

# Machine Remaining Useful Life Prediction via an Attention Based Deep Learning Approach

Zhenghua Chen<sup>1</sup>, Min Wu<sup>1\*</sup>, Rui Zhao<sup>2</sup>, Feri Guretno<sup>1</sup>, Ruqiang Yan<sup>3\*</sup>, and Xiaoli Li<sup>1</sup>

**Abstract**—For prognostics and health management (PHM) of mechanical systems, a core task is to predict machine remaining useful life (RUL). Currently, deep structures with automatic feature learning, such as long short-term memory (LSTM), have achieved great performances for RUL prediction. However, the conventional LSTM network only uses the learned features at last time step for regression or classification, which is not efficient. Besides, some handcrafted features with domain knowledge may convey additional information for the prediction of RUL. It is thus highly motivated to integrate both those handcrafted features and automatically learned features for RUL prediction. In this paper, we propose an attention based deep learning framework for machine RUL prediction. The LSTM network is employed to learn sequential features from raw sensory data. Meanwhile, the proposed attention mechanism is able to learn the importance of features and time steps, and assign larger weights to more important ones. Moreover, a feature fusion framework is developed to combine the handcrafted features with automatically learned features to boost the performance of RUL prediction. Extensive experiments have been conducted on two real datasets and experimental results demonstrate that our proposed approach outperforms the state-of-the-arts.

**Index Terms**—Machine RUL prediction, PHM, LSTM, attention mechanism, feature fusion, handcrafted features

## I. INTRODUCTION

**R**EMAINING useful life (RUL) prediction is crucial for prognostics and health management (PHM) of mechanical systems. With accurate RUL prediction, maintenance schedules can be designed to keep good working conditions for machines (or components) and thus abrupt system failures can be avoided [1]. To achieve this objective, many advanced solutions have been developed, which generally can be divided into two categories, i.e., model-based and data-driven. For model-based solutions, they require to accurately model the dynamics of mechanical systems (or components) [2], [3]. However, due to the rapid development of industry, the mechanical systems become more and more complicated with complex interactions

This work is supported by the A\*STAR Industrial Internet of Things Research Program under the RIE2020 IAF-PP Grant A1788a0023, and partially supported by National Natural Science Foundation of China (No. 51835009). (Min Wu and Ruqiang Yan are the corresponding authors.)

<sup>1</sup> Zhenghua Chen, Min Wu, Feri Guretno and Xiaoli Li are with the Institute for Infocomm Research, A\*STAR, Singapore (Email: chen0832@e.ntu.edu.sg, wumin@i2r.a-star.edu.sg, guretnof@i2r.a-star.edu.sg, xlli@i2r.a-star.edu.sg).

<sup>2</sup> Rui Zhao is with harveston asset management company, Singapore (Email: rzhao001@e.ntu.edu.sg).

<sup>3</sup> Ruqiang Yan is with School of Mechanical Engineering, Xi'an Jiaotong University, China (Email: yanruqiang@xjtu.edu.cn).

between each other. Therefore, accurate modeling of these systems is not realistic, even for experts. Besides, the flexibility and transferability of model-based solutions are poor, due to the distinct mechanisms for different mechanical systems.

Recently, data-driven solutions have attracted more and more attention for RUL prediction [4]. For data-driven solutions, it is not compulsory to know the detailed operation mechanism of mechanical systems. Instead, one only needs to collect some data from the systems, where the conditions of the systems can be identified based on data-driven algorithms. Generally, data-driven solutions can be further divided into statistic degradation modeling and artificial intelligence [4].

A typical statistic degradation modeling approach is the Cox's regression [5]. It models the hazard probability of an object based on the historical data about the life span of objects and their associate covariates. The hazard probability at time step  $t$  can be expressed as

$$\lambda(t|\mathbf{z}) = \lambda_0(t) \exp(\beta^\top \mathbf{z}) \quad (1)$$

where  $\mathbf{z}$  is the covariates also known as features,  $\lambda_0(t)$  is the baseline hazard function which changes over time,  $\beta$  is the regression coefficients, and  $\top$  is the transpose operation. Many RUL prediction systems have been developed based on the Cox's regression. Pham et al. presented a RUL prediction system which combines the Cox's proportional hazard model with support vector machine (SVM) [6]. Liao et al. compared the Cox's regression model with logistic regression model for predicting the RUL of an individual unit [7]. In [8], You et al. proposed a two zone Cox's regression model for equipment RUL prediction.

For artificial intelligence based RUL prediction, the objective is to directly build a relationship between the RUL of an object and the features which can be extracted manually from sensor measurements or automatically learned by deep learning algorithms. With the rapid development of artificial intelligence [9], in this work, we mainly focus on artificial intelligence based RUL prediction, which includes the popular shallow machine learning and deep learning algorithms.

For shallow model based RUL prediction, it normally consists of feature extraction and inference. With domain knowledge on mechanical systems, some representative features can be extracted from raw sensory data which may be noisy and not representative for RUL prediction [10], [11]. After extracting some informative features, conventional machine learning algorithms, such as artificial neural network (ANN) [12], [13], extreme learning machine (ELM) [14], SVM [15], [16], neural networks [17] and random forest (RF) [18], can be employed to predict RUL.

Another popular algorithm for artificial intelligence based RUL prediction is deep learning. Instead of manual feature extraction, deep learning is able to automatically learn representative features from raw sensory data [19], [20]. Besides, it can jointly optimize feature learning and RUL inference, leading to a better generalization performance for RUL prediction. For machine RUL prediction, most of sensory data are time series with temporal dependency. The deep learning approach of long short-term memory (LSTM) which is designed for sequential data analytics can be naturally suitable for RUL prediction. A very good performance has been achieved for RUL prediction using the LSTM approach [21], [22]. However, it still has several limitations for the conventional LSTM in RUL prediction, which are presented as follows:

- 1) For the conventional LSTM network, it only uses the learned features at last time step for regression or classification. We argue that the learned features at other time steps may also have some contribution. And the learned features may have different contribution for RUL prediction. Therefore, an efficient operation is to assign larger weights to more important features and time steps.
- 2) Although the features learned by the LSTM network have been shown to be effective for RUL prediction, some handcrafted features with domain knowledge also convey important information. The design of a network that can take both automatically learned features from the LSTM and some meaningful handcrafted features into consideration for RUL prediction may boost the performance of RUL prediction.

To solve these limitations, we propose an attention based deep learning framework for RUL prediction. The proposed approach firstly exploits the LSTM network to learn representative sequential features from raw sensory data. Then, an attention network is developed to learn the importance of features and time steps, and assign larger weights to more important ones. Finally, we propose a feature fusion framework to make full use of all available information by combining the automatically learned features and some useful handcrafted features for the RUL prediction of machine RUL. To verify the effectiveness of the proposed approach for RUL prediction, we use real datasets for evaluation and compare it with various state-of-the-art methodologies for machine RUL prediction. The main contributions of this paper are summarized as follows:

We propose an attention based deep learning framework for machine RUL prediction. The attention network is able to learn the importance of features and time steps and assign larger weights for more important ones to boost the performance of RUL prediction.

Since some handcrafted features with domain knowledge may convey meaningful information for RUL prediction, we propose a feature fusion framework to combine both automatically learned features and some meaningful handcrafted features for RUL prediction of machine RUL.

Real datasets are leveraged to evaluate the performance of the proposed approach. The results show that the proposed approach can significantly improve the performance

of RUL prediction.

## II. RELATED WORKS

In this section, we review some advanced deep learning algorithms for machine RUL prediction.

Deep learning is able to learn representative features and perform inference simultaneously, resulting remarkable performance for RUL prediction. Babu et al. proposed a deep convolutional neural network (CNN) for RUL prediction [23].

The experimental results on two datasets indicated its superior performance for RUL prediction when compared with some shallow learning algorithms. Zhu et al. proposed a multi-scale convolutional neural network (MSCNN) to predict RUL [24].

Firstly, a wavelet transform was conducted on raw sensory data to get time frequency representation (TFR). Then, the TFR was used as the input of MSCNN for the prediction of RUL.

Deutsch and He presented a deep belief network feedforward neural network (DBN-FNN) for RUL prediction [25]. The DBN was used to learn representative features, and the FNN was employed to perform RUL prediction with the learned features.

Zheng et al. presented a long short-term memory (LSTM) based RUL prediction. The experimental results on three datasets showed that the LSTM performs much better than some shallow learning algorithms and CNN [21]. Zhang et al. proposed a bi-directional long short-term memory (BD-LSTM) approach to predict RUL [22]. They firstly defined a

network based on a perceptron approach. The BD-LSTM was then utilized to track the variation of the HI for RUL prediction. In [26], the authors proposed a multiobjective deep belief networks ensemble (MODBNE) approach for the prediction of RUL. They applied a multiobjective evolutionary algorithm to train DBNs with two conflict objectives, i.e., accuracy and diversity. The evolved DBNs were then combined to form an ensemble model for RUL prediction.

Since the sensory data for PHM are time series with temporal dependency, the LSTM network which performs well for sequential data modeling is naturally suitable for RUL prediction. However, different features at different time steps learned by the LSTM network will have equal contribution

for the RUL prediction, which is not effective. A more effective operation is to assign larger weights for more important features and time steps. Therefore, in this work, we

propose an attention mechanism with LSTM to automatically assign larger weights to more significant features and time steps to boost the performance of RUL prediction. Meanwhile,

some handcrafted features may convey useful information for the prediction of RUL. Hence, we develop a feature fusion framework to combine these handcrafted features with the features learned by the attention based LSTM to further improve the performance of machine RUL prediction.

## III. METHODOLOGY

### A. Long Short-Term Memory

To predict the RUL of machines, a number of sensors, such as vibration, temperature, acoustic, etc., should be deployed.

Generally, the sensor measurements are time series with temporal dependency. Recurrent neural network (RNN) whose

nodes are connected along a sequence was designed to model contribution for nal RUL prediction. Therefore, we intend temporal dependency in time series [27]. Therefore, RNN design an attention mechanism to learn the importance of is naturally suitable for machine RUL prediction leveraging features and time steps. The details will be introduced in the sequential sensor measurements. However, the conventional following paragraphs.

RNN often suffers from the problem of gradient vanishing or exploding during network training, which greatly degrades its performance on modeling long-term dependencies [28]. To solve this problem, Hochreiter and Schmidhuber proposed a new architecture, named long short-term memory (LSTM), which can be treated as a memory cell that consists of a few gates [29]. The gates, which can allow or prevent the passing of information along a sequence, can capture long-term dependencies. Owing to its unique property, the LSTM network has achieved great successes in the analysis of time series data, such as occupancy estimation [30], video analysis [31] and nature language processing [32]. Recently, it also has achieved great performance for RUL prediction [21], [22], [33].

Fig. 1. The structure of LSTM

A typical LSTM network is shown in Fig. 1. It consists of a forget gate to discard the unnecessary information from previous time steps, an input gate to select useful information from inputs, and an output gate to control the outputs of the current LSTM network. Assuming that  $x^t$  is the input at time step  $t$ ,  $h^t$  is the hidden state at time step  $t$ ,  $C^{t-1}$  is the memory cell state,  $w^f$ ,  $w^i$ ,  $w^c$  and  $w^o$  are the weights,  $b^f$ ,  $b^i$ ,  $b^c$  and  $b^o$  are the biases, and  $\sigma(\cdot)$  and  $\tanh(\cdot)$  are the sigmoid and tanh functions respectively, the LSTM network can be expressed as

$$\begin{aligned}
 f^t &= \sigma(w^f [h^{t-1}; x^t] + b^f); \\
 i^t &= \sigma(w^i [h^{t-1}; x^t] + b^i); \\
 C^t &= \tanh(w^c [h^{t-1}; x^t] + b^c); \\
 C^t &= f^t \odot C^{t-1} + i^t \odot C^t; \\
 o^t &= \sigma(w^o [h^{t-1}; x^t] + b^o); \\
 h^t &= o^t \odot \tanh(C^t);
 \end{aligned}
 \tag{2}$$

Due to the strong sequential modeling ability of the LSTM network, it has been successfully used for machine RUL prediction in [21], [22]. However, they applied the standard LSTM which only uses the learned features at last time step for regression, which is shown in Fig. 2. We argue that the learned features at other time steps may also have some contribution. And the learned features may have different

Fig. 2. The standard LSTM for regression problems.

### B. Attention mechanism

The attention mechanism is firstly proposed for the task of image processing, inspired by human vision system [34], [35]. Human always pays attention to a certain region of an image during recognition, indicating that different weights will be assigned to different regions of an image. The attention mechanism has been successfully applied for a number of applications, such as language translation [36] and time series prediction [37].

For the task of machine RUL prediction, an efficient operation is to focus on different region of interest by assigning different weights for different features at different time steps. In this task, since no prior information is available, we leverage a self-attention mechanism to learn the importance of features and time steps. Assume that the learned features

by the LSTM network for one sample can be expressed as  $H = [h_1; h_2; \dots; h_d]^T$ ,  $\odot$  is the transpose operation. Here,  $h_i \in \mathbb{R}^n$ , where  $n$  is the number of sequential steps of the features. Based on the self-attention mechanism, the importance for different sequential steps with input  $h_i$  can be expressed as

$$s_i = (\mathbf{W}^T h_i + b); \tag{3}$$

where  $\mathbf{W}$  and  $b$  are the weight matrix and the bias vector, respectively,  $\sigma(\cdot)$  is the score function which can be designed as an activation function in neural networks, such as sigmoid and linear. After obtaining the score for  $i$ -th feature vector, it can be normalized using softmax function as follows:

$$a_i = \text{softmax}(s_i) = \frac{\exp(s_i)}{\sum_j \exp(s_j)}; \tag{4}$$

The final output feature  $O$  of the attention mechanism can be expressed as

$$O = H \odot A; \tag{5}$$

where  $A = [a_1; a_2; \dots; a_d]^T$ , and  $\odot$  is a new operation defined as element-wise multiplication. Given vectors  $a = [a_1; a_2; \dots; a_n]^T$  and  $c = [c_1; c_2; \dots; c_n]^T$ ,  $b \odot c = [b_1 c_1; b_2 c_2; \dots; b_n c_n]^T$ .

C. Attention based deep learning for RUL prediction

1) Handcrafted features For sensory data based RUL prediction, some intuitive handcrafted features can be extracted, such as mean and trend coefficient of linear regression. The mean value shows the magnitude of sensory data, and the trend coefficient indicates the degradation of sensory data. These two simple handcrafted features have been shown to be effective for RUL prediction in [16]. An example of these two features is shown in Fig. 3. Note that, the features have been standardized for normalization. It can be found that these two features well indicate the properties of the raw sensory data.

To make full use of all the available information, we propose a feature fusion framework to combine the features learned by deep structures with some meaningful handcrafted features to boost the performance of RUL prediction. The details of the proposed approach will be shown in the following paragraphs.

Fig. 3. An example for the features of mean and trend coefficient.

2) Proposed framework: Fig. 4 shows the proposed attention based deep learning framework for the prediction of machine RUL. Firstly, the raw sensory data are fed into the LSTM network for feature learning. The learned sequential features are treated as the inputs of the attention model, whose outputs (attention weights) indicate the importance of features and time steps. Then, the learned sequential features are merged with the weights generated by the attention model. After that, two fully connected layers (FCLs) are performed to obtain more abstract features. Meanwhile, the handcrafted features are extracted from the raw sensory data, and then fed into a FCL to obtain more abstract features. To make full use of these two types of features, we concatenate them to form a complete feature set. Finally, a regression layer is used for RUL prediction. Table I summarizes the inputs of some key modules, including the LSTM network, the attention layer, the merge layer and the fully connected layer (i.e., the FCL at right-hand side in Fig. 4).

Since the prediction of RUL is a typical regression problem, the loss function of the proposed approach is set to be the mean square error (MSE) loss which is defined as  $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ . Given the predicted RULs and the true RULs, the MSE losses over training data can be calculated and backpropagated to generate the error gradients for each layer (see Fig. 4). The masking probability for dropout is set to 0.5.

Fig. 4. The proposed attention based deep learning approach for RUL prediction.

TABLE I  
INPUTS OF SOME KEY MODULES

Module	Input
LSTM network	Raw sensory data
Attention layer	Features learned by the LSTM
Merge layer	Features learned by the LSTM; Weights generated by the attention layer
Fully connected layer (Right-hand side in Fig. 4)	Handcrafted features

IV. EVALUATION

A. Data description

The widely used C-MAPSS (Commercial Modular Aero Propulsion System Simulation) dataset [40] is adopted for the evaluation of the proposed approach. This dataset describes the degradation process of aircraft engine whose diagram is shown in Fig. 5. The engine consists of fan, low pressure compressor (LPC), high pressure compressor (HPC), combustor, low pressure turbine (LPT) and high pressure turbine (HPT). Twenty-one on-board sensors, measuring temperature, pressure and speed, are deployed at different locations to monitor the condition of the engine.

Fig. 5. The diagram of aircraft engine [40].

The entire dataset contains four sub-datasets with varying number of operation conditions and faulty types. We choose two typical sub-datasets: the simplest one FD001 which contains one operation condition and one faulty type and the most complicated one FD004 which contains six operation conditions and two faulty types. For FD001 and FD004, they both contain two files for training and testing. The training file records sensor data at each running cycle in the run-to-fail experiments for certain number of engines. The testing file only contains the sensor measurements to certain running cycles for another certain number of engines. The objective is to predict the RUL of each engine in the testing file with the given sensor measurements.

Another widely used dataset for evaluation is the PHM 2008 dataset [41], which has the same data structure with the C-MAPSS dataset, but different number of training and testing engines. The description of the datasets is shown in Table II.

TABLE II  
THE DESCRIPTION OF THE DATASETS FOR EVALUATION [40], [41].

Dataset	C-MAPSS		PHM 2008
	FD001	FD004	
No. of training engines	100	249	218
No. of testing engines	100	248	218
Working conditions	1	6	6
Faulty types	1	2	2

For the twenty-one sensors (indices from 1 to 21 in training and testing files), the sensors with indices 1, 5, 6, 10, 18, and 19 always have constant values during the run-to-fail experiments. This means that these sensors are not related to the degradation of engines. Hence, these sensors are removed from the two datasets as did in [26], [42]. Finally, fourteen sensors are used for the RUL prediction. Since different operating conditions will influence RUL, we treat operating conditions as measurement signals for RUL prediction. Therefore, the operating conditions and sensor measurements are the inputs of prediction models.

B. Data preprocessing

A sliding window is commonly used for data segmentation [23], [26]. An example of data segmentation for training samples is shown in Fig. 6. For the run-to-fail experiments, we assume that the number of total running cycles of an engine is  $T$ , the window size is  $w$  and the step size is  $s$ . Each sample will have a size of  $fs$ , where  $n$  is the number of sensors. According to Fig. 6, the RUL of the  $(i+1)$ -th sample (window) is  $T - s \cdot i$ . Note that, a piece-wise linear RUL [23], [26] is used instead of the true RUL for training data, which means that if the true RUL is larger than the maximal RUL, it will be set to the maximal RUL. Here, we choose a window size of 30 and a step size of 1 which are the same as these in [26]. Under these settings, the numbers of training samples for FD001, FD004 and PHM 2008 are 17731, 54028 and 39596 respectively. For testing, only one data window to the last sensor measurements for each engine is used as the testing sample. Hence, the numbers of testing samples for FD001, FD004 and PHM 2008 are 100, 248 and 218 respectively, which are the numbers of testing engines in the three datasets. Note that, the actual RUL values of testing samples for FD001 and FD004 are available to the public, while the actual RUL values of testing samples for PHM 2008 are not available.

C. Evaluation criteria

To evaluate the performance of RUL prediction, two widely used evaluation criteria, i.e., root mean square error (RMSE) and scoring function, are adopted [23], [26]. The definition of the RMSE can be found in Appendix. In machine RUL prediction, the late prediction refers to that the predicted RUL is larger than the actual RUL. Late predictions will generally lead to more severe consequences than the early prediction, however, this fact cannot be reflected by the criterion of RMSE. Thus, the scoring function [23] (See the Appendix for the definition) can be utilized. Based on the definition, more penalties will be given to late predictions, which is consistent with our common sense. Both criteria are useful to evaluate the performance of RUL prediction. Fig. 7 compares the RMSE with the scoring function. Both criteria are useful to evaluate the performance of RUL prediction. Our experimental results will be reported based on these two criteria.

D. Experimental setup

To verify the effectiveness of the proposed approach, we firstly perform an initial test on the training data of FD001 which is split for training and testing. A comparison has been made between the proposed approach and some widely used benchmark approaches including SVR, RF, CNN and LSTM. We have also compared the proposed approach with some state-of-the-art approaches on the testing data of FD001, FD004 and PHM 2008 datasets. Note that, the state-of-the-arts used the same data for training and testing. For the proposed attention based deep learning approach, some structural parameters, i.e., the number of hidden nodes, should be tuned based on the given training data. A cross-validation is performed with the training data to determine

Fig. 6. An example of data segmentation for RUL prediction.

TABLE III  
TRAINING AND TESTING TIME OF THE BENCHMARK APPROACHES AND  
THE PROPOSED APPROACH.

	SVR	RF	CNN	LSTM	Proposed
Training time (s)	2.43	6.80	27.77	81.75	110.15
Testing time (s)	0.45	0.010	0.11	0.31	0.42

Fig. 7. RMSE vs scoring function [23].

the parameters of the proposed approach. Specifically, the number of hidden units for the LSTM network is 50. The hidden nodes of two FCLs after the merge layer are set to be [50, 10]. And the hidden nodes of the FCL on the handcrafted features is set to be 10. The learning rate for the optimization algorithm of Adam is set to be 0.001. For the algorithms of SVR, RF and CNN, we use the parameters in [23] and [26], which have carefully tuned the parameters of the models. For the algorithm of conventional LSTM, the parameters have been carefully tuned by using cross-validation on the training data. Specifically, the number of hidden nodes for the LSTM network is set to be 50. Two FCLs with the sizes of [50, 10] are applied for RUL prediction. Considering the randomness in parameter initialization of some algorithms, i.e., CNN, LSTM and the proposed approach, these algorithms are run 10 times for each dataset and average results are reported. The source code of the proposed approach is available at the open-source GitHub (<https://github.com/ZhenghuaNTU/RUL-prediction-using-attention-based-deep-learning-approach>).

#### E. Initial test on FD001

In order to have more samples for testing, an initial test is firstly performed with the training data in FD001. By using sliding windows, 17731 samples can be obtained. Here, we use first 70% of data (12412 samples) for model training and the rest (5319 samples) for testing. A comparison

has been made with some widely used benchmark algorithms for RUL prediction, including SVR, RF, CNN and LSTM. The results are shown in Fig. 8. We also indicate the 95 percent confidence interval of all the results in the figure. It can be found that, due to the powerful feature learning ability of deep structures, deep learning based algorithms perform better than shallow models, i.e., SVR and RF. Due to the sequential modeling ability of LSTM, it has a superior performance over the CNN. Owing to the proposed attention mechanism and the feature fusion framework, the proposed approach outperforms all the benchmark approaches.

We also present the training and testing time for different algorithms on a workstation which has twelve core CPUs of Intel i7-8700 3.20GHz and a GPU of NVIDIA GeForce GTX1080Ti. The results are shown in Table III. Obviously, deep learning based approaches require much more training and testing time than shallow models, because of much more parameters to be optimized. The proposed approach has the longest training time which is 110.15 seconds. Considering that the model training only requires to be done once and it is off-line, this amount of time for training is still acceptable. The testing time of the proposed method for all the 5319 samples is only 0.42 seconds, which means that the testing time for one sample is  $7.8 \times 10^{-5}$  that can be neglected for most of industrial applications. Hence, it can be claimed that the proposed method can be used for real-time implementations. To better interpret the attention mechanism for RUL prediction, an attention matrix of one testing sample is illustrated in Fig. 9. Since we set 50 hidden nodes for the LSTM network and the sliding window contains 50 time steps, the learned features by the LSTM will have a dimension of 50. In the conventional LSTM without the attention mechanism, only the learned features at last time step of a sliding window will be used for classification or regression problems. We argue that the learned features at other time steps may also have some contribution for RUL prediction. From Fig. 9,

(a) RMSE (b) Score

Fig. 8. The experimental results of all the benchmark approaches and the proposed approach in the initial test on FD001.

(a) FD001: RMSE (b) FD001: Score

Fig. 9. Attention matrix of one sample.

it can be found that more recent time steps will have larger attention weights, indicating that more recent steps will be more important for RUL prediction. This is consistent with our common sense. Although the learned features at each time step cannot be explicitly interpreted because they are high-level features learned by the LSTM network, we can still concur that they will have different contribution for RUL prediction, which has been achieved by using the attention mechanism.

F. Results on the testing data of FD001 and FD004

Here, we firstly analyze the impact of window size on the performance of RUL prediction. Then, an ablation study is performed to demonstrate the effectiveness of the proposed attention mechanism and feature fusion. Finally, we compare the proposed approach with some state-of-the-art methods on RUL prediction which use the same data for model training and testing.

1) Analysis on window size For data preprocessing, window size is one of the most important parameters. To evaluate the impact of this parameter, we implement the proposed approach with different window sizes, i.e., 10, 20, 30, 40, 50, and 60, on the two datasets. The results are shown in

(c) FD004: RMSE (d) FD004: Score

Fig. 10. The experimental results of the proposed approach with different window sizes on the two datasets.

Fig. 10. For the simple FD001 dataset, it can be found that the performance is enhanced when increasing window size at the beginning. This is because more information is included for RUL prediction. However, after increasing the window size over a certain value, i.e., 30, the performance of model degrades on this dataset. This may be caused by the overfitting of the algorithm with too much information for this simple dataset. For the much more complicated FD004 dataset, similarly, the performance of the model is enhanced when the number of hidden nodes increases from 10 to 30. After 30, the performance slightly degrades, which may be because the other model parameters are tuned under 30, resulting a slight degradation of the performance with different number of hidden nodes. Moreover, we interestingly find that when we further increase the number of hidden nodes, the performance in the model improves a lot on this complicated dataset.

This means that the complicated dataset may require more information for accurate RUL prediction. The performance of the model under the window sizes of 50 and 60 is similar, which should be caused by the limited modeling capacity of the model with large amount of information. Therefore, we may be able to further enhance the performance of the model by using more hidden nodes or deeper structure.

In summary, it can be concluded that more complicated datasets require a larger window size to include more information for accurate RUL prediction. In real applications, we should choose different window sizes for different datasets by using cross-validation on the training data. To give a fair comparison with state-of-the-arts which used a sliding window of 30 in their works [22], [26], a window size of 30 is chosen in this work.

2) Ablation study of the proposed approach. To evaluate the effectiveness of the proposed attention mechanism and feature fusion, an ablation study of the proposed approach is performed. Specifically, we implement the original LSTM, LSTM with attention, the LSTM plus handcrafted features and the proposed framework. The results can be found in Table 11 and Fig. 11.

(a) FD001: RMSE

(b) FD001: Score

(c) FD004: RMSE

(d) FD004: Score

Fig. 11. The experimental results of the ablation study.

Generally, the LSTM with attention and the LSTM plus feature fusion outperform the original LSTM, which indicates the effectiveness of the two proposed schemes. Moreover, it can be found that the LSTM plus feature fusion has a superior performance than the LSTM with attention mechanism in terms of RMSE, Score and their variances. This means that the proposed feature fusion is more effective than the proposed attention mechanism to enhance the performance of the original LSTM for RUL prediction. In other words, the handcrafted features well compensate the automatically learned features for the task of RUL prediction. The proposed approach with these two effective schemes achieves the best performance on the two datasets in terms of error and variance under the two criteria.

3) Compare to state-of-the-arts. The experimental results of the proposed approach and some state-of-the-art approaches on the two datasets are shown in Table V. Overall, all the approaches perform better on the FD001 than that on the FD004. This is because the FD001 is relatively simple with only one operation condition and one faulty type. Besides, the number of engines for testing in the FD004 is 248 which is much larger than that in the FD001. Therefore, the scores which are summations over all the engines of the FD004 and the FD001 are under different magnitude. According to Table V, deep learning based methods, such as LSTM, DBN and MODBNE, perform better than statistic learning algorithm of Cox's regression and shallow learning algorithms, i.e., MLP, SVR, RVR, ELM and RF. This indicates the powerful feature learning ability of deep structures. The MODBNE which is an ensemble of DBN outperforms the other benchmark approaches, due to the feature learning of DBN and the ensemble structure. However, the ensemble structure of DBN will have a much higher computational complexity than other shallow and deep structures. Owing to the proposed attention mechanism and feature fusion framework, the proposed approach outperforms all these state-of-the-art approaches, including the powerful ensemble deep learning approach of MODBNE.

The predicted RUL on the two datasets is shown in Fig. 12. For both datasets, the predicted RUL matches very well with the true RUL, which indicates the feasibility of RUL prediction. Since FD001 is relatively simple, the prediction performance on FD001 is better than that on FD004, which is consistent with our previous analysis based on Table V.

(a) FD001

(b) FD004

Fig. 12. The true RUL and the predicted RUL by the proposed approach on the two datasets.

### C. Results on the PHM 2008 dataset

For the PHM 2008 dataset, the actual RUL values of testing samples are not available. The results need to be uploaded to the NASA Data Repository website (<https://ti.arc.nasa.gov/tech/dash/groups/pcoe/prognostic-data-repository/>) where a score value (See Equation 8) will be given.

TABLE IV  
THE RESULTS OF THE ABLATION STUDY

Methods	FD001		FD004	
	RMSE	Score	RMSE	Score
LSTM + 2 FCs + Regression (Ori-LSTM)	15.43	410.60	29.12	12551.44
LSTM + Attention + 2 FCs + Regression (LSTM-Attention)	15.27	362.92	28.04	7850.40
[LSTM + 2 FCs, Handcrafted Feat + FC] + Regression (LSTM+HF)	14.81	360.72	27.66	7812.37
[LSTM + Attention + 2 FCs, Handcrafted Feat + FC] + Regression (Proposed)	14.54	322.44	27.08	5649.14

TABLE V  
THE EXPERIMENTAL RESULTS ON THE FD001 AND FD004 DATASETS.

	Criterion	Statistic Cox's regression	Shallow learning					Deep learning				Hybrid
			MLP [23]	SVR [23]	RVR [23]	ELM [26]	RF [26]	CNN [23]	LSTM [21], [22]	DBN [26]	MODBNE [26]	Proposed
FD001	RMSE	45.10	37.56	20.96	23.80	17.27	17.91	18.45	15.42	15.21	15.04	14.53
	Score	28616	17972	1381.5	1502.9	523.00	479.75	1286.7	410.60	417.59	334.23	322.44
FD004	RMSE	54.29	77.37	45.35	34.34	38.43	31.12	29.16	29.12	29.88	28.66	27.08
	Score	1164590	5616600	371140	26509	121414.47	46567.63	7886.4	12551.44	7954.51	6557.62	5649.14

The results on the PHM 2008 dataset are shown in Table VI. For machine learning based RUL prediction, the basic assumption is that the underline patterns between training and testing data are the same. However, if the training and testing data are collected under different environments, working conditions, or machines, the underline patterns between training and testing data may be distinct, which will hinder the performance of machine learning based methods [43]. To solve this issue, transfer learning which is able to transfer the knowledge learned from one domain to another [44] can be adopted. In our future works, we will investigate transfer learning based RUL prediction with varying environments, working conditions, or machines.

TABLE VI  
THE EXPERIMENTAL RESULTS ON THE PHM 2008 DATASET.

Method	Score
Cox's regression	65558
MLP [23]	3212
SVR [23]	15886
RVR [23]	8242
CNN [23]	2056
LSTM [21], [22]	1862
Proposed	1584

V. CONCLUSION AND FUTURE WORKS

In this paper, we proposed an attention based deep learning framework for machine RUL prediction. Firstly, we employ the deep learning algorithm of long short-term memory (LSTM) for automatic feature learning from raw sensory data. Then, an attention mechanism was proposed to learn the importance of features and time steps, and automatically assign larger weights to more important ones. Meanwhile, some handcrafted features with domain knowledge may convey additional information for RUL prediction. Hence, a feature fusion framework was designed to combine the handcrafted features and the automatically learned features to boost the performance of RUL prediction. The proposed approach was evaluated using real datasets. Since the size of sliding window is important for RUL prediction, the impact of different window size on prediction performance was investigated. Then, we verified the effectiveness of the proposed attention mechanism and feature fusion for machine RUL prediction. Finally, a comparison has been made with various state-of-the-art approaches. The proposed approach outperforms these state-of-the-arts under two popular evaluation criteria.

REFERENCES

- [1] L. Liao, "Discovering prognostic features using genetic programming in remaining useful life prediction," *IEEE Transactions on Industrial Electronics* vol. 61, no. 5, pp. 2464–2472, 2014.
- [2] Y. Qian, R. Yan, and R. X. Gao, "A multi-time scale approach to remaining useful life prediction in rolling bearing," *Mechanical Systems and Signal Processing* vol. 83, pp. 549–567, 2017.
- [3] Q. Zhai and Z.-S. Ye, "Rul prediction of deteriorating products using an adaptive wiener process model," *IEEE Transactions on Industrial Informatics* vol. 13, no. 6, pp. 2911–2921, 2017.
- [4] Y. Lei, N. Li, L. Guo, N. Li, T. Yan, and J. Lin, "Machinery health prognostics: A systematic review from data acquisition to rul prediction," *Mechanical Systems and Signal Processing* vol. 104, pp. 799–834, 2018.
- [5] C. R. Davidet al., "Regression models and life tables (with discussion)," *Journal of the Royal Statistical Society* vol. 34, no. 2, pp. 187–220, 1972.
- [6] H. T. Pham, B.-S. Yang, T. T. Nguyen et al., "Machine performance degradation assessment and remaining useful life prediction using proportional hazard model and support vector machine," *Mechanical Systems and Signal Processing* vol. 32, pp. 320–330, 2012.
- [7] H. Liao, W. Zhao, and H. Guo, "Predicting remaining useful life of an individual unit using proportional hazards model and logistic regression model," in *RAMS'06. Annual Reliability and Maintainability Symposium, 2006*, pp. 127–132. IEEE, 2006.
- [8] M.-Y. You, L. Li, G. Meng, and J. Ni, "Two-zone proportional hazard model for equipment remaining useful life prediction," *Journal of manufacturing science and engineering* vol. 132, no. 4, p. 041008, 2010.
- [9] X.-S. Si, W. Wang, C.-H. Hu, and D.-H. Zhou, "Remaining useful life estimation—a review on the statistical data driven approaches," *European Journal of Operational Research* vol. 213, no. 1, pp. 1–14, 2011.
- [10] T. H. Loutas, D. Roulias, and G. Georgoulas, "Remaining useful life estimation in rolling bearings utilizing data-driven probabilistic e-support vectors regression," *IEEE Transactions on Reliability* vol. 62, no. 4, pp. 821–832, 2013.

[11] T. Benkedjouh, K. Medjaher, N. Zerhouni, and S. Rechak, "Remaining useful life estimation based on nonlinear feature reduction and support vector regression," *Engineering Applications of Artificial Intelligence* vol. 26, no. 7, pp. 1751–1760, 2013.

[12] N. Gebraeel, M. Lawley, R. Liu, and V. Parmeshwaran, "Residual life predictions from vibration-based degradation signals: a neural network approach," *IEEE Transactions on Industrial Electronics* vol. 51, no. 3, pp. 694–700, 2004.

[13] Z. Tian, "An artificial neural network method for remaining useful life prediction of equipment subject to condition monitoring," *Journal of Intelligent Manufacturing* vol. 23, no. 2, pp. 227–237, 2012.

[14] K. Javed, R. Gouriveau, and N. Zerhouni, "A new multivariate approach for prognostics based on extreme learning machine and fuzzy clustering," *IEEE Transactions on Cybernetics* vol. 45, no. 12, pp. 2626–2639, 2015.

[15] P. G. Nieto, E. Garcia-Gonzalo, F. S. Lasheras, and F. J. de Cos Juez, "Hybrid psvm-based method for forecasting of the remaining useful life for aircraft engines and evaluation of its reliability," *Reliability Engineering & System Safety* vol. 138, pp. 219–231, 2015.

[16] R. Khelif, B. Chebel-Morello, S. Malinowski, E. Laajili, F. Fnaiech, and N. Zerhouni, "Direct remaining useful life estimation based on support vector regression," *IEEE Trans. Industrial Electronics* vol. 64, no. 3, pp. 2276–2285, 2017.

[17] F. Yang, M. S. Habibullah, T. Zhang, Z. Xu, P. Lim, and S. Nadarajan, "Health index-based prognostics for remaining useful life predictions in electrical machines," *IEEE Transactions on Industrial Electronics* vol. 63, no. 4, pp. 2633–2644, 2016.

[18] D. Wu, C. Jennings, J. Terpenney, R. X. Gao, and S. Kumara, "A comparative study on machine learning algorithms for smart manufacturing: tool wear prediction using random forests," *Journal of Manufacturing Science and Engineering* vol. 139, no. 7, p. 071018, 2017.

[19] L. Liao, W. Jin, and R. Pavel, "Enhanced restricted boltzmann machine with prognosability regularization for prognostics and health assessment," *IEEE Transactions on Industrial Electronics* vol. 63, no. 11, pp. 7076–7083, 2016.

[20] R. Zhao, R. Yan, Z. Chen, K. Mao, P. Wang, and R. X. Gao, "Deep learning and its applications to machine health monitoring," *Mechanical Systems and Signal Processing* vol. 115, pp. 213–237, 2019.

[21] S. Zheng, K. Ristovski, A. Farahat, and C. Gupta, "Long short-term memory network for remaining useful life estimation," *Prognostics and Health Management (ICPHM), 2017 IEEE International Conference on*, pp. 88–95. IEEE, 2017.

[22] J. Zhang, P. Wang, R. Yan, and R. X. Gao, "Long short-term memory for machine remaining life prediction," *Journal of manufacturing systems* vol. 48, pp. 78–86, 2018.

[23] G. S. Babu, P. Zhao, and X.-L. Li, "Deep convolutional neural network based regression approach for estimation of remaining useful life," in *International conference on database systems for advanced applications* pp. 214–228. Springer, 2016.

[24] J. Zhu, N. Chen, and W. Peng, "Estimation of bearing remaining useful life based on multiscale convolutional neural network," *IEEE Transactions on Industrial Electronics* vol. 65, no. 2, pp. 1539–1548, 2018.

[25] J. Deutsch and D. He, "Using deep learning-based approach to predict remaining useful life of rotating components," *IEEE Transactions on Systems, Man, and Cybernetics: Systems* vol. 48, no. 1, pp. 11–20, 2018.

[26] C. Zhang, P. Lim, A. Qin, and K. C. Tan, "Multiobjective deep belief networks ensemble for remaining useful life estimation in prognostics," *IEEE Transactions on Neural Networks and Learning Systems* vol. 28, no. 10, pp. 2306–2318, 2017.

[27] J. A. Bullinaria, "Recurrent neural networks," *Neural Computation: Lecture* vol. 12, 2013.

[28] R. Zhao, D. Wang, R. Yan, K. Mao, F. Shen, and J. Wang, "Machine health monitoring using local feature-based gated recurrent unit networks," *IEEE Transactions on Industrial Electronics* vol. 65, no. 2, pp. 1539–1548, 2018.

[29] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation* vol. 9, no. 8, pp. 1735–1780, 1997.

[30] Z. Chen, R. Zhao, Q. Zhu, M. K. Masood, Y. C. Soh, and K. Mao, "Building occupancy estimation with environmental sensors via cdblstm," *IEEE Transactions on Industrial Electronics* vol. 64, no. 12, pp. 9549–9559, 2017.

[31] K. Zhang, W.-L. Chao, F. Sha, and K. Grauman, "Video summarization with long short-term memory," in *European Conference on Computer Vision*, pp. 766–782. Springer, 2016.

[32] H. Palangi, L. Deng, Y. Shen, J. Gao, X. He, J. Chen, X. Song, and R. Ward, "Deep sentence embedding using long short-term memory

networks: Analysis and application to information retrieval," *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)* vol. 24, no. 4, pp. 694–707, 2016.

[33] Y. Cheng, H. Zhu, J. Wu, and X. Shao, "Machine health monitoring using adaptive kernel spectral clustering and deep long short-term memory recurrent neural networks," *IEEE Transactions on Industrial Informatics* 2018.

[34] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *International Conference on Machine Learning* pp. 2048–2057, 2015.

[35] W. Du, Y. Wang, and Y. Qiao, "Recurrent spatial-temporal attention network for action recognition in videos," *IEEE Transactions on Image Processing* vol. 27, no. 3, pp. 1347–1360, 2018.

[36] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473* 2014.

[37] Y. Wu, H. Tan, L. Qin, B. Ran, and Z. Jiang, "A hybrid deep learning based traffic flow prediction method and its understanding," *Transportation Research Part C: Emerging Technologies* vol. 90, pp. 166–180, 2018.

[38] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980* 2014.

[39] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhudinov, "Improving neural networks by preventing co-adaptation of feature detectors," *arXiv preprint arXiv:1207.0580* 2012.

[40] A. Saxena, K. Goebel, D. Simon, and N. Eklund, "Damage propagation modeling for aircraft engine run-to-failure simulation," *Prognostics and Health Management, 2008. PHM 2008. International Conference on*, pp. 1–9. IEEE, 2008.

[41] A. Saxena and K. Goebel, "PHM08 challenge data set. NASA AMES prognostics data repository," *Moffett Field, CA, Tech. Rep.* 2008.

[42] P. Lim, C. K. Goh, K. C. Tan, and P. Dutta, "Estimation of remaining useful life based on switching kalman filter neural network ensemble," in *Prognostics and Health Management, 2014. PHM 2014. International Conference on* pp. 2–9. IEEE, 2014.

[43] R. Yan, F. Shen, C. Sun, and X. Chen, "Knowledge transfer for rotary machine fault diagnosis," *IEEE Sensors Journal* 2019.

[44] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on Knowledge and Data Engineering* vol. 22, no. 10, pp. 1345–1359, 2010.

## APPENDIX

Loss Function: The MSE loss is defined as follows:

$$\text{MSE Loss} = \frac{1}{N} \sum_{i=1}^N (\hat{r}_i - r_i)^2; \quad (6)$$

where  $r_i$  and  $\hat{r}_i$  are the true RUL and the predicted RUL, respectively, and  $N$  is the total number of samples.

Evaluation Metrics: The definition of the RMSE is as follows:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{r}_i - r_i)^2}; \quad (7)$$

where  $r_i$  and  $\hat{r}_i$  are the true RUL and the predicted RUL, respectively. The scoring function [23] is defined as

$$S = \begin{cases} P & \text{if } \hat{r}_i \geq r_i \\ \frac{1}{10} \left( e^{\frac{\hat{r}_i - r_i}{13}} + 1 \right) & \text{when } \hat{r}_i < r_i \end{cases} \quad (8)$$

