques [4]. Deep learning techniques have also been of particular relevance to the area in recent years [15, 16, 17, 18, 19, 20, 21] due to their success in other fields. These solutions are not suited to deployment on resource-limited hard-

ware due to their high computational cost at run-time and/or their memory requirements.

Anomaly detection using limited hardware has also been explored, using remote processing [22], look-up tables [23], model-based comparison [24], and in other heavily band-limited system [25]. These solutions, however, are not suited to gas turbine anomaly detection due to the complexity of modelling the highly non-linear asset behaviours.

Our system’s novelty is derived from its combination of a flexible modelling approach (allowing prediction of any input-output path within a gas turbine), a confidence prediction for quantifying uncertainty in the models, and software that allows for real-time operation on embedded hardware. Furthermore, our solution is especially flexible as it allows drop-in, and drop-out of models without retraining.

### 3 Methods

We use Chandola *et al’s* [26] definition of an anomaly: “Anomalies are patterns in data that do not conform to a well defined notion of normal behavior”

That is, we do not know how anomalies will present themselves ahead of time and therefore must instead define normal in order to determine when anomalies occur. We therefore take the data-driven Digital Twin approach of selecting a number of inputs and an output based on engineering knowledge, so that the error between the model prediction and true signal value may be evaluated to determine whether data is normal.

This viewpoint is also particularly useful in the gas turbine context because their high reliability means that the vast majority of available data is of nominal operation. A prediction approach allows a proxy measurement for detection in the form of prediction error. This is useful because it provides an indicator of model performance even when there is limited data with faults available.

Each model is designed to predict only a single output, multiple models with different input-output subsets are then required to run in parallel to provide wide coverage of GTE operation. This allows each model to be made lean and computationally efficient

while simultaneously providing a simple way to scale to the computational resources available.

The complexity of gas turbines coupled with the variety of operating conditions and environmental factors that affect their operation presents an additional problem for anomaly detection: data acquisition. It is difficult to acquire enough data to well define normal for all possible conditions, environments, and modes of operation. Therefore, we designed our models following the advice of Nix and Weigend [9], so that each model predicts a probability distribution characterised by a mean and confidence, see Section 3.2.1.

### 3.1 Anomaly Detection System

The architecture of the Anomaly Detection System is shown in Figure 1 which allows for scalable, real-time operation:

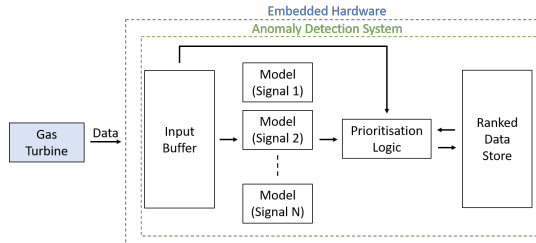


Figure 1: System architecture. Data from the gas turbine is piped in and buffered before being run through the currently loaded models. The most unusual data is then stored ready for later review.

The architecture allows streaming of data from sensors attached to a gas turbine in real-time into a local input buffer. The buffer is then used to feed  $N$  models, where  $N$  is selected based on a trade-off of coverage with computational resources. Data prioritisation is performed based on the Standardised Euclidean Distance metric described in Equation 10 (Section 3.2.4) and the most anomalous data stored in a local data store. The data stored can be any sensor information from the gas turbine and need not be limited to a model input-output subset.

Data in the store may be accessed whenever a data

link is available, for example, via satellite link when in-flight, or airport connectivity when landed. The store may be sized to match memory requirements of the hardware.

## 3.2 Data modelling

The anomaly detection scheme presented in the previous section is reliant on the existence of models that can map observations of the input data to observations of output data. Consider a single input-output data set  $\mathcal{D} = \{X, \mathbf{y}\}$ , where

$$\mathbf{y} = [y_1, y_2, \dots, y_N]^T \in \mathbb{R}^{N,1},$$

$$X =$$

where  $N$  is the length of the data series and  $m$  is the number of channels in the input data. The input-output mapping is approximated with a function  $f(\cdot)$ , such that the  $n$ 'th element of the output can be modelled as

$$y_n = f(\mathbf{x}_n) + e(\mathbf{x}_n), \quad (1)$$

where  $e(\cdot)$  is additive noise that is dependant on the input data.

The estimation of the distribution of the output data is central to the anomaly detection scheme, allowing the quantification of uncertainty on the estimate. Importantly, the uncertainty is a function of the input data  $\mathbf{x}_n$  to quantify the uncertainty in the model in different regions of the input space. The output distribution is described by a multivariate Gaussian with mean  $\mu_n = \mu(\mathbf{x}_n)$  and variance  $\sigma_n = \sigma(\mathbf{x}_n)$  given by

$$\mathcal{N}(y_n | \mu_n, \sigma_n^2) = \frac{1}{\sqrt{2\pi\sigma_n^2}} e^{-\frac{(y_n - \mu_n)^2}{2\sigma_n^2}}. \quad (2)$$

In this work, the input dependant mean and variance are estimated simultaneously using a deep neural network, the design, training, and real-time implementation of which are discussed in the remainder of this section.

### 3.2.1 Neural Network Design

The choice of a recurrent neural is often favoured when dealing with time-series data, however, their

high computational and memory costs make them generally unsuitable to implement on current embedded hardware [27]. We therefore chose to use a convolutional Neural Network architecture in order to minimise the computational cost of the forward passes.

The large number of sensors and complexity of the gas turbine lead to our decision to use a flexible network design that can be adapted to any target input-output data set. Further we also engineered the design to be flexible in terms of the number of time-steps the network can take as an input. Typically this requires a large number of extra weights to scale the number of time-steps. Our design uses a constrained convolutional layer (the Temporal Fold) to minimise the computational impact of inputting extra time-steps. This is especially important as, while some input-output subsets are expected to have short time-lag relationships, others e.g. temperatures, may have much longer time constants which a convolutional network will not be able to represent if the input data is unavailable.

The design is flexible in terms of the number of inputs, which allows for deployment on new data by non-specialists; since no redesign of the architecture is necessary to accommodate for different input-output subsets or timescales. The network architecture takes advantage of skip connections to improve performance at minimal computational cost, dropout regularisation to reduce over-fitting [28], and Leaky Rectified Linear Units (LReLU) throughout. The conceptual architecture is shown in Figure 2 which is described in more detail below;

**Input** Incoming data is treated as a 2D window with length  $T$  and width  $S$ , where  $T$  and  $S$  are time steps and the number of signals for a particular input-output subset respectively. This simplifies implementation and visualisation as it means the data can be treated in a similar manner to images.

**Temporal Fold** The stage of the network is a 1D convolution in the time dimension, explicitly enforcing feature extraction from individual signals at the lowest level, here named the Temporal Fold. Denoting the convolution operator as  $*$ , the discrete 1D convolution operation is defined as;

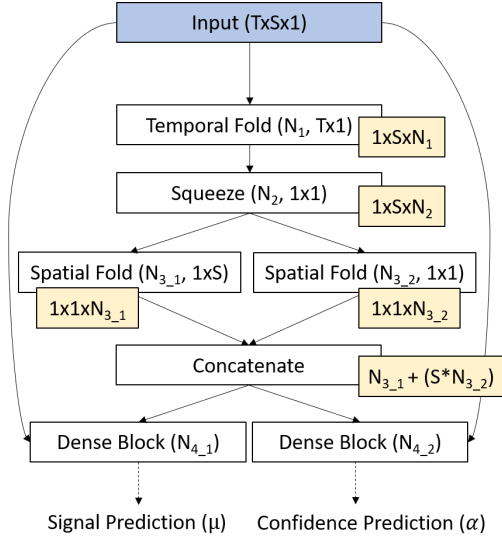


Figure 2: Overview of the network architecture. Number of hidden units and stride size given in brackets, output size given in yellow. The design is flexible such that “T” matches the number of incoming timesteps in each window and “S” is the number of signals. The “Dense Blocks” are 3 dense layers each followed by a dropout layer with  $rate = 0.5$ . The number of hidden units ( $N_X$ ) in different layers may be adjusted based on the complexity of the input-output set and available computation.

$$s_j^{[l]} = (\mathbf{z}_j^{[l-1]} * \mathbf{K}) = \sum_i \mathbf{z}_j^{[l-1]}(i) \mathbf{K}(i) \quad (3)$$

where  $\mathbf{z}_j^{[l-1]}$  is the output of the  $j$ 'th neuron at the  $l$ 'th layer of the network and  $\mathbf{K}$  is a 1D kernel of appropriate size. The temporal fold takes the network input as input  $\mathbf{z}_j^{[l-1]}(i) = \mathbf{x}_{k-T:k}$  with  $\mathbf{K} = \mathbf{K}_{TF}$  is the Temporal Folder Kernel with with height  $T$ , width 1, and depth equal to the number of input features. The output of the convolution can be seen as a weighted average over the past  $T$  data points for each input channel, where the weighting is defined by the elements of  $\mathbf{K}_{TF}$ .

The output of the convolution is passed through a

LReLU activation function defined as,

$$z_j^{[l]} = g(s_j^{[l]}) \quad (4)$$

$$g(s_j^{[l]}) = \begin{cases} s_j^{[l]} & \text{if } s_j^{[l]} > 0, \\ 0.01s_j^{[l]} & \text{otherwise} \end{cases} \quad (5)$$

**Squeeze** The next stage of the network is “Squeeze” layer. It consists of a  $1 \times 1$  Squeeze convolution which outputs to a  $1 \times S$  convolution and a  $1 \times 1$  convolution, both described by Equation 3 with appropriately sized kernels and with LReLU activation functions described by Equation 5. This method is inspired by Squeezenet [29], and replaces the  $c \times c$  convolution for a  $1 \times S$  convolution in order to significantly reduce the number of hyper-parameters to be learned.

**Spatial Fold** The two outputs of the Squeeze layer are passed through a 1D convolution described by Equation 3, this time in the spatial dimension (across features generated at the previous layer) and thus named the Spatial Fold. A LReLU activation function is again used. The layer determines cross-channel features and inter-filter features. The outputs of the convolutions are flattened and concatenated.

**Dense layers - Mean Prediction:** A skip connection is performed, where the output of the Spatial Fold is concatenated with the flattened original input. The concatenated outputs are then passed into a series of three dense layers described by;

$$z_j^{[l]} = g(Wz_j^{[l-1]} + \mathbf{b}) \quad (6)$$

where  $W$  and  $\mathbf{b}$  are hyper-parameters to be learned and  $g$  is the  $L - ReLU$  activation function defined by Equation 5. A dropout layer with  $rate = 0.5$  is placed after each dense layer to reduce over-fitting during training. The output is a single value that is an mean estimate of the output for the given network input.

**Dense layers - Confidence Prediction:** A further skip connection is performed where, again, the output of the Spatial Fold is concatenated with the flattened original input but here the mean estimate is also concatenated. This is then passed through into a series of three dense layers, as above for the mean

estimation, and outputs an estimate of the output variance.

### 3.2.2 Neural Network Optimisation

Optimisation of the model weights of a deep neural network requires an optimisation target. Considering the model described by Equation 2, it is simple to construct a likelihood function as

$$\mathcal{L}\mathcal{L} = \prod_{n=1}^B \mathcal{N}(y_n | \mu_n, \sigma_n^2), \quad (7)$$

where  $B$  is the batch size. Maximisation of the likelihood is equivalent to the minimisation of the negative log-likelihood

$$\begin{aligned} \mathcal{N}\mathcal{L}\mathcal{L} &= -\frac{1}{B} \sum_{n=1}^B \ln \left( \frac{1}{\sqrt{2\pi\sigma_n^2}} e^{-\frac{1}{2} \left( \frac{y_n - \mu_n}{\sigma_n} \right)^2} \right) \quad (8) \\ &= \frac{1}{B} \sum_{n=1}^B \left( \frac{1}{2} \left( \frac{y_n - \mu_n}{\sigma_n} \right)^2 - \ln(\sigma_n) - \frac{1}{2} \ln(2\pi) \right) \quad (9) \end{aligned}$$

which is used as the target for optimisation [9].

Optimisation was performed using the Rectified ADAM (RADAM) algorithm [30, 31] combined with early stopping on a separate validation data split. The parameters of the Rectified ADAM algorithm were varied depending on the input-output subset based on experiments with the validation set. RADAM was selected based on its performance relative to a small pool of candidate optimisers.

The network outputs the mean,  $\mu_n$ , and the confidence  $\alpha_n = \ln(\sigma_n)$ . The re-parameterisation enforces  $\sigma_n > 0$  to ensure that the network output is always valid.  $\sigma_n$  is regenerated after the forward pass using  $\sigma_n = e^{\alpha_n}$ . The assumption of normally distributed residuals is reasonable, given that we expect random fluctuations from sensor readings.

This confidence estimation approach is based on the work of Nix and Weigend [9] and similar to more recent work e.g. bounding box estimation in computer vision [32]). Other approaches to uncertainty estimation in neural networks include using dropout at run-time [33], using ensembles as a prediction scatter and Bayesian neural network solutions. We utilise

the explicit estimation approach since it allows us to enforce a Gaussian confidence estimation while retaining the greater flexibility of non-Bayesian neural networks and fast run-time of a single network solution.

### 3.2.3 Model and Hyper-parameter Selection

In order to ensure a neural network was an appropriate model choice, a number of simpler models were tested as a benchmark. These models include Lasso, Ridge, Elastic Net linear regression and Regression Forests. Each was found to perform worse than the neural network solution and naturally lacked the ability to estimate confidence. Typically it was possible to train these models to work well at a local scale, performing the same or better than the neural network over small portions of the data, but considerably worse elsewhere indicating a tendency to over-fit to the training data and a lack of generalisation.

Network hyper-parameters were selected based on starting values from the literature followed by jittered grid search and manual tuning. Validation set (not test set) data was used to evaluate performance and inform selection. More information on the hyper-parameters are included in Appendix A.

### 3.2.4 Model Comparison and Data Prioritisation

The system-level assessment of abnormality relies on being able to compare measurements to the outputs of any number of different models where each model may have a different output range. This is achieved by treating outputs of the network as a multivariate Gaussian and adopting the Mahalanobis distance [34]. Since each model is independent, such that the co-variance matrix is diagonal, the Mahalanobis distance reduces to Standardised Euclidean Distance (SED). We then refactor the equation for novelty ranking and drop the square root since it is unnecessary for calculating the relative ranking of each sample. This leaves Equation (10) that retains the intended functionality, but can be computed quickly in an online context,

$$\tilde{d}(y_n, \mu_n) = \frac{(y_n - \mu_n)^2}{\sigma_n^2}, \sigma = e^\alpha \quad (10)$$

This setup relies on models being up-to-date, that is, if an engine maintenance is performed then it is reasonable to expect some input-output relationships to change. Therefore, an out-of-date model may correctly indicate anomalous data for its input-output set, however, this is likely to be undesirable in practice. Methods for keeping models up-to-date are beyond the scope of this paper, and are discussed in our forthcoming work.

### 3.2.5 Real-time Implementation

Deployment to embedded hardware was achieved via an application written in C++. The Tensorflow Lite [35] library was employed to ensure fast forward-passes of all models involved while the application took care of the input and output buffering as well as data prioritisation tasks.

Design assurance for the application was achieved via a combination of automated code analysis (via CLang Format), automated testing, memory profiling (via Valgrind) and static analysis.

In order to evaluate performance, internal timers were used around the critical code sections in order to review data load and data evaluation times (model runs combined with data prioritisation). This allowed monitoring of real-world performance.

## 4 Data Sets

Two different data sets were used to validate that the system produced appropriate predictions and gave acceptable novelty detection performance.

Each data set was split into three sets: Training, Validation, and Test. Test data was used solely for evaluation and, unless otherwise specified, consisted of the most recent data. Training and Validation data was drawn from earlier data randomly. This split setup mimics real world applications where only past data is available and models must be deployed to an aircraft. In all cases, the Training set was used for determination of scaling and the only data directly

available to the networks for learning. The Validation set was used for early stopping and for hyperparameter tuning while the Test set was kept solely for evaluation.

**NASA TurboFan Degradation Dataset** [36]. This data comes from a prognostics competition aimed at assessing the quality of remaining useful life prediction algorithms. Specifically, subset “FD002” was used, that contains data on synthetically modelled degradation over six different operating conditions (combinations of altitude, air speed and thrust settings). The data simulates one measurement snapshot of each measured signal per flight to characterize an engine’s performance during that flight. For the selected data set, 519 (split between training and testing sets) different engines with different initial health conditions were simulated over increasing synthetic degradation, each engine is termed a simulation ‘run’. Our training data (marked for testing in the original) contained only low degradation and no failures, while the test set (marked for training in the original) had progressed to failure conditions near the end. Further, only the first half of each run was used during model training to help prevent the model learning the degradation.

**In-Service Fault Case Study.** This case study, provided by Rolls-Royce, includes continuous data from a twin engine aircraft across 116 flights at 1Hz. This data, in contrast to the NASA data set, contains engine dynamics and real world disturbances, but a relatively low level of degradation and only one engine. There is one known anomaly in this data (validated by domain experts) as well as several other potential anomalies. Only one engine from the aircraft is used for modelling so to represent the real aircraft avionics, which do not allow real-time access to signals from both engines by design. The data presented here has been normalised and anonymised to protect the interests of parties that supplied it.



## 5 Results and Discussion

### 5.1 NASA TurboFan Degradation

For each sensor signal in the data a model was trained using the three “operating condition” signals as input in order to ensure a model of nominal operation was learnt. Since the data consists of one measurement snapshot per flight, the window length was set to 1 i.e. only considering most recent data and without information on engine physical dynamics.

The model was validated with a Mean Absolute Error (MAE) of 0.029 +/- 0.034 (normalised range) across all sensors on all the engines in the validation data. The growth in plant-model mismatch, measured by Mean Standardised Euclidean Distance (MSED), is due to the increase in degradation towards the end of each run, as observed in Figure 3 A.

The high variance in the MSED signal (see Figure 3 A) is due to changes in the estimated confidence, which is similar across all signal models due to them using the same three operating condition inputs. The operating conditions in the MSED troughs of Figure 3 B are different to those found in the training data which generates a wide confidence interval (since the model is less likely to be a good predictor) and therefore returns a small NSED even if the prediction error is relatively high. Therefore, the model correctly indicates that its predictions may be less accurate in these regions.

In the test data, it is known that the Turbofan degrades to a failure state at the end of each run. It is therefore expected (and desirable) that the model will better represent the data at the start of the run and diverge from the plant behaviour at the end of each run. In order to test this property, the distribution of the MSED is compared for the 1st and 2nd half of each run contained in the test data, for each turbofan.

If the system is functioning correctly, the average MSED over the second half of each run would be greater indicating successful prioritisation of degraded (more unusual) data. The distributions were compared using a Mann-Whitney U test (since there is no guarantee of normality) under the null hypothesis that both sets came from the same distribution.

The arithmetic means of the distributions were then compared, a larger deviation from 0 indicating more unusual data.

The two sets of distributions were determined to have different shapes with a p value  $\approx 0$  ( $< 1e^{-308}$ ). The mean of the earlier portion of the runs was 0.980 while the latter portion was 3.466 thus demonstrating that the system correctly prioritises this data.

Not all runs are as clear-cut, however, as demonstrated by the Mann-Whitney U test, the trend of having higher distances near to faults is consistent.

Each model in this setup used 74,946 parameters making it suitable for deployment in real-time on embedded hardware (see Section 5.3).

These results illustrate that the modelling process and prioritisation methodology are applicable to a time series of static snapshot generated by a time-varying process<sup>1</sup>. The model flexibility allows the approach to be applied without structural changes to dynamic data, as shown in the next section.

### 5.2 In-Service Fault Case Study

This case study concerns the detection of a real fault that was detected in service and has been confirmed by other health monitoring systems and domain experts at our industrial partner. The fault was detectable in real-time using the current generation of health monitoring techniques, but no warning was given. The symptoms were noted as being most prominent in the P30 reading and related systems. The initial spike is most prominent, although not unique in the dataset, and after the spike, engine behaviour changes in an anomalous way for the rest of the flight, see Figure 4.

A model was trained using the scheme described in Section 3.2.1 for each of the four signals shown in Figure 4. The input to each model are the signals not predicted for, see Appendix B for more details. Each model has between 42,754, and 49,154 parameters. This is a tunable parameter, but was only affected by the number of input signals in this case.

The prediction performance was evaluated across the test set in terms of Mean Absolute Error (MAE)

---

<sup>1</sup>The code to reproduce these models and results on the NASA data-set is made available alongside this paper.

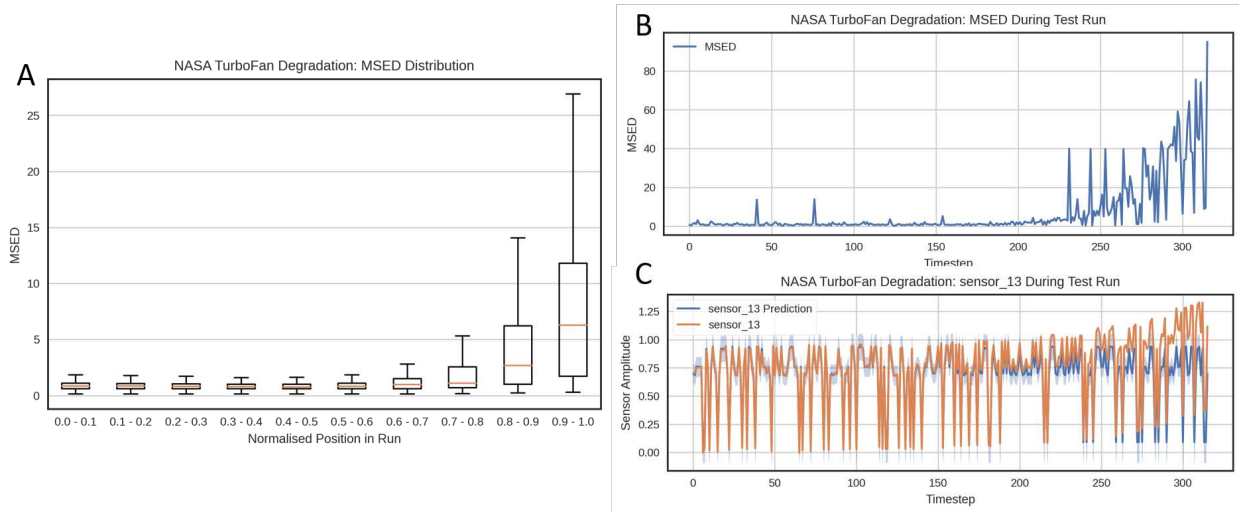


Figure 3: A: Distribution of MSEDs based on normalised position within test run demonstrating higher mean and spread towards the end of runs where faults occur. B: Example distribution of MSED over a test run demonstrating much higher values in the latter section where a fault occurs. C: Prediction for dominant (highest SED) model during the run, changes in shaded 95% confidence interval demonstrate cause of MSED oscillation is due to lack of confidence.

in the original signal range. This demonstrates that the models performed well in terms of pure predictive error:

- The N1 model achieved an MAE of 0.085 (% of max speed)
- The P30 model achieved an MAE of 0.715psi
- The TGT model achieved an MAE of 1.754°K
- The T30 model achieved an MAE of 0.872°K

Several interesting properties of the system are observed from the models' signal predictions at the time region around the known fault, see Figure 5. Firstly, the fault, and the ensuing anomalous behaviour, is clearly detected by the system. Secondly, the dominant SED (highest magnitude) for this fault is observed in the model predicting N1, despite the fault being most prominent in the P30 signal. This indicates that it is possible for the system to detect anomalous behaviour in signals not being directly predicted. Further, this indicates that the system

could monitor a wide range of signals with a small number of models if the set of inputs and outputs is well designed.

In contrast, the T30 model (which also uses P30 as an input) is significantly more accurate in prediction, and therefore rates the data as less anomalous. This indicates that during the fault, T30 still maintains the same relationship with P30, unlike the other signals. It also indicates that, while using a signal as a model input may allow detection of faults that manifest in that signal, it does not guarantee it.

Thirdly, the confidence intervals for all the models remain self-similar, with the exception of the P30 and TGT model during the initial spike. This indicates that the training data covers a sufficient range of operating conditions such that data in this flight is expected to be accurately predicted. The large confidence interval around the initial spike in these models demonstrates that this is the only section (from those tested) that includes data outside previously seen ranges. A large confidence interval may be an indication that something unusual is occurring, how-

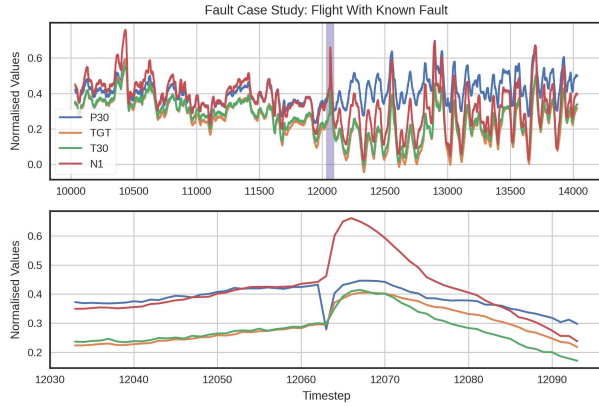


Figure 4: Normalised signals for the flight with the confirmed fault. The highlighted region in the top graph is zoomed in on in bottom graph. The spike in the internal pressure signal (P30) in the bottom graph is a known fault as is all behaviour afterwards. The temperature T30 is an internal temperature detected in the same locale as P30, N1 is the low speed shaft, and TGT is the Turbine Gas Temperature.

ever, since these regions do not appear in training data, it is not reasonable for the system to identify them as highly anomalous unless the error is also very high.

The fault shown in Figures 4 and 5 is the only fault in the data that has been previously detected, however, the system was able to identify multiple potential precursor events that show the same signature as the known fault, but without the initial spike. These were detected in the flight prior, 3 flights prior, and 6 flights prior, to the known fault. These events were not detected by current health monitoring systems, however, comparison to the other (assumed correctly functioning) engine on the aircraft, they present the same fault signature [37]. This indicates that they are highly likely to be related. The system’s ability to detect them lends further credence to its effectiveness. An example from the flight 3 flights before the known fault is shown in Figure 6.

The MSED for the flight containing the known fault (flight 113), as well as the following flights, are observed to be distributed differently from the normal

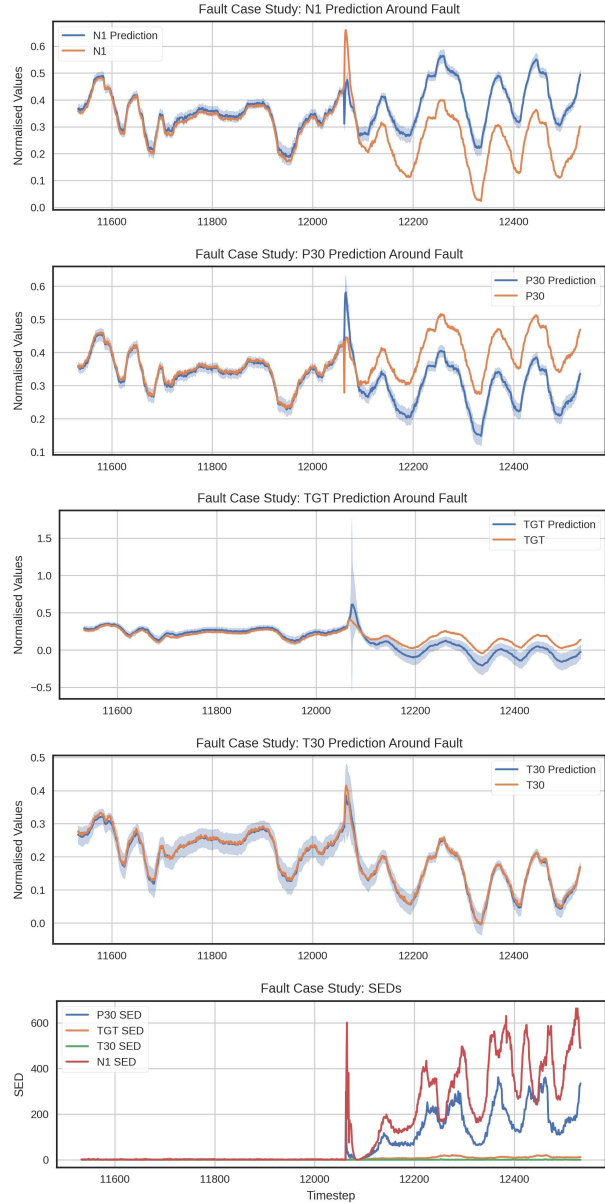


Figure 5: Top 4 graphs: Prediction around the fault, in descending order of average Standardised Euclidean Distance (SED). Highlighted regions indicates a 95% confidence interval (predicted). Bottom graph: SED of each model around the fault.

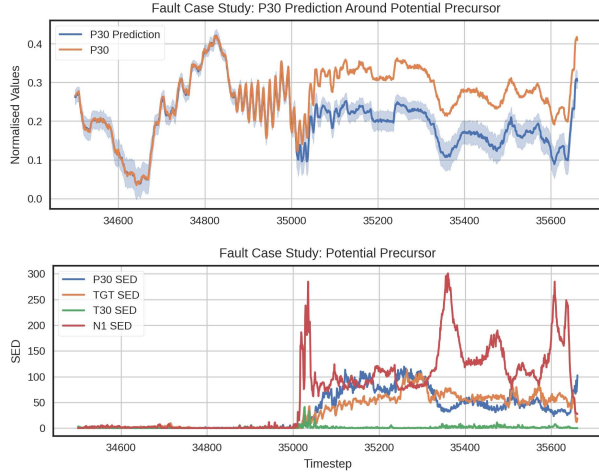


Figure 6: Potential precursor fault, 3 flights before known fault.

flights, see Figure 7. The precursor faults also show a similar shaped distribution with a second peak at  $MSED > 100$ . This indicates that the system is correctly prioritising unusual data on real flight data.

Flight 111 also demonstrates an unusual distribution; this is due to unusual behaviour in the shaft speed signal which is not found in other flights. It is unclear whether this is another fault or an unusual but harmless behaviour. This highlights the benefit of data prioritisation, over binary labelling since an operator can review the data to make the final determination on this sort of edge case.

### 5.3 Real-time Performance

The system was deployed on a Xilinx Zynq XC7Z020-2CLG484I, an industrial grade device based on an ARM Cortex A7 microprocessor, and on-wing system target. In order to test real-time operation, the models designed above were deployed to the device and constrained to  $\leq 20\%$  of CPU time. A synthetic data set of 25 signals over 100,000 time points was sent to the device and an output buffer of size 50 windows was used so that the 50 most unusual windows would be returned per ground data transmission.

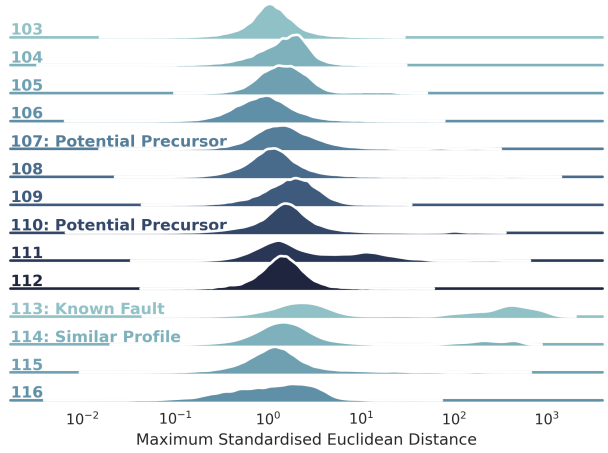


Figure 7: MSED Distribution for test flights. Known fault has a large amount of mass at a much higher MSED than other flights. Potential precursors show smaller mass at similar MSED ( $> 100$ ). Flight 111 also exhibits an odd profile due to irregularity in the shaft speed signal not found in other flights.

The actual processing time for each run was recorded; a run consisted of all currently uploaded models being run once, the MSED being determined and the output buffer being updated if necessary. The other primary resource-intensive operation, input retrieval and clean-up, was ignored since its computational expense is heavily dependent on the data source. Although, in our application its expense was orders of magnitude smaller than that of a run ( $\sim 25\mu s$ ).

Figure 8 shows the distribution of run durations when the 20 models used to evaluate the NASA data set were used and Figure 9 when the 4 models used to evaluate the In-Service Fault Case Study data set were used.

Both histograms have a long tail caused by the fact that the software runs on a non-real-time operating system (unpredictable system interrupts) which is why they are curtailed at the 95th quantile. These results show that the system is capable of operating in real-time e.g. all 20 of the NASA models could be run simultaneously on data at up to  $\sim 50Hz$  and

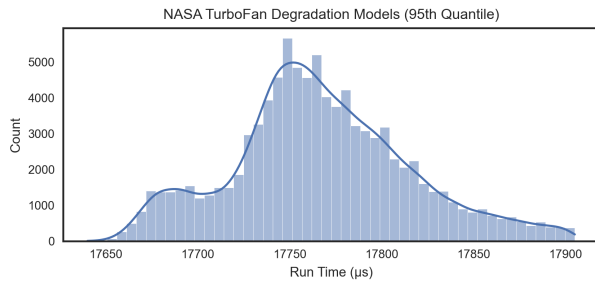


Figure 8: 95th quantile of run durations for the NASA TurboFan Degradation data. Smaller peak likely caused by cache hits.

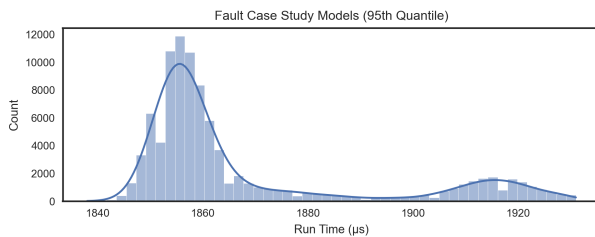


Figure 9: 95th quantile of run durations for the Fault Case Study data. Smaller number of models and overall smaller model size relative to NASA data lead to the lower average run duration and likely the reversal of shape due to more cache hits.

the Fault Case Study models could be run at up to  $\sim 500Hz$ , including all overhead.

## 6 Conclusion

We have presented a data-driven anomaly detection system that runs in real-time on low-power embedded hardware. The system prioritises data based on the Maximum Standardised Euclidean Distance of a flexibly sized ensemble of models. A flexible neural network design is presented that is used to build models for the ensemble from any input-output set and window length without the need for architecture changes. The efficacy of the system has been demonstrated on both real time-series and synthetic snapshot data, plus its real-time operation on embedded

hardware has been proven and characterised.

The system represents a practically realisable, computationally inexpensive method for detection of unusual behaviour in real-time. Further the architecture is flexible and allows for deployment on new systems with minimal need for in-depth domain knowledge.

The system will be deployed on-board in an upcoming flight test program. Future work will focus on analysis of the data returned by this test and how the process described here-in can be applied to higher frequency data.

## 7 Bibliography

### References

- [1] Y. G. Li and P. Nilkitsaranont. Gas turbine performance prognostic for condition-based maintenance. *Applied Energy*, 86(10):2152–2161, October 2009.
- [2] Martha A. Zaidan, Andrew R. Mills, Robert F. Harrison, and Peter J. Fleming. Gas turbine engine prognostics using Bayesian hierarchical models: A variational approach. *Mechanical Systems and Signal Processing*, 70-71:120–140, March 2016.
- [3] Sandra Hofmann, Svetoslav Duhovnikov, and Dominic Schupke. Massive Data Transfer From and to Aircraft on Ground: Feasibility and Challenges. *IEEE Aerospace and Electronic Systems Magazine*, 36(5):6–14, May 2021. Conference Name: IEEE Aerospace and Electronic Systems Magazine.
- [4] Ningbo Zhao, Xueyou Wen, and Shuying Li. A review on gas turbine anomaly detection for implementing health management. In *Turbo Expo: Power for Land, Sea, and Air*, volume 49682, page V001T22A009. American Society of Mechanical Engineers, 2016.
- [5] Edward Glaessgen and David Stargel. The Digital Twin Paradigm for Future NASA

- and U.S. Air Force Vehicles. In *53rd AIAA/ASME/ASCE/AHS/ASC Structures, Structural Dynamics and Materials Conference*, page 118. American Institute of Aeronautics and Astronautics, 2012.
- [6] Fei Tao, He Zhang, Ang Liu, and Andrew YC Nee. Digital twin in industry: State-of-the-art. *IEEE Transactions on Industrial Informatics*, 15(4):2405–2415, 2018. publisher: IEEE.
- [7] Hamid Asgari, XiaoQi Chen, and Raazesh Sainuddin. Modelling and simulation of gas turbines. *International Journal of Modelling, Identification and Control*, 20(3):253–270, 2013. publisher: Inderscience Publishers Ltd.
- [8] V. Panov. Gasturbolib: Simulink library for gas turbine engine modelling. In *Turbo Expo: Power for Land, Sea, and Air*, volume 48821, pages 555–565, 2009.
- [9] D. A. Nix and A. S. Weigend. Estimating the mean and variance of the target probability distribution. In *Proceedings of 1994 IEEE International Conference on Neural Networks (ICNN'94)*, volume 1, pages 55–60 vol.1, June 1994.
- [10] S. Willskyi. *A Survey of Design Methods for Failure Detection in Dynamic Systems\**.
- [11] Paul M. Frank and Xianchun Ding. Frequency domain approach to optimally robust residual generation and evaluation for model-based fault diagnosis. *Automatica*, 30(5):789–804, May 1994.
- [12] Erik Frisk and Lars Nielsen. Robust residual generation for diagnosis including a reference model for residual behavior. *Automatica*, 42(3):437–445, March 2006.
- [13] Tom Brotherton, Al Volponi, Rob Luppold, and Donald L. Simon. eSTORM: Enhanced Self Tuning On-board Real-time Engine Model. Technical report, INTELLIGENT AUTOMATION CORP POWAY CA, January 2003. Section: Technical Reports.
- [14] Takahisa Kobayashi and Donald L. Simon. Hybrid Kalman Filter Approach for Aircraft Engine In-Flight Diagnostics: Sensor Fault Detection Case. pages 745–755. American Society of Mechanical Engineers Digital Collection, September 2008.
- [15] M. Amozegar and K. Khorasani. An ensemble of dynamic neural network identifiers for fault detection and isolation of gas turbine engines. *Neural Networks*, 76:106–121, April 2016.
- [16] Shisheng Zhong, Dan Liu, Lin Lin, Minghang Zhao, Xuyun Fu, and Feng Guo. A novel anomaly detection method for gas turbines using weight agnostic neural network search. In *2020 Asia-Pacific International Symposium on Advanced Reliability and Maintenance Modeling (APARM)*, pages 1–6. IEEE, 2020.
- [17] Weizhong Yan and Lijie Yu. On accurate and reliable anomaly detection for gas turbine combustors: A deep learning approach. *arXiv preprint arXiv:1908.09238*, 2019. tex.ids: yanAccurateReliableAnomaly2019a.
- [18] Weizhong Yan. Detecting gas turbine combustor anomalies using semi-supervised anomaly detection with deep representation learning. *Cognitive Computation*, 12(2):398–411, 2020. tex.ids: yan-DetectingGasTurbine2020a publisher: Springer.
- [19] Geunbae Lee, Myungkyo Jung, Myoungwoo Song, and Jaegul Choo. Unsupervised anomaly detection of the gas turbine operation via convolutional auto-encoder. In *2020 IEEE International Conference on Prognostics and Health Management (ICPHM)*, pages 1–6. IEEE, 2020.
- [20] Jiao Liu, Jinfu Liu, Daren Yu, Myeongsu Kang, Weizhong Yan, Zhongqi Wang, and Michael G. Pecht. Fault detection for gas turbine hot components based on a convolutional neural network. *Energies*, 11(8):2149, 2018. publisher: Multidisciplinary Digital Publishing Institute.
- [21] Hui Luo and Shisheng Zhong. Gas turbine engine gas path anomaly detection using deep

- learning with Gaussian distribution. In *2017 Prognostics and System Health Management Conference (PHM-Harbin)*, pages 1–6. IEEE, 2017.
- [22] Praween Amontamavut, Yuki Nakagawa, and Eiichi Hayakawa. Separated Linux process logging mechanism for embedded systems. In *2012 IEEE International Conference on Embedded and Real-Time Computing Systems and Applications*, pages 411–414. IEEE, 2012.
- [23] Douglas H. Summerville, Kenneth M. Zach, and Yu Chen. Ultra-lightweight deep packet anomaly detection for Internet of Things devices. In *2015 IEEE 34th international performance computing and communications conference (IPCCC)*, pages 1–8. IEEE, 2015.
- [24] Maroua Ben Attia, Chamseddine Talhi, Abdelwahab Hamou-Lhadj, Babak Khosravifar, Vincent Turpaud, and Mario Couture. On-device anomaly detection for resource-limited systems. In *Proceedings of the 30th Annual ACM Symposium on Applied Computing*, pages 548–554, 2015.
- [25] Gianluca Furano, Gabriele Meoni, Aubrey Dunne, David Moloney, Veronique Ferlet-Cavrois, Antonis Tavoularis, Jonathan Byrne, Léonie Buckley, Mihalis Psarakis, Kay-Obbe Voss, and Luca Fanucci. Towards the Use of Artificial Intelligence on the Edge in Space Systems: Challenges and Opportunities. *IEEE Aerospace and Electronic Systems Magazine*, 35(12):44–56, December 2020. Conference Name: IEEE Aerospace and Electronic Systems Magazine.
- [26] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3):1–58, 2009. publisher: ACM New York, NY, USA.
- [27] Nesma M. Rezk, Madhura Purnaprajna, Tomas Nordström, and Zain Ul-Abdin. Recurrent Neural Networks: An Embedded Computing Perspective. *IEEE Access*, 8:57967–57996, 2020. arXiv: 1908.07062.
- [28] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014. publisher: JMLR. org.
- [29] Forrest N. Iandola, Song Han, Matthew W. Moskewicz, Khalid Ashraf, William J. Dally, and Kurt Keutzer. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5 MB model size. *arXiv preprint arXiv:1602.07360*, 2016.
- [30] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [31] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*, 2019.
- [32] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You Only Look Once: Unified, Real-Time Object Detection. *arXiv:1506.02640 [cs]*, May 2016. arXiv: 1506.02640 version: 5.
- [33] Yarín Gal and Zoubin Ghahramani. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. *arXiv:1506.02142 [cs, stat]*, October 2016. arXiv: 1506.02142.
- [34] R. De Maesschalck, D. Jouan-Rimbaud, and D. L. Massart. The Mahalanobis distance. *Chemometrics and Intelligent Laboratory Systems*, 50(1):1–18, January 2000.
- [35] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, and Michael Isard. Tensorflow: A system for large-scale machine learning. In *12th symposium on operating systems design and implementation*, pages 265–283, 2016.

- [36] A. Saxena, K. Goebel, D. Simon, and N. Eklund. Damage propagation modeling for aircraft engine run-to-failure simulation. In *2008 International Conference on Prognostics and Health Management*, pages 1–9, October 2008.
- [37] William Jacobs, Huw Edwards, Visakan Kadirkamanathan, and Andrew Mills. Inter-engine variation analysis for health monitoring of aerospace gas turbine engines. *PHM Society European Conference*, 4(1), July 2018. Number: 1.

## A Appendix 1: Network Reproduction

The loss function (Equation 9) was simplified and reordered for batch-wise implementation since the model’s confidence output  $\alpha = \log(\sigma)$  rather than  $\sigma$ :

$$\mathcal{L} = \frac{1}{2} \sum_{n=1}^N \left( \frac{y_n - \mu}{e^{\alpha_n}} \right)^2 - \sum_{n=1}^N \alpha_n - \frac{1}{2} \log(2\pi) \quad (11)$$

where  $N$  is the batch size. The final term is constant and could therefore be dropped for computational efficiency, although it was kept during our work to keep the meaning of the equation intact.

Important network hyper-parameters not mentioned elsewhere:

- LReLU  $\alpha = 0.3$
- Dropout Rate = 0.5
- Convolutional layers
  - Each used 64 filters in the NASA study
  - Each used 32 filters in the Fault Case Study
    - \* An optimisation found to cause minimal performance loss while reducing computation
  - Padding irrelevant for 1x1 blocks, “same” used internally

- No padding for Temporal and Spatial Fold blocks
  - \* Used to reduce the size of the input
  - \* Equivalent to “valid” padding
- Each dense layer had 64 hidden units for the signals tested
- Early stopping patience was 10 epochs without improvement
- RADAM parameters (adjusted based on validation data for each signal)
  - Learning rate  $1e - 3$  to  $1e - 5$
  - Min learning rate  $1e - 5$  to  $1e - 7$
  - Total steps 50,000 to 200,000
  - Warm-up proportion 0.1 to 0.5

An example implementation of the architecture used on our testing for 7 inputs, looking at 10 seconds of data per window is shown in Figure 10.

## B Appendix 2: Notes on Data Preprocessing

The same process was applied to each dataset:

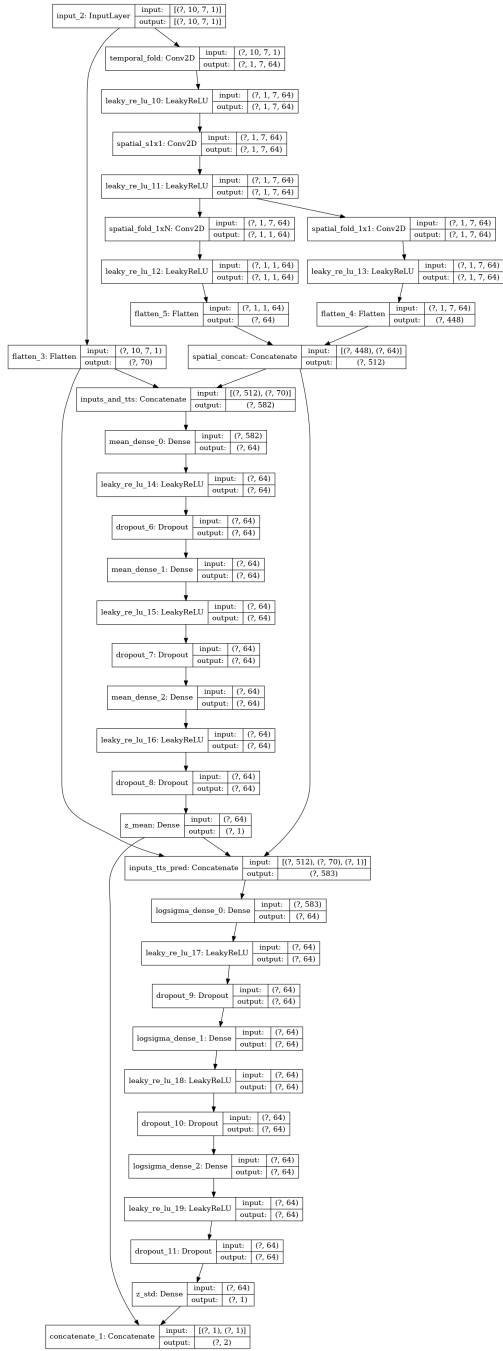
- Data was scaled to the range [0, 1]
  - Based on 2% to 98% quantiles in the training set

## C Appendix 3: In-Service Fault Case Study

The four models were:

- P30 (internal pressure) from outside temperature, outside pressure, and N1
- TGT (Turbine Gas Temperature) from outside temperature, outside pressure, P30, and N1
- T30 (internal temperature) from outside temperature, outside pressure, P30, and N1





- N1 (speed of low speed shaft) from outside temperature, outside pressure, and P30

Figure 10: Neural network architecture with 7 input signals and a 10 second window.