# Multi-Dimensional Self Attention based Approach for Remaining Useful Life Estimation

Zhi Lai[a], Mengjuan Liu[a,*], Yunzhu Pan[a], Dajiang Chen[a]

[a] *School of Information and Software Engineering, University of Electronic Science and Technology of China, Jianshe North Road, Chengdu, 610057, Sichuan, China*

## Abstract

Remaining Useful Life (RUL) estimation plays a critical role in Prognostics and Health Management (PHM). Traditional machine health maintenance systems are often costly, requiring sufficient prior expertise, and are difficult to fit into highly complex and changing industrial scenarios. With the widespread deployment of sensors on industrial equipment, building the Industrial Internet of Things (IIoT) to interconnect these devices has become an inexorable trend in the development of the digital factory. IIoT, combined with data science, provides efficient solutions for the predictive maintenance of complex machines with multiple sensors, and the RUL prediction algorithm is the core component of these solutions. Using the device's real-time operational data collected by IIoT to get the estimated RUL through the RUL prediction algorithm, the PHM system can develop proactive maintenance measures for the device, thus, reducing maintenance costs and decreasing failure times during operation. This paper carries out research into the remaining useful life prediction model for multi-sensor devices in the IIoT scenario. We investigated the mainstream RUL prediction models and summarized the basic steps of RUL prediction modeling in this scenario. On this basis, a data-driven approach for RUL estimation is proposed in this paper. It employs a Multi-Head Attention Mechanism to fuse the multi-dimensional time-series data output from multiple sensors, in which the attention on features is used to capture the interactions between features and attention on sequences is used to learn the weights of time steps. Then, the Long Short-

---

*Mengjuan Liu is the corresponding author.

*Email addresses:* `laizhi727@126.com` (Zhi Lai), `mjliu@ustc.edu.cn` (Mengjuan Liu), `2294523327@qq.com` (Yunzhu Pan), `djchen@uestc.edu.cn` (Dajiang Chen)

Term Memory Network is applied to learn the features of time series. We evaluate the proposed model on two benchmark datasets (C-MAPSS and PHM08), and the results demonstrate that it outperforms the state-of-art models. Moreover, through the interpretability of the multi-head attention mechanism, the proposed model can provide a preliminary explanation of engine degradation. Therefore, this approach is promising for predictive maintenance in IIoT scenarios.

## 1. Introduction

The advent of low-cost micro-sensors and high-bandwidth wireless networks allows us to access machine status data in real-time across physical space. Meanwhile, big data and artificial intelligence technology can mine valuable information from massive data to improve productivity and reduce the risk of failure. Driven by these rising technologies, the Industrial Internet of Things (IIoT) emerged as the times require[1], transforming data from a by-product of the manufacturing process to a strategic resource of great concern to companies. One highly promising application in IIoT is to apply big data technology to the predictive maintenance of complex mechanical equipment[2].

Mechanical failures do not occur regularly for complex mechanical equipment, resulting in higher maintenance time and cost, and these failures usually cause catastrophic consequences, such as aircraft engine damage[3]. Predictive maintenance monitors equipment's operating parameters through sensors embedded in it, uses big data to analyze the cycle, area, and failure types, which is used to predict equipment's RUL, thus, can repair it in advance, ensuring the stability and safety of the system. So far, predictive maintenance has been widely used in many fields such as aerospace, naval, power, gradually replacing event-driven breakdown maintenance and time-driven scheduled maintenance.

RUL estimation algorithm is the core algorithms for predictive maintenance and PHM systems. An accurate RUL prediction can significantly reduce repair and maintenance costs while also minimize the risk of unplanned shutdowns[4]. Moreover, RUL estimation also plays an essential role in other areas, such as product reuse and recycling management, energy consumption,

and landfills[5]. As a result, the significance of RUL estimation has gone beyond PHM to become a key algorithm in several industrial fields. RUL estimation algorithms can be categorized into physics model-based, data-driven, and hybrid algorithms. The physical model-based approach builds a precise model with physical laws, while the data-driven approach makes estimation by analyzing the data[6]. Data-driven algorithms can be further categorized into statistical-based methods and machine-learning-based methods. Statistical approaches include Markov models, Weibull distribution, Kalman filters, etc. In addition, machine learning approaches, including logistic and regression neural networks, train the models to convergence by gradient descent. Finally, the hybrid approach introduces physical modeling into data-driven methods, achieving better performance in scenarios where expertise is accessible[7].

In recent years, data-driven RUL prediction algorithms have received increasing attention. Due to the modern mechanical systems becoming more and more complex, it is nearly impossible to build accurate physical prediction models based on failure mechanisms. Instead, data-driven algorithms generally use an end-to-end design that enables to learn trends in machine degradation directly from sensor data and obtain results directly without concern for intermediate processes[8]. Moreover, data-driven approaches do not require prior expertise in specific fields, thus, are more conducive to industrial deployment. Also, data-driven deep learning approaches are more suitable for massive amounts of data in IIoT scenarios. As the amount increases, deep learning algorithms yield more accurate results than traditional machine learning methods[9][10].

Traditional machine learning algorithms (such as SVM[11], ANN[12], and DBN[13]) have achieved good results in RUL prediction. However, these methods are gradually displaced by neural networks because they are challenging to extract abstract features and do not scale well in big data scenarios. Convolutional neural networks (CNNs) have been proved to be efficient in extracting features from 1D and 2D data and are widely used in the industry for fault detection and RUL prediction[14]. Another prevailing structure is the recurrent neural networks (RNNs)[15]. In the RUL prediction task, the data are generally composed of time series output from multiple sensors. RNNs do well in extracting feature information of the time series, which is not available in CNNs.

However, these networks all treat the time series of multiple sensors equally, but in fact, different sensors have different contributions to the

RUL values. The attention mechanism was used in machine translation tasks[16][17], showing a good performance. Its interpretable weight settings have led to its widespread use in areas such as images[18] and recommendations[19]. The attention mechanism first adds a layer of trainable parameters to the last step's network and then regularizes them into interpretable weights using the softmax function. In this way, the features of the previous step are weighted, and in the meantime, interpretability is achieved. Some models incorporating attention mechanisms have been proposed for the RUL prediction task. According to where the attention mechanisms are used, they are divided into three types: time weighting, feature (sensor) weighting, and two dimensions weighting together. The above methods all use basic attention mechanisms. Nevertheless, in RUL degradation, there may be multiple patterns of sensor degradation, and the basic attention mechanism may not be able to take these degradation paths into account simultaneously. Furthermore, the basic attention mechanism cannot model the intrinsic connection of features.

A deep learning method based on the multi-head self-attention mechanism is proposed in this paper to solve those problems. Unlike the traditional attention mechanism, the self-attention mechanism does not require external information and is better at capturing the internal correlations of features. Additionally, the multi-head attention mechanism uses the attention weights of multiple dimensions to aggregate different contributions of features from multiple perspectives. Therefore, it is more generalizable. The proposed method learns both feature interactions between sensors and weights between sequences. First, it utilizes a self-attention mechanism on features to learn interactions between model features and a self-attention mechanism on sequences to learn the influence weights of different time steps. Then, an LSTM network is set up for learning time-series features. To validate the effectiveness of our approach, we conducted experiments on two benchmark datasets and compared our method with some state-of-the-art models.

The main contributions of this paper are summarized as follows:
- We propose an advanced data-driven model for RUL estimation, which introduces the multi-dimensional multi-head self-attention mechanism to learn data features from various aspects. The proposed model learns interactions between features by a self-attention mechanism on features, then learns the influence weights of different time steps by another self-attention mechanism on sequences, and finally uses LSTM to process multi-dimensional sequence features.
- We conduct extensive experiments on two widely used datasets. The

4

experimental results demonstrate that the proposed model outperforms the state-of-art models. In addition, we present a preliminary analysis of the interpretability of the proposed model.

- We share our experience in dealing with some practical issues in our experiments and discuss and analyze the settings and effects of some critical parameters in detail. Furthermore, we open-source the experiment code, which we believe will help the community develop.

## 2. Related Works

Early RUL prediction algorithms usually used model-based approaches. In [20], the RUL of the aircraft engine bearings is estimated using a bearing spall propagation model with the particle filter-based approach, which utilized vibration and online oil debris sensors to detect spalling of bearings. Coppe et al.[21] use a Paris model with an assumed analytical stress-intensity factor to estimate the remaining useful life of a system experiencing fatigue crack growth. Cheng and Pecht[22] proposed a hybrid model for predicting the RUL of electronic products. A data-driven approach is used to detect faults and extract features, and the PoF model with a data-driven method is used to make predictions.

In recent years, machine learning methods have been widely applied to problems in manufacturing and industrial systems, and have achieved excellent results in RUL prediction. For example, support vector machines (SVM)[11], artificial neural networks (ANN)[12][8] and deep belief network (DBN)[13] algorithms have been widely used for RUL prediction.

With the growing development of machine learning algorithms, neural networks have successfully become the mainstream, among which the first to stand out is the CNN. Sateesh Babu et al.[14] introduced a deep CNN-based regression approach for the turbine engine's RUL estimation and proved the CNN outperformed several traditional algorithms such as MLP, SVR, and RVR on two publicly available datasets. Li et al.[23] also proposed a deep convolution neural network(DCNN) method. Different from the previous method using a 1D CNN, they used a 2D deep CNN to predict RUL. In RUL prediction, CNNs that use large convolutional kernels can extract abstract features directly from the raw data, and the sparse connection and weight sharing strategies allow the network to stack many layers without worrying about gradient vanishing so that the ability to extract features is improved

while training parameters are reduced. Thus, CNNs achieve better results than traditional machine learning methods.

However, CNNs ignore the temporal characteristics of the data. In CNNs, the previous input is entirely unrelated to the following input, while most of the sensor data are time-series data, so RNN, a network designed to process time-series data, is naturally used by many researchers. RNNs can remember every time step's information, and the hidden layer at each time step is determined not only by the input layer at that moment but also by the hidden layer at the previous time step. Zheng et al.[24] used LSTM for the RUL prediction task, and the experimental results obtained on three datasets confirmed the better performance of LSTM over traditional methods as well as CNNs. Liao et.al[25] proposed an LSTM based on Bootstrap for uncertainty prediction of RUL estimation. Elsheikh et al.[26] replaced LSTM with bidirectional LSTM to achieve RUL prediction for random starting short sequences of monitored observations and proposed a safety-oriented objective function to train the network to favor safer early predictions rather than later predictions. Al-Dulaimi et al.[27] proposed a method called HDNN, which mixed LSTM and CNN networks, using LSTM to extract temporal features and CNN to extract spatial features, and the two networks were integrated in parallel.

Graph structure and graph neural networks develop rapidly recently and are widely used in image recognition, recommendation, and other fields. Compared with the hierarchical structure of traditional neural networks, graph structures can express more complex relationships. Some studies have also used graph neural networks in the RUL prediction task, where good results have been achieved. Li et al.[28] proposed a graph structure-based method to predict RUL. The method used directed acyclic graphs to combine LSTM with CNN networks organically instead of simply stacking or integrating LSTM and CNN. Narwariya et al.[29] used graph structures to capture the complex structures inside the device to improve the performance of RUL prediction. The approach used a gated graph neural network (GGNN) to model the internal module structure of the turbine engine, thus dividing the multi-dimensional time-series data into meaningful subsets.

Most deep learning algorithms do not consider that different input features contribute differently to the RUL values. However more important features should be paid more attention to help the model focus more on the critical data so that the results can be more accurate. The Attention mechanism is used to achieve this and has been widely applied to NLP, rec-

ommendation, and many other areas with excellent results.

Das et.al [30] introduced an attention-based bidirectional LSTM network model that weights the time steps. They considered that LSTM only used the output of the last time step for RUL prediction, which might lose the information from the earlier time steps. Xia et al.[31] implemented the self-attention mechanism of Transformer's Encoder module to weigh the time steps of the input data. The attention layers and linear layers are alternately arranged and connected using shortcut connections. Ragab et al.[32] proposed an LSTM based sequence to sequence model with an Attention mechanism for RUL prediction. Cao et al.[33] proposed a model named TCN-RSA for RUL estimation of rolling bearings, which consists of a temporal convolutional network append with a residual self-attention mechanism in the time dimension. Liu and Wang[34] proposed an RUL prediction model based on CNN and self-attention using two paralleled network structures. The left side is a DNN, and the right side is first a CNN and position encoding to process the temporal data, then the time series are weighted using the attention mechanism. Xu et al.[35] applied a dual-stream self-attention model, which uses two attention networks in parallel to extract features. In [36], a feature-attention-based model (AGCNN) is presented to predict the RUL of turbofan engines. The author proposed to apply the attention mechanism directly to the input data, where the more critical features were given greater weights. Then the weighted data were passed into Bi-GRU and CNN to extract long-term dependencies and capture local features.

The articles referenced above only considered the importance of either time steps or features. However, different time steps and features have different RUL contributions so that weighting can extract even more critical information and more accurate results. Chen et al.[37] proposed an attention-based LSTM deep learning framework. They used an attention mechanism to weight the 2D sequence data in time dimension output by LSTM, achieving that time steps and features were weighted simultaneously.

Song et al.[38] presented a distributed attention-based temporal convolutional network (TCN) for RUL prediction. The first attention mechanism passed sensor data collected simultaneously into the softmax function to calculate the weight of each feature at that moment, and the final weight of each feature is the average of the weights obtained from all time steps in the time window. Then the data are multiplied by the weights and passed into TCN for feature extraction, and the outputs are weighted by another attention mechanism in the time dimension. [39] applied the fully Transformer's

structure to RUL prediction. A position encoding layer is used before the encoder to parse the timing information. Then two encoder modules process the input in parallel that process the input from both time and sensor dimension. Finally, the two Encoder outputs are connected and sent into the Decoder and then sent to the linear layer to output the result.

Most of the above attention mechanisms directly regularize the parameters to interpretable weights by softmax. With the research of attention mechanisms going further and further, many fundamental models of attention mechanisms have been proposed, among which multi-head attention is one of the best. Multi-head attention mechanism transforms the input linearly and then feeds it into different attention networks separately, without sharing parameters between different heads, so that data features can be learned from multiple perspectives. The self-attention mechanism is another variation of the attention mechanism. Self-attention focuses more on the interactions between features and can explore feature interactions that ordinary attention cannot capture. Our model is based on the multi-head self-attention mechanism and weights both time steps and features. Thus we can better learn the feature interactions and degradation trends among multi-sensor data.

## 3. Methodology

### 3.1. Multi-Head Attention

The multi-head attention mechanism was proposed in [40] and has achieved remarkable results in machine translation. The attention mechanism can be abstracted as a mapping function from a $Query$ to a series of key-value pairs. The similarity between $Query$ and $Key$ is calculated using the scaled dot product, reflecting the importance of $Value$, i.e., attention weight. The formula of scaled dot-production attention is as follows.

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V \qquad (1)$$

where $d_k$ denotes the dimension of $Key$, he scaling factor $\sqrt{d_k}$ is used to solve the vanishing gradient problem when $d_k$ is large, which helps distribute attention weights and get a better generalization effect.

Multiple scaled dot-production attentions are first connected in multi-head attention, and then a weight matrix is used to map the output to the
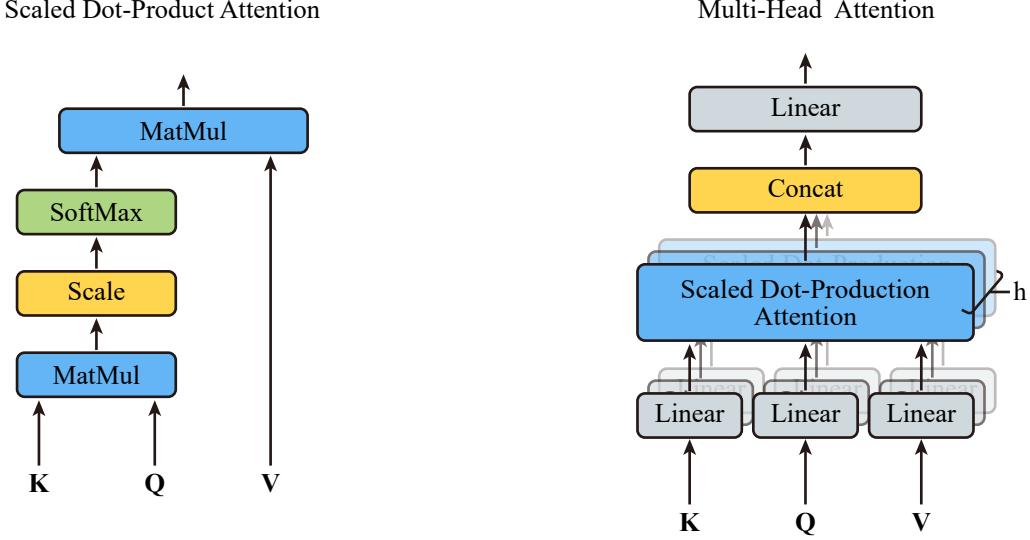
Figure 1: Scaled Dot-Product Attention (left) and Multi-Head Attention Mechanism (right) diagram[40].

size of a single head. The formula is defined as

$$MultiHead(Q, K, V) = Concat(head_1, head_2, ..., head_h)W^O$$
$$where\ head_i = Attention(QW_i^Q, KW_i^K, VW_i^V) \tag{2}$$

where $h$ denotes the number of heads, $W^Q$, $W^K$, $W^V$, and $W^O$ are all trainable parameter matrices. Each *head* is a scaled dot product attention, but the weights obtained in training are different. So multi-head attention can be considered as ensemble learning, aggregating the results of various attention mechanisms with different perspectives. A schematic representation of the scaled dot-production attention and the multi-head attention is shown in Figure 1.

### 3.2. Proposed Method

An overview of the proposed model is shown in Figure 2. First, the model input is a multi-dimensional time series that is produced by multiple sensors. The inputs can be represented as an $F \times T$ matrix representing $T$-step time series from $F$ sensor inputs.

The input matrix with the dimension of $F \times T$ is first fed to the multi-head self-attention mechanism. The $T$ time steps' data of each sensor can
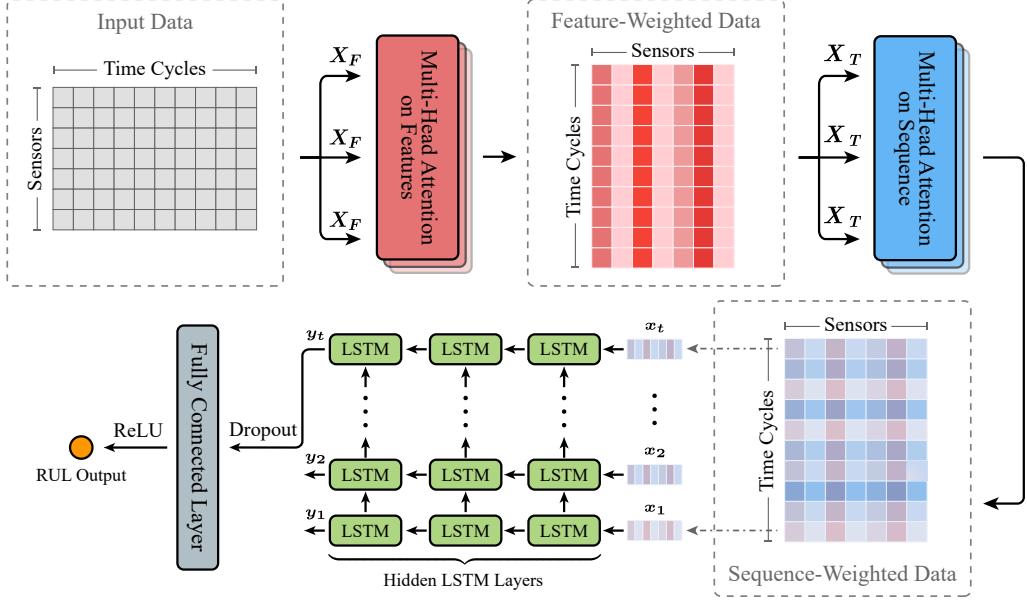
9

Figure 2: Overall structure of the proposed method.

be considered as the embedding of this sensor in the $T$ dimensions, which is expressed as

$$X_F = [S_1, S_2, ..., S_F], \ \ S_i = [s_1^i, s_2^i, ..., s_T^i] \tag{3}$$

where $F$ is the number of sensors, $T$ is the number of time steps, $S_i$ indicates the time series of a single sensor, and $s_j^i$ represents the data of the $i$-th sensor at the $j$-th time step. $j = 1, 2, ..., T$. In self-attention, $Q = K = V$, the self-attention in the sensor dimension can be represented as

$$Output = MultiHead(X_F; X_F; X_F) \tag{4}$$

where $MultiHead$ represents the multi-head attention mechanism introduced in section 3.1, and $X_F$ represents the input of the attention mechanism. The feature map weighted by the sensor is obtained after the attention layer.

The output of the attention mechanism in sensor dimension won't change its dimension and then is weighted by the attention in sequence dimension. In this step, we consider all sensor data in a time step as the embedding of this time step, so the matrix of the input can be expressed as

$$X_T = Out_F = [O_1, O_2, ..., O_T], \ \ O_j = [o_1^j, o_2^j, ..., o_F^j] \tag{5}$$

10

And the attention mechanism can be defined as

$$Output_T = MultiHead(X_T; X_T; X_T) \tag{6}$$

After these two steps, we get a time series weighted in both sequence and sensor dimensions.

The weighted data are then fed into a deep LSTM network to learn the sequence features. By stacking LSTM cells vertically, the LSTM network can learn more high-order sequence features. For each time step of LSTM network, the input of upper LSTM cells is the output of lower cells. The output of the last LSTM layer at the last time step is used as the feature vector, which has contained the information of all previous time steps.

The feature vector is sent through a multilayer perceptron (MLP) and mapped to individual neuron outputs to obtain the final RUL estimation. The dropout[41] was used to mitigate the overfitting problem. In the dropout technique, neurons are randomly masked with probability $P$ during training, while during testing, all neurons are involved in the output.

Considering that RUL estimation is a typical regression problem, we use the mean square error as the loss function.

$$L = \frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2 \tag{7}$$

where $y_i$ and $\hat{y}_i$ denote the real and prediction values of RUL, respectively, and $N$ is the total number of train set samples.

The proposed model is an end-to-end model with all parameters trained jointly. Batch gradient descent algorithms are employed to update the parameters, and the early-stop mechanism is used to prevent overfitting. The details of hyper-parameters will be presented in section 4.3.2.

## 4. Experiment Setting

### 4.1. Dataset

The turbofan engine degradation simulation dataset(C-MAPSS dataset)[42] is one of the most popular public datasets for RUL estimation tasks provided by NASA. The C-MAPSS dataset includes four sub-datasets that collect degradation data for aircraft engines simulated by C-MAPSS in two fault modes and six operating conditions. The PHM08 dataset comes from the

11

challenge at the first PHM conference in 2008 and has the same format as the C-MAPSS dataset, except that the C-MAPSS dataset reveals the results of the test set, while the PHM08 dataset requires the researcher to upload the results to the official server via a webpage[1] to get the test results so that the results of PHM08 are more objective and fair compared to C-MAPSS.
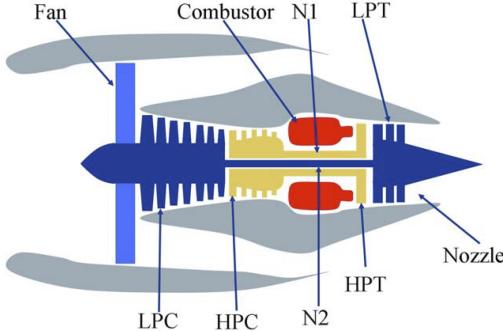


Figure 3: Simplified diagram of turbofan engine simulated in C-MAPSS[43]

Table 1: C-MAPSS Dataset Statistics

| Dataset | C-MAPSS | | | | PHM08 |
|---|---|---|---|---|---|
| | FD001 | FD002 | FD003 | FD004 | |
| Train trjectories | 100 | 260 | 100 | 248 | 218 |
| Test trjectories | 100 | 259 | 100 | 249 | 218 |
| Conditions | 1 | 6 | 1 | 6 | 6 |
| Fault modes | 1 | 1 | 2 | 2 | 1 |

The turbofan engine consists of several components, which are illustrated in figure 3. All datasets provide data for 21 sensors and three environment parameters, which change with time step. The values of the three environment parameters determine which specific operating condition the engine is in. The statistics of the C-MAPSS and PHM08 datasets are shown in Table 1. Each sensor and environment parameter explanation is given in Table 2. For each engine, the engine is healthy at the beginning of the simulation.

---

[1]https://ti.arc.nasa.gov/tech/dash/groups/pcoe/prognostic-data-repository/

Table 2: Sensor and Environment parameters explanation in C-MAPSS Dataset[42]

| Sensor Symbol | Description | Units |
|---|---|---|
| T2 | Total temperature at fan inlet | R |
| T24 | Total temperature at LPC outlet | R |
| T30 | Total temperature at HPC outlet | R |
| T50 | Total temperature at LPT outlet | R |
| P2 | Pressure at fan inlet | psia |
| P15 | Total pressure in bypass-duct | psia |
| P30 | Total pressure at HPC outlet | psia |
| Nf | Physical fan speed | rpm |
| Nc | Physical core speed | rpm |
| epr | Engine pressure ratio (P50/P2) | – |
| Ps30 | Static pressure at HPC outlet | psia |
| phi | Ratio of fuel flow to Ps30 | pps/psi |
| NRf | Corrected fan speed | rpm |
| NRc | Corrected core speed | rpm |
| BPR | Bypass Ratio | – |
| farB | Burner fuel-air ratio | – |
| htBleed | Bleed Enthalpy | – |
| Nf_dmd | Demanded fan speed | rpm |
| PCNfR_dmd | Demanded corrected fan speed | rpm |
| W31 | HPT coolant bleed | lbm/s |
| W32 | LPT coolant bleed | lbm/s |

| Environment Parameters | Description | Units |
|---|---|---|
| Altitude | – | ft. |
| TRA | Throttle resolver angle | deg. |
| Mach number | – | – |

At some random time step, a failure occurs, which propagates, causing the engine to degrade until it is completely damaged. The dataset consists of run-to-failure trajectories for multiple engines. The train set gives complete data of the trajectories, while in the test set, only part of the data is given, and the researcher is required to predict the amount of time steps the engine can still run.
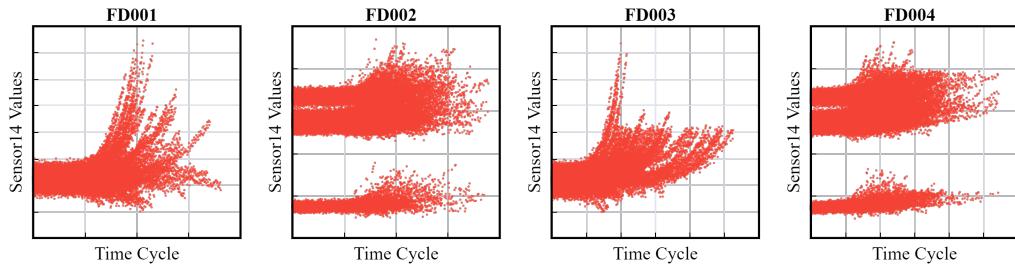


Figure 4: Scatter plot of Sensor 14 data distribution on the C-MAPSS dataset.

Four sub-datasets of C-MAPSS are four combinations of fault modes and working conditions. Since working conditions significantly impact sensor data, more engines are given for the FD002 and FD004 datasets with six working conditions. Figure 4 shows the data distribution of the 14th sensor on FD001 – FD004. It can be found that the engine degradation proceeds in different directions due to the different fault modes and working conditions in the four data sets, which leads to different data distribution in each dataset.

*4.2. Data Pre-Processing*

Data pre-processing has a crucial impact on the performance of the RUL prediction model. In the experiments, a time window technique is used to segment the data. Piece-wise RUL modeling is applied to the degradation process, and regularization is used to normalize the data. For datasets with multiple working conditions, a k-means clustering is also performed to environment parameters.

*4.2.1. Piece-Wise RUL*

Since the engine fails at some unknown time and is healthy before the failure, it is also difficult to predict the exact RUL of a healthy engine. Thus a linear degradation model is essential for the convergence of the model. In piece-wise RUL, shown in Figure 5, a constant value $R_{max}$ is set to indicate
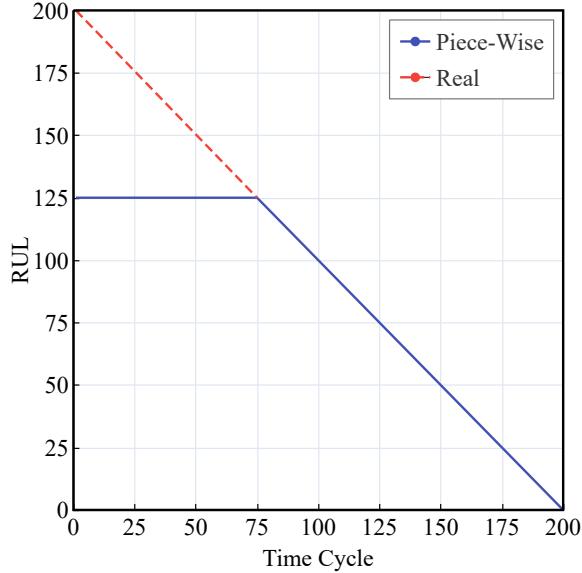
14

Figure 5: Piece-wise RUL Function.

the remaining useful life of a healthy engine. In the beginning, the RUL of the engine is constant and then begins to degrade linearly. Piece-wise RUL is proved to be effective for the C-MAPSS dataset as it helps the network to converge better[23][37]. Following other studies, $R_{max}$ is set to 125 in this experiment. We will discuss the effect of the $R_{max}$ setting on the model performance in subsequent experiments.

*4.2.2. Cluster and Normalization*

The outputs of different sensors need to be normalized because of the different data units. For C-MAPSS and PHM08, the two most commonly used methods are min-max normalization[23] and z-score normalization[25]. In the absence of prior expertise, reasonable max-min values for sensors are difficult to determine, so we use z-score normalization, which is defined as follows:

$$z_i = \frac{x_i - \mu_i}{\sigma_i}, i \in S \tag{8}$$

where $S$ denotes the collection of sensors, and $z_i$ represents the normalized value of the $i$-th sensor, $\mu_i$ and $\sigma_i$ represent the average and standard deviation of the $i$-th sensor, respectively.

Various working conditions have a substantial influence on the sensor
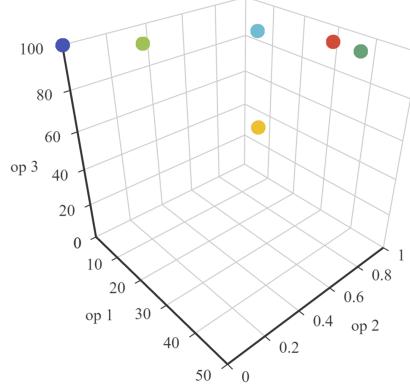
15

Figure 6: Clustering results of environment parameters on FD004 dataset.

data. So it is essential to process the sensor data separately according to the different operational settings. Take the 14th sensor of an engine on FD002 as an example. Figure 7.a) shows the original data of the sensor in degradation progress, where the series data violently oscillates, resulting in the overall trend becoming obscure. Figure 7.b) is the data normalized by formula 8, which changed in the data range but did not change in data distribution. Therefore, we normalize the data separately according to the working conditions at the time of data recording so that the normalized data can reflect the trend of degradation.

We can easily get the classification of six conditions by clustering based on the three operational settings given in the datasets. The results of clustering three operational settings using the K-means algorithm are shown in Figure 6. The values of the different working conditions are very centralized and separated to a large extent, so the clustering results are highly reliable.

After getting the working conditions of the data for each time step, we can normalize the data according to the working condition separately as the following equation:

$$z_{i,j} = \frac{x_{i,j} - \mu_{i,j}}{\sigma_{i,j}}, i \in S, \ j \in C \tag{9}$$

where $C$ represents the set of conditions, $i$ represents the $i$-th sensor, and $j$ represents the $j$-th conditions. Figure 7.c) illustrates the data colored by working conditions, and 7.d) is the data normalized separately by working conditions, which shows that the data trend with the time step is distinct from the unnormalized data.
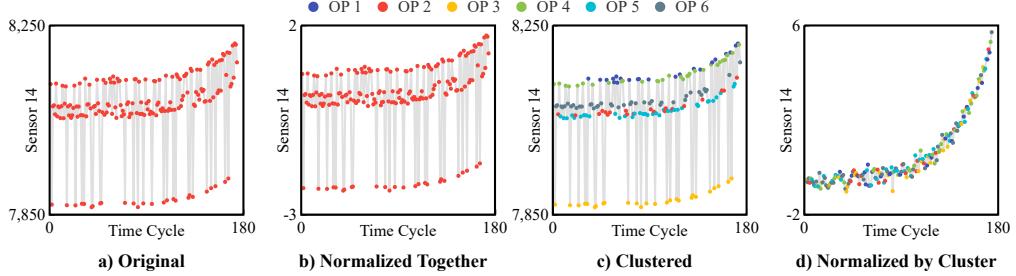
16

Figure 7: Original and normalized trajectory for the 14th sensor of the No.7 engine in the FD002 dataset.

### 4.2.3. Time Window

Time window is a common technique for data segmentation, which is illustrated in Figure 8. For each run-to-failure time series, a time window slides from the beginning to the end. Let the total sequence length be $T_{total}$, the window length be $T$, and the step size of the window sliding is fixed to 1. Then, the RUL of the engine at $t_i$ can be expressed as

$$RUL_i = T_{total} - ti \qquad t_i = 1, 2, 3, ..., T_{total} \tag{10}$$

The number of samples obtained for each sequence can be expressed as

$$Sample\ number = \begin{cases} T_{total} - T + 1 & T \leq T_{total} \\ 1 & T > T_{total} \end{cases} \tag{11}$$

As the test set does not give the whole sequence, the sequence length T may be smaller than the window length. The number of samples at this point is one, and the sequences are filled forward to the window length using the data of the first time step.

### 4.3. Setup

### 4.3.1. Metrics

Two widely adopted metrics were used to evaluate the proposed method's performance: the score function and the root mean square error(RMSE). The score function is the evaluation metric used in the PHM08 challenge
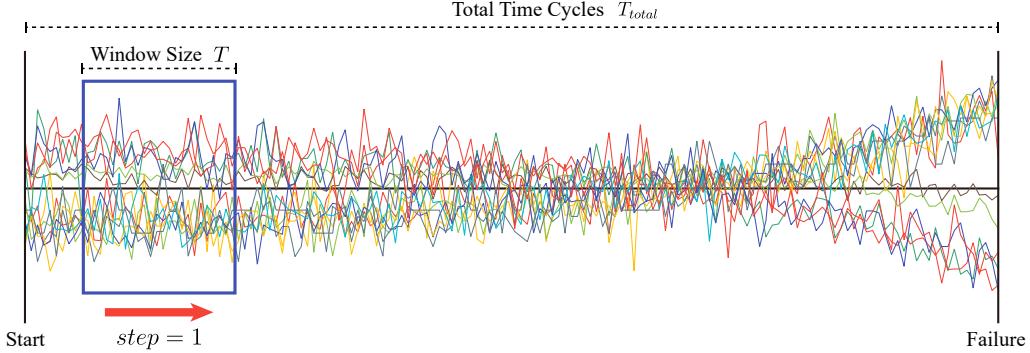
17

Figure 8: Illustration of splitting data with the time window method.

dataset[42], which is formulated as follows:

$$Score = \begin{cases} \sum\limits_{i=1}^{N} e^{-\left(\frac{d}{a_1}\right)} - 1 & for \ d < 0 \\ \sum\limits_{i=1}^{N} e^{\left(\frac{d}{a_2}\right)} - 1 & for \ d \geq 0 \end{cases} \tag{12}$$

$$where \ d = \widehat{RUL} - RUL, \quad a_1 = 13, \quad a_2 = 10$$

The other metric is the RMSE, which has the following equation.

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} d_i^2} \tag{13}$$

A comparison of the two metrics is shown in Figure 9. With increasing error between the prediction and truth values, the RMSE increases linearly, while the Score score increases exponentially. Moreover, the score function penalizes more severely when the prediction is larger than the real RUL, since predicting a larger RUL for safety-critical industrial equipment usually results in more severe consequences than predicting a smaller RUL.

### 4.3.2. Hyper-Parameters

Experiments were conducted on all four sub-datasets of the C-MAPSS dataset and PHM08 dataset. The window length was uniformly set to 30, and all 24 features (3ops+21sensors) provided by the dataset were used.
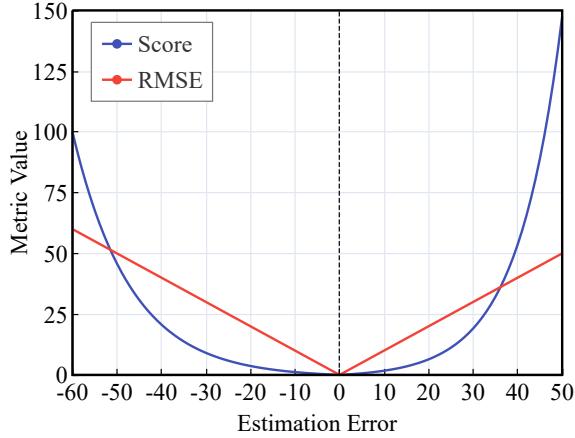
Figure 9: Comparison of RMSE and score functions.

The Adam[44] algorithm was applied to optimize the proposed model. The learning rate of Adam was set to 0.0002. The number of training batches was set to 128, and the early stop mechanism was employed to stop training after 50 rounds without better results.

For the proposed model, the LSTM uses three hidden layers with 100 nodes per layer. The MLP layer has a hidden layer with 100 nodes, the activation function is ReLU[45], and the dropout is set to 0.5. The final output layer uses a single node to predict the RUL value. We open-source the experiment's code on GitHub[2], hoping to contribute to the community for better improvement.

## 5. Results Analysis

This section provides a complete analysis of the experiment results. All experiments are repeated 30 times independently to ensure the accuracy of the results.

### 5.1. Parameter Study

We first investigate the effect of different parameters on the model performance, which include the number of feature heads, the number of sequence heads, the window length, and the piece-wise RUL.

---

[2]https://github.com/LazyLZ/multi-head-attention-for-rul-estimation
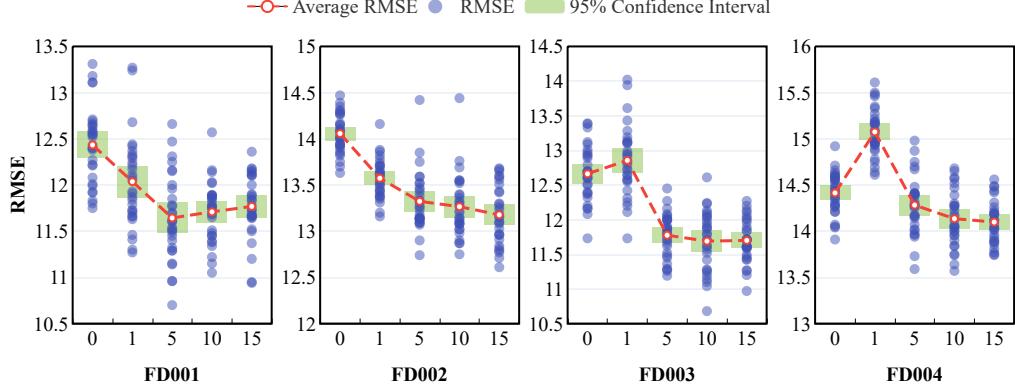
## 5.1.1. Impact of Feature Head



Figure 10: Feature head study on the C-MAPSS dataset.

First, we use a feature attention only with one layer and set the different number of heads to evaluate the impact of feature heads on performance. Note that the head number must be a factor of the embedding dimension. Take the head number as {0, 1, 5, 10, 15} where head=0 represents no attention mechanism, i.e., a simple deep LSTM network. After 30 times independent repeated trials, the results are shown in Figure 10. The blue points are the results (RMSE) of each trial, the red points are the average values, and the green area is the 95% confidence interval of the experimental results. The horizontal coordinate represents the number of feature heads.

On the FD001 and FD002 datasets with only one fault mode, the RMSE decreases as the number of heads increases. While on the FD003 and FD004 datasets with two fault modes, the single-head attention's performance is not even as good as the LSTM. However, the model performance is still remarkably improved as the head number increases, indicating that the single-head attention is incapable of handling the degradation trend of multiple fault modes simultaneously, while the multi-head attention is competent. On all four datasets of C-MAPSS, better results are obtained on head=10 or 5, while the gains obtained by further increasing the head count are not significant. Overall, experiments on all four datasets demonstrate that the multi-head attention mechanism outperforms traditional attention mechanisms and significantly contributes to the prediction performance.
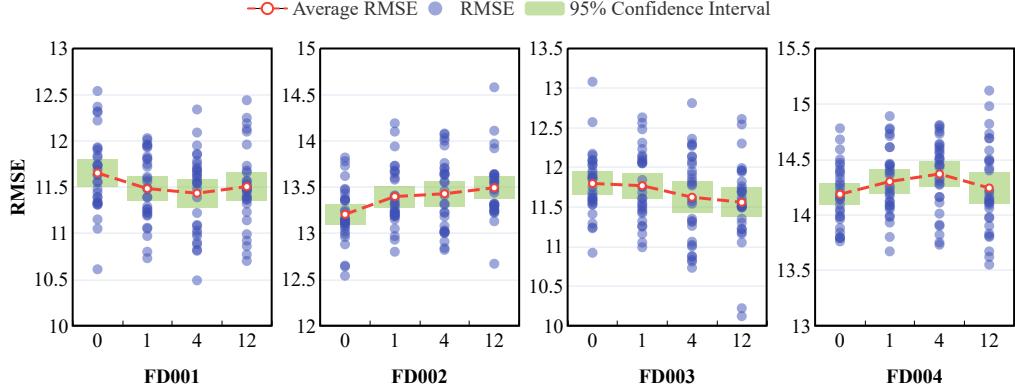
Figure 11: Sequence head study on the C-MAPSS dataset.

## 5.1.2. Impact of Sequence Head

We also study the impact of sequence attention. Based on the study on feature attention, we set the head of feature attention to 5 on FD001 and FD003, 15 on FD002 and FD004, respectively. Then we adjust the sequence attention to {0, 1, 4, 12}, where 0 represents only feature attention.

The results are shown in Figure 11. First, the improvement in sequence attention is generally less pronounced than feature attention because the data is generated with noise. For the feature vector consisting of time series, time sequence performs noise reduction, while for sequence attention, the feature vector consisting of data from 24 sensors does not have this effect.

On the FD001 dataset, the best results are obtained at head=4. A higher number of heads increases the training time, but with worse results. The results for FD003 are similar to FD001's, except that the performance has not yet reached the inflection point at head=12. For FD002 and FD004, the performance becomes worse regardless of the number of heads. Therefore, sequence attention has no significant effect on the performance for data with multiple conditions.

## 5.1.3. Impact of RNN Cell

There are many variants of RNN. In this paper, we select four of the most commonly used structures: basic RNN, GRU, LSTM, and Bidirectional LSTM, to study the effects of different RNN cells on performance. We use the hyper-parameters of the attention layer obtained in the subsections above and only change the cell structure of the RNN. The experiment results are
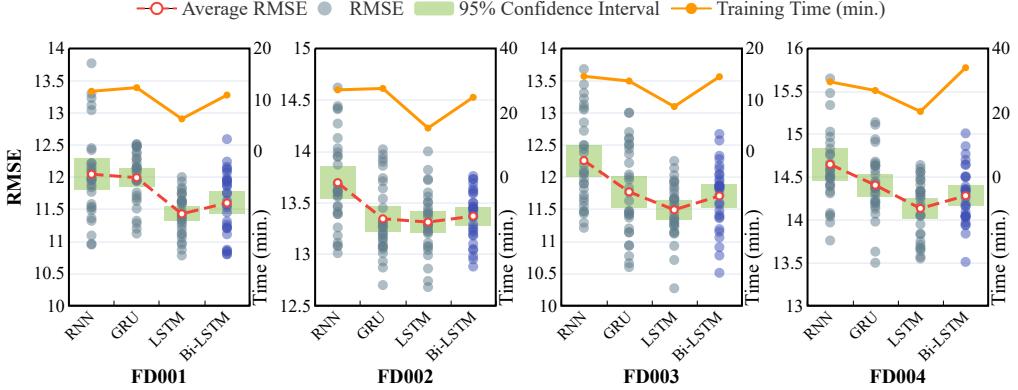
Figure 12: RNN cell study on the C-MAPSS dataset.

shown in Figure 12, where the green intervals, red dots, dashed lines, and blue dots have the same meaning as the experiments above. The orange line represents the average training time of the model for each experiment.

Among them, RNN has the worst performance, followed by GRU. Meanwhile, the repeated trials of these two models are not as stable as others. Although the structure of these models is simpler with less computation, it also leads to difficulties in convergence and more training epochs. Therefore its training time is longer than LSTM and even longer than Bi-LSTM in FD001 and FD002 datasets. Bi-LSTM performs better than GRU but worse than LSTM because the bidirectional LSTM increases the model complexity. Further, the Bi-LSTM is intended to infer the following content through the bidirectional relationship of the sequences. In contrast, the signal sequence of the sensor in the RUL estimation scenario is meaningless when reversed. Therefore BiLSTM did not achieve the desired results.

### 5.1.4. Impact of Window Size

Most studies using the C-MAPSS dataset set the window size and $R_{max}$ uniformly to 30 and 125. So we do the same for comparison with other studies. However, we believe that window size and $R_{max}$ are still valuable to investigate.

We adjust the window size to {1, 10, 20, 30, 40, 50, 60} on four datasets of C-MAPSS, with other parameters as above, and obtain the results in Figure 13 after 30 independent repeated trials for each configuration. The RMSE first decreases and then increases as the window size increases on all four
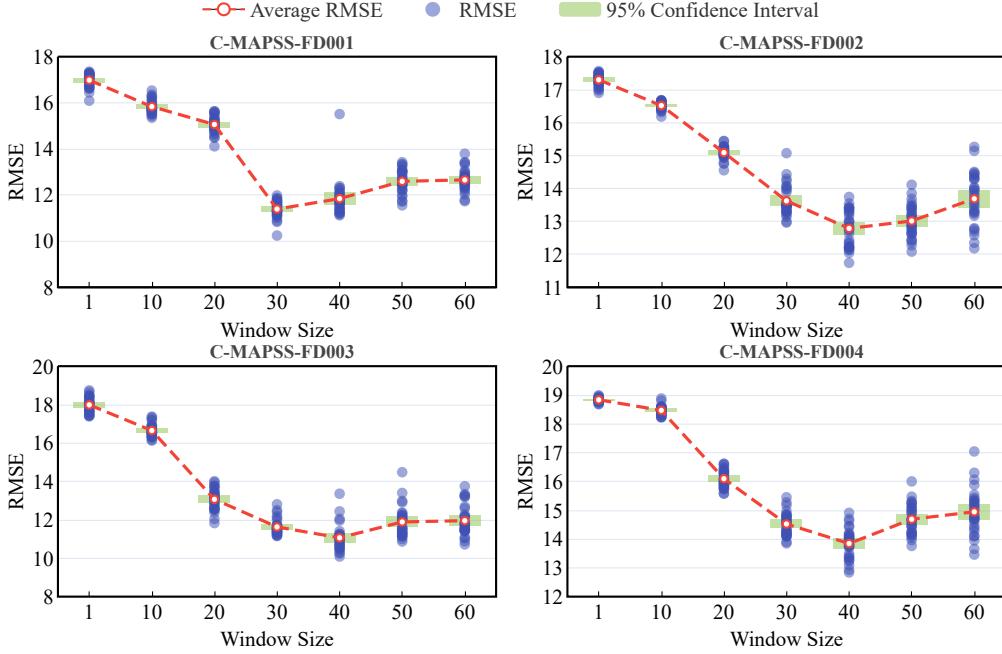
Figure 13: Window size study on the C-MAPSS dataset.

data sets, which indicates that an excessively long window size does not help the performance. This is because the current health status of the engine is mainly correlated with the data from the most recent period, and the correlation decreases as the time interval becomes larger. Taking these data into consideration may instead introduce too much noise, which is partially evidenced in the experiments. On the other hand, the variance of the trial results gets larger as the window size increases, indicating that the model training is indeed disturbed by noise. On FD001, the best results are achieved with a window size of 30, while on the other three datasets, the best results appear at 40.

### 5.1.5. Impact of Piece-Wise RUL

Similar to the experiments of window lengths. We set the $R_{max}$ to {100, 110, 120, 125, 130, 140, 150}. Each experiment is set to use or not use the piece-wise RUL for the test set separately as a comparison. The results are shown in Figure 14, where the red dots and lines are for the test set without piece-wise RUL and blue is for using it.

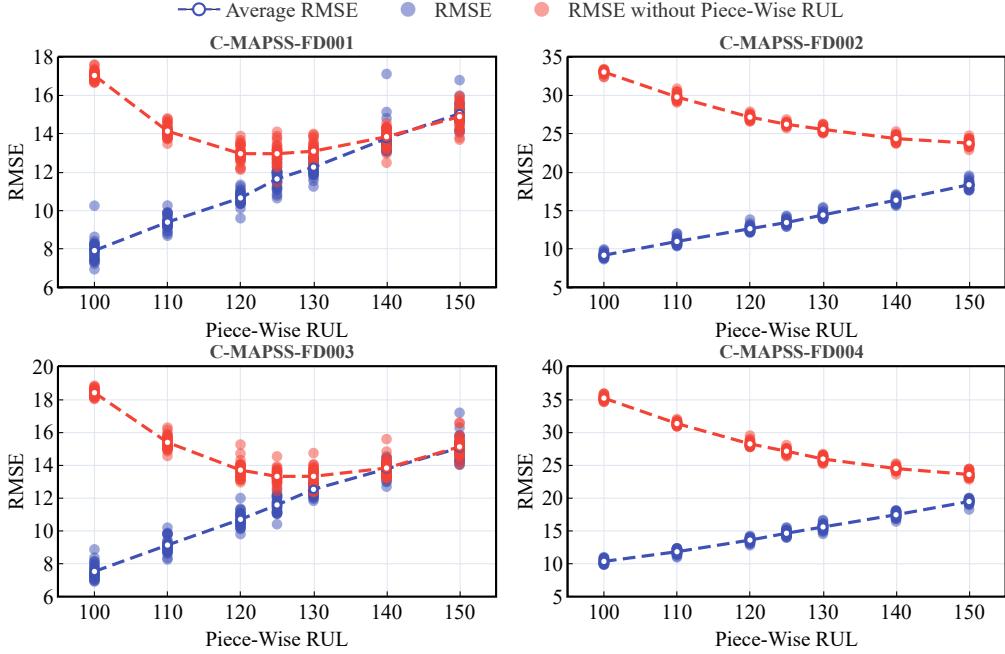On FD001 and FD003, the results without piece-wise RUL decreases and

Figure 14: Piece-wise RUL study on the C-MAPSS dataset.

then rises as $R_{max}$ goes from small to large, while on FD002 and FD004, the RMSE of the red line decreases and the blue line increase as the $R_{max}$ rises. This suggests that a smaller numerical RMSE obtained by setting a very small $R_{max}$ does not indicate that the model has good predictive ability. Because this is achieved by setting too many engines' RUL to $R_{max}$ on the test set. As an extreme example, if $R_{max}$ is set to 0, the test set results will all be 0, and the model will be trained as a model with a constant output of 0, which gives a very low RMSE but does not demonstrate the prediction ability of the model. On the FD001 and FD003 datasets, the best results are indeed obtained when $R_{max}$=125, while on the other two datasets, setting $R_{max}$ to greater than or equal to 150 may be a better choice, although this may not look good numerically.

*5.2. Ablation Study*

In this section, we design an ablation experiment to verify the effectiveness of each module of the proposed model. The experimental results on the four sub-datasets of C-MAPSS are shown in Figure 15. Where $L$ denotes the single LSTM network, $A$ denotes the single head attention mechanism on
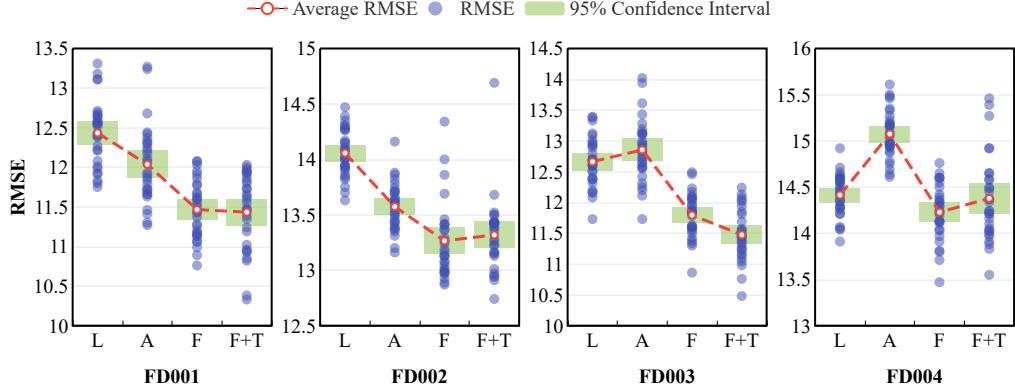
Figure 15: Ablation study on the C-MAPSS dataset.

features, $F$ denotes using only multi-head attention on features, and $F + T$ denotes using multi-head attention both on features and sequences.

From the perspective of the fault mode, LSTM performs the worst on FD001 and FD002 with only one fault mode. Single-head attention has a significant performance improvement than LSTM, and multi-head attention mechanism has another enormous improvement compared to single-head attention. As for FD003 and FD004 with multiple fault modes, the single-head attention mechanism is even worse than the simple LSTM, but the model with multi-head attention on feature and the combined model of the two types of attention outperform both the LSTM and the single-head attention model. In terms of working conditions, the model with two types of attention performs best for FD001 and FD003 with only one working condition. On FD002 and FD004 with multiple working conditions, the model with two types of attention is slightly worse than the model with multi-head attention on features.

From the above observations, we can conclude that the results of multi-head attention-based models are better than the single-head attention-based model on all four datasets of C-MAPSS since the multi-head attention mechanism learns the weights of the features from different perspectives by using a linear transformation on the inputs. In contrast, single-head attention mechanism cannot consider multiple degradation trends, therefore, the results are not good on the datasets with multiple failure modes (FD003, FD004). The models with two types of attention mechanisms perform poorly on the dataset with multiple working conditions caused by the changes in working

25

conditions that make the sequence data oscillate, resulting in excessive noise learned by the attention mechanism on the sequence.
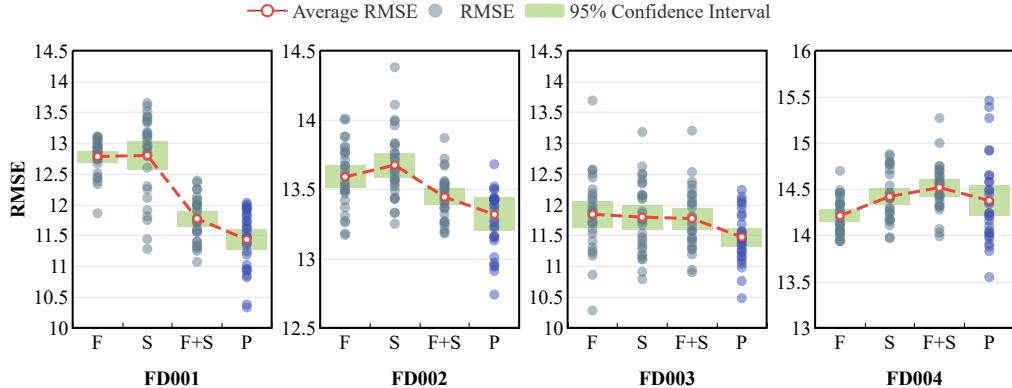
## 5.3. Attention Scheme Study



Figure 16: Attention scheme study on the C-MAPSS dataset.

In order to investigate the performance of our model's attention scheme compared to other attention schemes, we referred to the attention mechanism used in [38], weighting the data in different ways and compared it with our model. The experiment results are shown in Figure 16. Where P is the proposed model; F and S are basic attention weighted in feature and time dimension, respectively; F+S represents first weighted by features and then weighted by time. The results show that the basic attention mechanism weighted on features achieves the best results on FD004, probably due to the simple model convergence better. On the other three datasets, the models weighted for both dimensions (F+S, P) are better than those weighted for a single dimension, indicating that attention usage in multiple dimensions contributes to the model performance. Furthermore, our model obtained better results than the model using the basic attention, which demonstrates that the multi-head self-attention performs better than that of the basic attention model.

## 5.4. Case Study

Figure 17 shows the prediction results of all engines on the FD001 – FD004 test sets. The engines are sorted by real RUL ascendingly in order to observe the model's performance among different RULs better. The model's
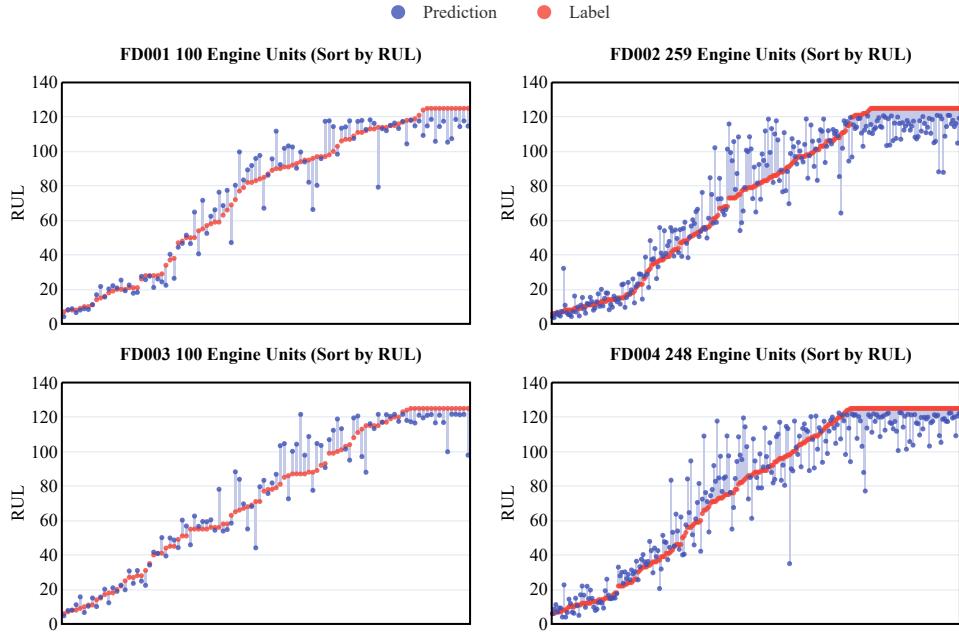
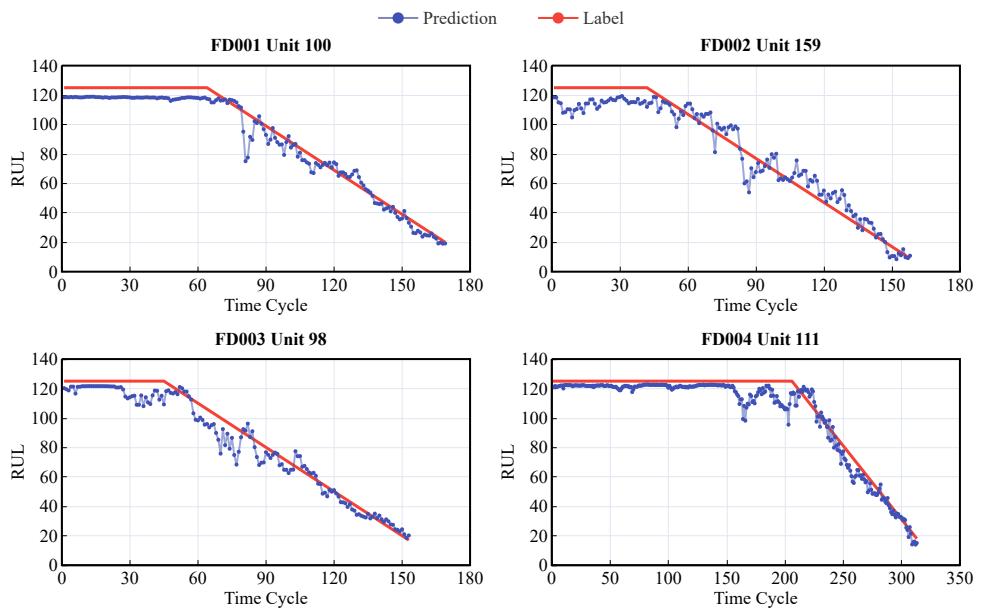Figure 17: Predictions for all engines on the four datasets.



Figure 18: Example of continuous prediction for a single engine.

predictions are generally accurate when RUL is greater than 120 and less than 50 but are comparatively inaccurate between 50 and 120. In the case of a damaged engine, the data of each sensor converge more as the fault spreads, which helps the model predict the RUL more accurately. When the RUL is greater than 120, the model considers the engine to be healthy because of the piece-wise RUL mechanism, which predicts the model's RUL to be $R_{max}$.

Although the test set only gives partial sequences, we can calculate the RUL for each time step of the given sequences based on the final RUL. Figure 18 shows the results of all predictions for the four engines selected from FD001~FD004. The proposed model results are stable when continuously predicting the RUL of a standalone engine. As the engine approaches failure, the prediction becomes more accurate, proving the effectiveness of our model in estimating RUL.

*5.5. Interpretability of Attention*

The weights of the attention mechanism are interpretable because it uses the softmax function to output probabilities so that we can get some interpretations of the model predictions by the weights of attention.
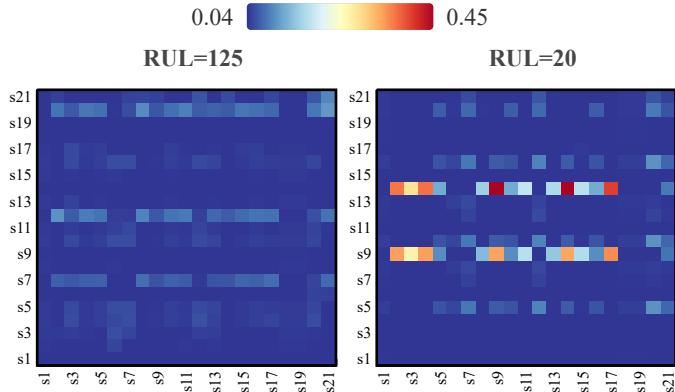


Figure 19: Heat map of the self-attention weight matrix for the No.100 engine of the FD001 test set

Figure 19 shows the heat map of the weights of the self-attention mechanism for the 21 features, where each column represents the corresponding sensor's attention to each sensor. The attention weights are assigned differently for different sensors, which indicates that the self-attention mechanism

learns the relationship between sensors. Moreover, at the early stage of degradation, the attention among sensors is relatively even, while the attention of specific sensors increases sharply at the end of degradation. Also, these sensors are commonly attended by other sensors, such as s9, s14 in Figure 19, indicating that these sensors are vital to the degradation of the engine, and these sensors with higher mutual attention weights may also be correlated in their physical structure.

It is noted that the most degraded component is not always the component from which the fault starts. As faults propagate, they always spread to multiple components, and some components are more sensitive to the engine's condition so that they also receive more attention weight.
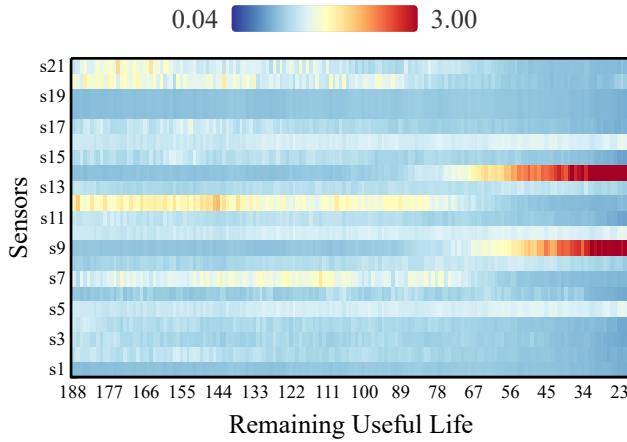


Figure 20: Heat map of the weight sum for the No.100 engine of the FD001 test set

For each time step, by summing the attention of each sensor, we can obtain the trend of the attention of the sensors over the time step. As shown in Figure 20, the attention is distributed relatively evenly at the beginning. While at the end of the degradation, the attention shifts rapidly to several sensors. This may indicate the propagation process of degradation, and the sensors with attention may be more indicative of the increased effect on RUL, i.e., more severe degradation. Moreover, we can find that the shift in attention always occurs in the last 30-40 time steps, which supports the reason why the model achieves the best results when the window size is taken as 30-40 because the data from these last time steps are genuinely representative of the current engine state.
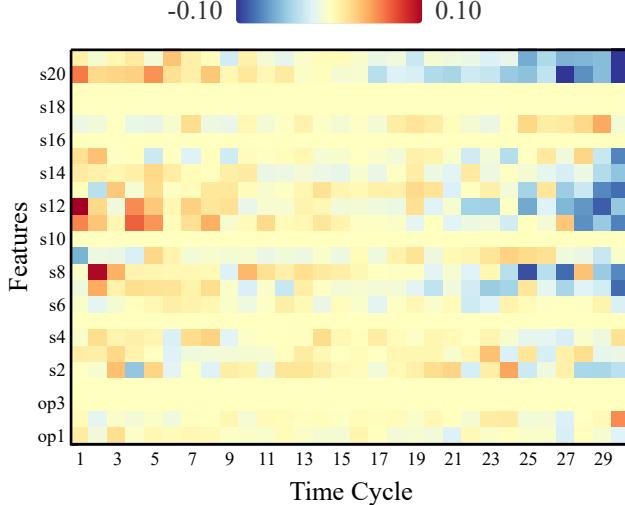
Figure 21: Heat map using SHAP analysis on FD001 test set

Finally, we use SHAP[46], a general interpretable method for machine learning, to explain our model. The model and experiment results of the FD001 test are used to analyze the contribution of each feature of the input to the prediction results, shown in figure 21. The x-axis and y-axis denote time steps and features, respectively. Each square of the heat map is the model input at the corresponding position. The bluer color of the box means that the input at that position is negatively correlated with the results, and contrarily, the redder the positive correlation. As illustrated by the heat map, the most concerned data of the model are the tails and heads of the input sequence, with the tail of the sequence being particularly essential, which is substantially consistent with the fact. From the sensor's perspective, the results of the SHAP analysis also corroborate the effectiveness of our attention mechanism. The attention mechanism concerns sensors such as S9 and S14, which are also those inputs in SHAP analysis that significantly impact the experiment results.

## 5.6. Comparison with other work

### 5.6.1. Results on C-MAPSS

Since studies of RUL prediction widely use the C-MAPSS dataset, we can make a side-by-side comparison of these works. The comprehensive comparison of the proposed method with state-of-art studies is presented in Table

30

3. Bold font indicates the best result, and underlined indicates the second-best result. The proposed method achieves the best results in most cases. The differences in the results obtained on the four datasets are due to the different fault modes and working conditions. First, FD002 and FD004 have significantly worse metrics than FD001 and FD003 in all studies because the variation in working conditions causes the serial data to lose its original trajectory with component degradation. Even by regularizing the data, the oscillations due to changes in operating conditions cannot be completely neutralized.

In all studies, the results for FD001 and FD003 are similar, which indicates that the fault modes have little effect on performance. Thus the model should achieve comparable results for a wide range of fault modes that may occur in practice. Our study achieves outstanding results for FD002 and FD004 due to normalization based on the working conditions separately so that the degradation trend that would have been disturbed by the working conditions can be recovered to some extent, which makes the model can learn the degradation trend of the sequence better.

### 5.6.2. Results on PHM08

We also validate our model on the PHM08 dataset. The test set results for the PHM08 dataset are not given directly but need to be submitted to the official website to receive the score values. The data for the PHM08 challenge dataset was generated from on fault mode and six working conditions, which is consistent with the FD002 dataset. Note that the test set used for the

Table 3: Results comparison of different methods on C-MAPSS datasets.

| Methods | Year | FD001 | | FD002 | | FD003 | | FD004 | |
|---|---|---|---|---|---|---|---|---|---|
| | | RMSE | Score | RMSE | Score | RMSE | Score | RMSE | Score |
| LSTM[24] | 2017 | 16.14 | 338 | 24.49 | 4450 | 16.18 | 852 | 28.17 | 5550 |
| DCNN[23] | 2017 | 12.61 | 273 | 22.36 | 10412 | 12.64 | 284 | 23.31 | 12466 |
| Bi-LSTM[26] | 2018 | 13.65 | 295 | 23.18 | 4130 | 13.74 | 317 | 24.86 | 5430 |
| DAG[28] | 2019 | 11.96 | 229 | 20.34 | 2730 | 12.46 | 284 | 22.43 | 3370 |
| HDNN[27] | 2019 | 13.02 | 245 | <u>15.24</u> | <u>1282</u> | 12.22 | 287 | <u>18.16</u> | <u>1527</u> |
| GNMR[29] | 2020 | 12.14 | **212** | 20.85 | 3196 | 13.23 | 370 | 21.34 | 2795 |
| Attn-LSTM[37] | 2020 | 14.53 | 322 | - | - | - | - | 27.08 | 5649 |
| AGCNN[36] | 2020 | 12.42 | 226 | 19.43 | 1492 | 13.39 | <u>227</u> | 21.5 | 3392 |
| DA-TCN[38] | 2020 | <u>11.78</u> | 229 | 16.95 | 1842 | **11.56** | 257 | 18.23 | 2317 |
| **Proposed (F)** | 2021 | **11.71** | <u>223</u> | **13.26** | **1077** | <u>11.69</u> | **192** | **14.14** | **1375** |
| **Proposed (F+T)** | 2021 | **11.43** | <u>209</u> | **13.32** | **1058** | <u>11.47</u> | **187** | **14.38** | **1618** |

Table 4: Results comparison of different methods on PHM08 dataset.

| Methods | Score |
|---------|-------|
| SVR[14] | 15886 |
| RVR[14] | 8242 |
| MLP[14] | 3212 |
| CNN[14] | 2056 |
| LSTM[24] | 1862 |
| Attn-LSTM[37] | <u>1584</u> |
| Proppsed | **1060** |

competition had 435 engines, while the current website provides a dataset of 218 engines for the evaluation. Table 4 shows the comparison of our model with other methods. Compared with CNN, LSTM, attention-based LSTM, and other methods, our model achieves better result, which is consistent with our experimental results on the C-MAPSS dataset.

## 6. Conclution

This paper proposes a data-driven deep learning model for predicting the remaining useful life of complex machines containing multiple sensors. The proposed model learns data features through a multi-dimensional self-attention mechanism. In particular, attention on feature dimension is used to learn interactions between model features, and attention on sequences is used to learn the influence weights at different time steps. An LSTM network is applied to learn the sequential features. Finally, an MLP is used to get the final RUL results.

Extensive experiments on the C-MAPSS dataset and PHM08 dataset demonstrate the effectiveness of our model in dealing with the RUL estimation tasks for multi-dimensional time series. Time windows, piece-wise RUL, z-score regularization, and K-mean clustering are utilized to pre-process the data. On FD002 and FD004 datasets, normalizing data by working conditions respectively make a vast improvement to performance.

Unlike [37][23][38], the proposed method does not require feature selection which makes it more generalizable. We investigate the impact of window size, the number of heads in the multi-head attention mechanism, and $R_{max}$ on

the model and use them to tune our model performance to the optimum. We also demonstrate that each part of our proposed model contributes to the performance through the ablation experiments. Preliminary studies on the interpretability of the attention mechanism show that the multi-head attention of the proposed model is able to focus on the degradation trend of the engine, which is very promising for interpretable prognosis and health maintenance systems.

In our future work, we will continue focusing on the application of the attention mechanism in RUL estimation, such as Transformer and its variant structures. Graph neural networks are also one of the rapidly developing branches of deep learning, and we will later try GNN-based model structures as well. Finally, we will try to apply our models to more industrial scenarios.

## Acknowledgment

## References

[1] W. Wang, N. Kumar, J. Chen, Z. Gong, X. Kong, W. Wei, H. Gao, Realizing the potential of the internet of things for smart tourism with 5g and ai, IEEE Network 34 (6) (2020) 295–301. `doi:10.1109/MNET.011.2000250`.

[2] M.-A. Koulali, S. Koulali, H. Tembine, A. Kobbane, Industrial internet of things-based prognostic health management: A mean-field stochastic game approach, IEEE Access 6 (2018) 54388–54395. `doi:10.1109/ACCESS.2018.2871859`.

[3] R. Kothamasu, S. Huang, W. VerDuin, System health monitoring and prognostics – a review of current paradigms and practices, Handbook of Maintenance Management and Engineering 28 (2006) 1012–1024. `doi:10.1007/s00170-004-2131-6`.

[4] X.-S. Si, W. Wang, C.-H. Hu, D.-H. Zhou, Remaining useful life estimation – a review on the statistical data driven approaches, European Journal of Operational Research 213 (1) (2011) 1–14.

doi:https://doi.org/10.1016/j.ejor.2010.11.018.
URL https://www.sciencedirect.com/science/article/pii/
S0377221710007903

[5] M. Mazhar, S. Kara, H. Kaebernick, Remaining life estimation of used components in consumer products: Life cycle data analysis by weibull and artificial neural networks, Journal of Operations Management 25 (6) (2007) 1184–1193. arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1016/j.jom.2007.01.021, doi:https://doi.org/10.1016/j.jom.2007.01.021.
URL https://onlinelibrary.wiley.com/doi/abs/10.1016/j.jom.2007.01.021

[6] M. Pecht, Prognostics and health management of electronics, Encyclopedia of structural health monitoring (2009).

[7] M. G. Pecht, A prognostics and health management roadmap for information and electronics-rich systems, IEICE ESS Fundamentals Review 3 (4) (2010) 4_25–4_32.

[8] L. Peel, Data driven prognostics using a kalman filter ensemble of neural network models, in: 2008 International Conference on Prognostics and Health Management, 2008, pp. 1–6. doi:10.1109/PHM.2008.4711423.

[9] P. Wen, Y. Li, S. Chen, S. Zhao, Remaining useful life prediction of iiot-enabled complex industrial systems with hybrid fusion of multiple information sources, IEEE Internet of Things Journal 8 (11) (2021) 9045–9058. doi:10.1109/JIOT.2021.3055977.

[10] L. Ren, Y. Liu, X. Wang, J. Lü, M. J. Deen, Cloud–edge-based lightweight temporal convolutional networks for remaining useful life prediction in iiot, IEEE Internet of Things Journal 8 (16) (2021) 12578–12587. doi:10.1109/JIOT.2020.3008170.

[11] X. Wang, H. Gu, L. Xu, C. Hu, H. Guo, A svr-based remaining life prediction for rolling element bearings, Journal of Failure Analysis and Prevention 15 (4) (2015) 548–554.

[12] Z. Tian, L. Wong, N. Safaei, A neural network approach for remaining useful life prediction utilizing both failure and suspension histories, Mechanical Systems and Signal Processing 24 (5)

(2010) 1542–1555, special Issue: Operational Modal Analysis. `doi:https://doi.org/10.1016/j.ymssp.2009.11.005`. URL `https://www.sciencedirect.com/science/article/pii/S088832700900377X`

[13] C. Zhang, P. Lim, A. K. Qin, K. C. Tan, Multiobjective deep belief networks ensemble for remaining useful life estimation in prognostics, IEEE Transactions on Neural Networks and Learning Systems 28 (10) (2017) 2306–2318. `doi:10.1109/TNNLS.2016.2582798`.

[14] G. Sateesh Babu, P. Zhao, X.-L. Li, Deep convolutional neural network based regression approach for estimation of remaining useful life, in: S. B. Navathe, W. Wu, S. Shekhar, X. Du, X. S. Wang, H. Xiong (Eds.), Database Systems for Advanced Applications, Springer International Publishing, Cham, 2016, pp. 214–228.

[15] F. O. Heimes, Recurrent neural networks for remaining useful life estimation, in: 2008 International Conference on Prognostics and Health Management, 2008, pp. 1–6. `doi:10.1109/PHM.2008.4711422`.

[16] D. Bahdanau, K. Cho, Y. Bengio, Neural machine translation by jointly learning to align and translate (2016). `arXiv:1409.0473`.

[17] Y. Kim, C. Denton, L. Hoang, A. M. Rush, Structured attention networks (2017). `arXiv:1702.00887`.

[18] J. Hu, L. Shen, S. Albanie, G. Sun, E. Wu, Squeeze-and-excitation networks (2019). `arXiv:1709.01507`.

[19] G. Zhou, X. Zhu, C. Song, Y. Fan, H. Zhu, X. Ma, Y. Yan, J. Jin, H. Li, K. Gai, Deep interest network for click-through rate prediction, in: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '18, Association for Computing Machinery, New York, NY, USA, 2018, p. 1059–1068. `doi:10.1145/3219819.3219823`. URL `https://doi.org/10.1145/3219819.3219823`

[20] N. Bolander, H. Qiu, N. Eklund, E. Hindle, T. Rosenfeld, Physics-based remaining useful life prediction for aircraft engine bearing prognosis, in: Annual Conference of the PHM Society, Vol. 1, 2009.

[21] A. Coppe, M. J. Pais, R. T. Haftka, N. H. Kim, Using a simple crack growth model in predicting remaining useful life, Journal of Aircraft 49 (6) (2012) 1965–1973.

[22] S. Cheng, M. Pecht, A fusion prognostics method for remaining useful life prediction of electronic products, in: 2009 IEEE International Conference on Automation Science and Engineering, 2009, pp. 102–107. `doi:10.1109/COASE.2009.5234098`.

[23] X. Li, Q. Ding, J.-Q. Sun, Remaining useful life estimation in prognostics using deep convolution neural networks, Reliability Engineering & System Safety 172 (2018) 1–11. `doi:https://doi.org/10.1016/j.ress.2017.11.021`.
URL `https://www.sciencedirect.com/science/article/pii/S0951832017307779`

[24] S. Zheng, K. Ristovski, A. Farahat, C. Gupta, Long short-term memory network for remaining useful life estimation, in: 2017 IEEE International Conference on Prognostics and Health Management (ICPHM), 2017, pp. 88–95. `doi:10.1109/ICPHM.2017.7998311`.

[25] Y. Liao, L. Zhang, C. Liu, Uncertainty prediction of remaining useful life using long short-term memory network based on bootstrap method, in: 2018 IEEE International Conference on Prognostics and Health Management (ICPHM), 2018, pp. 1–8. `doi:10.1109/ICPHM.2018.8448804`.

[26] A. Elsheikh, S. Yacout, M.-S. Ouali, Bidirectional handshaking lstm for remaining useful life prediction, Neurocomputing 323 (2019) 148–156. `doi:https://doi.org/10.1016/j.neucom.2018.09.076`.
URL `https://www.sciencedirect.com/science/article/pii/S0925231218311573`

[27] A. Al-Dulaimi, S. Zabihi, A. Asif, A. Mohammadi, A multimodal and hybrid deep neural network model for remaining useful life estimation, Computers in Industry 108 (2019) 186–196. `doi:https://doi.org/10.1016/j.compind.2019.02.004`.
URL `https://www.sciencedirect.com/science/article/pii/S0166361518304925`

[28] J. Li, X. Li, D. He, A directed acyclic graph network combined with cnn and lstm for remaining useful life prediction, IEEE Access 7 (2019) 75464–75475. `doi:10.1109/ACCESS.2019.2919566`.

[29] J. Narwariya, P. Malhotra, V. TV, L. Vig, G. Shroff, Graph neural networks for leveraging industrial equipment structure: An application to remaining useful life estimation, arXiv preprint arXiv:2006.16556 (2020).

[30] A. Das, S. Hussain, F. Yang, M. S. Habibullah, A. Kumar, Deep recurrent architecture with attention for remaining useful life estimation, in: TENCON 2019 - 2019 IEEE Region 10 Conference (TENCON), 2019, pp. 2093–2098. `doi:10.1109/TENCON.2019.8929267`.

[31] J. Xia, Y. Feng, C. Lu, C. Fei, X. Xue, Lstm-based multi-layer self-attention method for remaining useful life estimation of mechanical systems, Engineering Failure Analysis 125 (2021) 105385.

[32] M. Ragab, Z. Chen, M. Wu, C.-K. Kwoh, R. Yan, X. Li, Attention-based sequence to sequence model for machine remaining useful life prediction, Neurocomputing 466 (2021) 58–68. `doi:https://doi.org/10.1016/j.neucom.2021.09.022`. URL `https://www.sciencedirect.com/science/article/pii/S0925231221013801`

[33] Y. Cao, Y. Ding, M. Jia, R. Tian, A novel temporal convolutional network with residual self-attention mechanism for remaining useful life prediction of rolling bearings, Reliability Engineering & System Safety 215 (2021) 107813.

[34] Y. Liu, X. Wang, Deep & attention: A self-attention based neural network for remaining useful lifetime predictions, in: 2021 7th International Conference on Mechatronics and Robotics Engineering (ICMRE), IEEE, 2021, pp. 98–105.

[35] D. Xu, H. Qiu, L. Gao, Z. Yang, D. Wang, A novel dual-stream self-attention neural network for remaining useful life estimation of mechanical systems, Reliability Engineering & System Safety 222 (2022) 108444. `doi:https://doi.org/10.1016/j.ress.2022.108444`. URL `https://www.sciencedirect.com/science/article/pii/S0951832022001090`

[36] H. Liu, Z. Liu, W. Jia, X. Lin, Remaining useful life prediction using a novel feature-attention-based end-to-end approach, IEEE Transactions on Industrial Informatics 17 (2) (2021) 1197–1207. `doi:10.1109/TII.2020.2983760`.

[37] Z. Chen, M. Wu, R. Zhao, F. Guretno, R. Yan, X. Li, Machine remaining useful life prediction via an attention-based deep learning approach, IEEE Transactions on Industrial Electronics 68 (3) (2021) 2521–2531. `doi:10.1109/TIE.2020.2972443`.

[38] Y. Song, S. Gao, Y. Li, L. Jia, Q. Li, F. Pang, Distributed attention-based temporal convolutional network for remaining useful life prediction, IEEE Internet of Things Journal (2020) 1–1`doi:10.1109/JIOT.2020.3004452`.

[39] Z. Zhang, W. Song, Q. Li, Dual aspect self-attention based on transformer for remaining useful life prediction, arXiv preprint arXiv:2106.15842 (2021).

[40] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need (2017). `arXiv:1706.03762`.

[41] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, Dropout: a simple way to prevent neural networks from overfitting, The journal of machine learning research 15 (1) (2014) 1929–1958.

[42] A. Saxena, K. Goebel, D. Simon, N. Eklund, Damage propagation modeling for aircraft engine run-to-failure simulation, in: 2008 International Conference on Prognostics and Health Management, 2008, pp. 1–9. `doi:10.1109/PHM.2008.4711414`.

[43] D. K. Frederick, J. A. DeCastro, J. S. Litt, User's guide for the commercial modular aero-propulsion system simulation (c-mapss), Tech. Rep. Technical Manual TM2007-215026, NASA/ARL (2007).

[44] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization (2017). `arXiv:1412.6980`.

[45] X. Glorot, A. Bordes, Y. Bengio, Deep sparse rectifier neural networks, in: G. Gordon, D. Dunson, M. Dudík (Eds.), Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, Vol. 15 of Proceedings of Machine Learning Research, PMLR, Fort Lauderdale, FL, USA, 2011, pp. 315–323.
URL http://proceedings.mlr.press/v15/glorot11a.html

[46] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), Advances in Neural Information Processing Systems, Vol. 30, Curran Associates, Inc., 2017.
URL https://proceedings.neurips.cc/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf