



# BMSFormer: An efficient deep learning model for online state-of-health estimation of lithium-ion batteries under high-frequency early SOC data with strong correlated single health indicator



Xiaopeng Li <sup>a</sup>, Minghang Zhao <sup>a,c,\*</sup>, Shisheng Zhong <sup>a,b,c</sup>, Junfu Li <sup>d</sup>, Song Fu <sup>b</sup>, Zhiqi Yan <sup>e</sup>

<sup>a</sup> Department of Mechanical Engineering, Harbin Institute of Technology, Weihai, Shandong, 264209, China

<sup>b</sup> School of Mechatronics Engineering, Harbin Institute of Technology, Harbin, Heilongjiang, 150000, China

<sup>c</sup> Weihai Key Laboratory of Intelligent Operation and Maintenance, Harbin Institute of Technology, Weihai, Shandong, 264209, China

<sup>d</sup> School of Automotive Engineering, Harbin Institute of Technology, Weihai, Shandong, 264209, China

<sup>e</sup> College of Aeronautical Engineering, Civil Aviation University of China, Tianjin, 300300, China

## ARTICLE INFO

Handling editor: Prof X Ou

### Keywords:

State of health

Lithium-ion batteries

Efficiency estimation

Local-global fusion attention

Depthwise feature fusion

## ABSTRACT

The efficient and accurate state-of-health (SOH) estimation is crucial for reducing risks and ensuring effective application in battery management systems (BMS) of resource-limited devices. However, many recent state-of-the-art SOH estimation approaches rely on resource-consuming structures to obtain good performance. In this paper, an efficient deep learning model for SOH estimation, namely BMSFormer, is constructed. BMSFormer mainly integrates a Local-Global Fusion Attention module to capture both long-term and short-term dependencies while reducing computational complexity compared to traditional Softmax-based attention. Additionally, two kinds of depthwise separable convolution are embedded to fuse multi-scale and multi-channel features, enhancing feature diversity with fewer parameters than standard convolution. Three widely used battery datasets, each with different chemistries and operating conditions, are employed to evaluate the performance of BMSFormer. The experiments results illustrate that the proposed model achieves higher accuracy, lower computational consumption, and stabler performance across various hyperparameter combinations compared to alternative models.

## 1. Introduction

Nowadays, with the advantages of high energy density, low self-discharge rate and long service life, lithium-ion batteries have become one of the main energy storage devices in civil airlines, electric vehicles, mobile devices, aerospace and other domains [1]. Although, lithium-ion batteries bring convenience to people's lives, they also pose latent safety hazards. These hazards primarily arise from the complicated degradation processes brought by complex internal chemical structure, different working conditions and environment. Therefore, precise monitoring and state estimation are essential to enhance battery safety, reliability, and performance.

As an essential guard of battery safe and efficient operation, battery management system (BMS) can record real-time data, manage cell balance, estimate internal status and provide emergency protection [2–4]. The internal states, such as state-of-health (SOH), cannot be measured

directly. At the same time, due to the limited storage space of embedded devices, they often rely on online estimations using simple approaches, such as Kalman filter (KF) and its variants [5–7], classical physical models [8], equivalent circuit models (ECM) [9] and traditional machine learning (ML) algorithms like support vector machines (SVM) [10] and random forest (RF) [11]. However, on one hand, when faced with nonlinear and unstable data from online monitoring and historical cycles, traditional models struggle to provide high performance due to their structural constraints [12,13]; On the other hand, high-performance deep models are generally not feasible enough for convenient deployment on BMS due to their computationally expensive structures. Actually, since the emergence of AlexNet [14], the predominant trend in model architecture has been towards increasing depth and size, but many real-world applications need to be performed in real-time and/or on limited resource mobile devices [15]. Thereby, the high-performance deep models have a certain demand to be compact

\* Corresponding author. Department of Mechanical Engineering, Harbin Institute of Technology, Weihai, Shandong, 264209, China.

E-mail address: [zhaomh@hit.edu.cn](mailto:zhaomh@hit.edu.cn) (M. Zhao).

and exhibit low computational cost for practical applications in the future.

The recent SOH estimation approaches could be divided into two main categories: (1) model-based approaches, including electrochemical models and empirical models, (2) data-driven approaches, including machine learning and deep learning methods [16–19].

- (1) Model-based approaches simulate the electrochemical principles in batteries using mathematical equations or circuit components [20]. Electrochemical models (*e.g.*, single-particle model and pseudo-two-dimensional model) utilize a series of differential equations to represent the electrochemical reaction processes in the battery. While these models may offer high accuracy, they often require extensive computational resources, making online estimation challenging; and many parameters, such as particle radius and diffusion coefficients, are difficult to obtain [21,22]. In contrast, as a widely employed empirical model, equivalent circuit model (ECM), simulates battery internal behavior with circuit components such as resistors and capacitors [9,23]. ECMs (*e.g.*, Thevenin, PNGV, GNL and RC models) have relatively low computational complexity and easy parameter identification. However, they cannot fully capture battery internal changes, probably resulting in lower accuracy compared to electrochemical models.
- (2) Data-driven approaches can be trained using historical data and mostly have lower dependencies on the internal degradation mechanism of batteries compared to model-based approaches [23]. They employ various algorithms to establish the relationships between input (health indicators, HIs) and output (state-of-health, SOH), then use new data for SOH estimation. Thereby, since a few years ago, many researchers tend to focus on establishing stronger data-driven models. For instance, Fei et al. [24] crafted 42 HIs from five kinds of main recorded data including voltage, current, temperature, internal resistance and capacity, and then fed them into five representative machine learning (ML) models including the Gaussian process regression (GPR) [25], support vector machine (SVM) [10], random forest (RF) [26], gradient boosting regression tree (GBRT) [27] and neural network (NN) [28]. Benefiting from the prosperity of deep learning (DL) methods, in recent years, many works employ convolutional neural networks (CNNs), recurrent neural networks (RNNs), and Transformers to estimate the states of batteries. For example, depending on the powerful local feature extraction ability of CNNs, Yang et al. [29] utilized 1D-CNN to handle capacity degradation datasets, while Lee et al. [30] transformed the data into a two-dimensional format and employed 2D-CNN. Additionally, dense neural network [31] was also attempted in similar research.

However, traditional CNNs are not adept at dealing with long time-series data, as this requires numerous convolutional layers to fuse local features that are distant from each other. RNNs and their variants such as gated recurrent unit (GRU) and long short-term memory (LSTM) attempt to overcome this challenge with hidden layers, state cells and various gates. These components function like a “nerve cell” that can store information and pass it from one time step to the next for computation. Nonetheless, due to the sequential dependence, it is not convenient enough for traditional RNNs to make full use of modern multi-threaded GPU resources, even though recent works have achieved some improvements in computational efficiency like bidirectional long short-term memory (Bi-LSTM) [32,33], Attention-based Spatial-Temporal LSTM (ASTLSTM) [34], CNN-LSTM [35], CNN-ASTLSTM [36], and CNN-GRU [37], while the fundamental constraint of limited-sequential computation still remains.

To address the above problem, Vaswani [38] proposed the Transformer, a model architecture avoiding recurrence and instead relying on

a self-attention mechanism to capture global dependencies between input and output. Due to its significant advantage in parallel computing and global modeling, Transformer has already achieved immense success across multiple domains [39–41]. Since Hannan et al. [42] applied the Transformer to estimate SOH of batteries in 2021, subsequent models have tended to increase in complexity. For instance, Gomez et al. [43] introduced an improved temporal fusion Transformer, integrating a Bi-LSTM encoder-decoder layer to enhance performance. Jia et al. [44] proposed a hybrid prediction model by combining bidirectional gated recurrent unit and Transformer. In another study, Gu et al. [45] suggested a novel approach using data pre-processing methods and a CNN-Transformer framework. Bai et al. [46] introduced a 2D-CNN-Transformer model using matrixing data, while Chen et al. [47] applied a vision-Transformer network (ViT) for SOH estimation, involving extra data flattening operations. Most existing methods do not attempt to simplify the structures of the Transformer-based networks, instead, they integrated it directly with other modules or models, resulting in computationally complicated approaches.

Three main problems faced by many recent approaches for battery SOH estimation, can be summarized as follows.

- (1) **Computational Complexity Limitations:** In the field of battery state estimation, many existing deep learning models attempt to achieve higher accuracy by increasing the complexity of model structures or/and combining different classical models. Those approaches introduce numerous unnecessary parameter updates, thereby increasing computational resource demands. Such complexity may hinder future applications in resource-constrained devices, such as BMS.
- (2) **Estimation Accuracy Challenges:** Many current approaches face accuracy challenges due to inherent computational constraints within the model itself, insufficient representational capabilities of the extracted health indicators, or the input of many weakly correlated health indicators, which together negatively impact the accuracy of SOH estimation results.
- (3) **Performance Stability Issues:** The training and optimization process often require extensive hyperparameter tuning to achieve higher performance, which consumes a lot of time and computational resources before reaching the desired goals. Such consumption can be effectively reduced if the model is able to maintain stable and high performance across various hyperparameter configurations. However, most existing methods only present their optimal results without demonstrating the model's ability to solve this problem.

To address these challenges, this article develops an efficient deep learning model, namely BMSFormer, for SOH estimation of lithium-ion batteries, with the goal of optimizing computational resource usage, balancing accuracy and efficiency, and enhancing model stability for reliable performance. The main contributions are summarized as follows.

- (1) By progressively shortening the window size and step size, health indicators (HIs) with higher correlations to the battery health state are identified within two high-frequency SOC charge and discharge segments. This approach successfully achieves Pearson correlation coefficient (PCC) values averaging over 0.99 across three different kinds of battery.
- (2) A Local-Global Fusion Attention module is constructed to effectively capture both short-term and long-term features while reducing computational complexity compared to traditional Softmax-based attention mechanisms.
- (3) Two kinds of depthwise separable convolution module with small kernel size (DSConv-S) and large kernel size (DSConv-L) are embedded to fuse multi-scale and multi-channel features,

enhancing feature diversity with fewer parameters than standard convolution.

The remaining sections of this paper are arranged as follows. Section 2 describes the major steps of the developed SOH estimation approach and the basic information of three different kinds of battery datasets. Section 3 introduces the structure and modules of BMSFormer. Section 4 presents comprehensive validations and compares BMSFormer with four prevailing deep learning models. Finally, the conclusions are summarized in Section 5.

## 2. Preliminaries

### 2.1. Overview of the developed SOH estimation approach

The SOH is a critical metric for evaluating health status of batteries. While there is still no standardized definition for SOH, battery capacity remains the most utilized representation [31], as illustrated below:

$$\text{SOH} = \frac{C_{\text{current}}}{C_{\text{rated}}} \times 100\%, \quad (1)$$

where  $C_{\text{current}}$  and  $C_{\text{rated}}$  are the current capacity and rated capacity, respectively.

As illustrated in Fig. 1, the developed SOH estimation approach

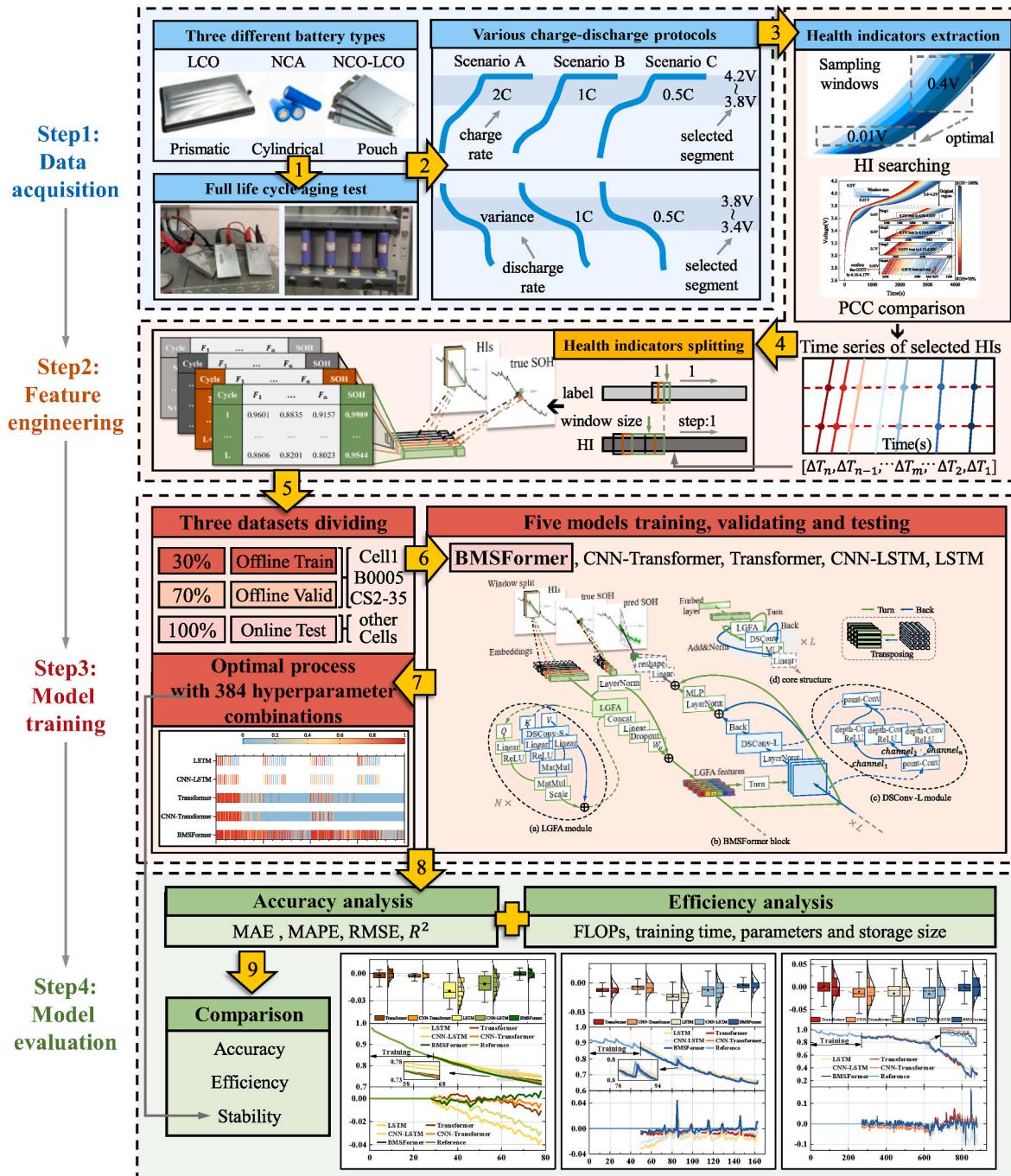


Fig. 1. Flowchart of developed SOH estimation approach.

consists of four steps.

- (1) Data acquisition. Three types of batteries with different chemical materials were subjected to full lifecycle aging tests under various charge/discharge protocols to evaluate the proposed model. The high frequency of the intermediate SOC during daily battery operation, charge segment from 3.8 V to 4.2 V and discharge segment from 3.8 V to 3.4 V are selected for HI extraction.
- (2) Feature engineering. The constant current charge and discharge times for each cycle are extracted as health indicators (HIs). Starting from the selected segments, the HI search process is refined using progressively smaller windows and steps until a 0.01V window is reached or no higher Pearson correlation coefficient (PCC) value is found. A sliding window is then used to divide the time series of HIs into subsets, moving forward one step at a time. The true SOH at the next step of each window serves as the label for its corresponding subset.
- (3) Model training. For convenience, Cell1, B0005, and CS2-35 are used as training sets, with the first 30 % of the data employed for training across 384 hyperparameter combinations including epochs, learning rate, number of blocks and layers, and dimension of dense and embedding layer. The remaining 70 % data is used for validation and comparison to select the best-performing model. Then the best perform model is tested directly on the entire data of other batteries in the corresponding dataset to evaluate the model generalization and compare the performance of BMSFormer with four different deep learning models.
- (4) Model evaluation. The accuracy, efficiency, and stability of different models are evaluated. Four typical evaluation metrics are used to measure the accuracy, four common computational complexity metrics are used to evaluate the training efficiency, and the  $R^2$  results of 384 hyperparameter combinations are used to assess the stability of training results.

## 2.2. Typical battery datasets acquisition

Three kinds of prevailing battery datasets from Oxford, NASA and CALCE are used in this paper, which are widely recognized and originate from leading battery research institutions, ensuring standardized experimental protocols and reliable data. They encompass various battery materials, specifications, and work conditions. Hence, the applicability and precision of the proposed approach can be evaluated through evaluation on a diverse range of batteries with different materials and subjected to various operating conditions. The main characteristics of battery datasets used in this work are summarized in Table 1; Fig. 2(a)–(c) illustrates the capacity degradation and increment phenomenon of three different batteries under various work conditions, while Fig. 2(d)–(f) displays the constant-current charge voltage curves in several cycle of the first battery from each dataset. Similar trends are observed within each battery type and distinct differences are noticeable between different types of batteries.

### 2.2.1. Oxford lithium battery dataset

The Oxford lithium-ion battery dataset is provided by the Battery Intelligence Laboratory at the University of Oxford including degradation data from eight Kokam SLPB533459H4 lithium-ion pouch cells, each with a 0.74 Ah nominal capacity and test at 40 °C. The negative electrode material of the Kokam pouch cells is graphite, the positive electrode material is a blend of lithium cobalt oxide (LCO) and lithium nickel cobalt oxide (NCO). The eight batteries are subjected to a constant-current and then constant-voltage charging profile up to 4.2 V, followed by a variance discharge current rate from the urban Artemis-derived profile process down to 2.7 V.

### 2.2.2. NASA lithium battery dataset

The NASA lithium battery dataset is provided by Prognostics Center

**Table 1**  
Description of three battery datasets.

| Data Sources                    | Oxford [48]                                     | NASA [49]   | CALCE [50]                        |
|---------------------------------|---|---|-----------------------------------|
| Manufacturer                    | Kokam   | LG Chem   | LG Chem                           |
| Cell types                      | Pouch   | 18650 Cylindrical                                 | Prismatic                         |
| Cathode material                | NCO-LCO   | NCA   | LCO                               |
| Selected battery series         | Cell (1–8)                                      | B00 (05–07 and 18)                                | CS2 (35–38)                       |
| Battery number                  | 8   | 4   | 4                                 |
| Rated capacity (Ah)             | 0.74  | 2   | 1.1                               |
| Capacity failure threshold (Ah) | 70 %  | 1.4 (70 %)  | 0.84 (76.4 %)                     |
| Charge protocol (rate)          | CC (2C)   | CC-CV (0.75C, 4.2V)                               | CC-CV (0.5C, 4.2V)                |
| Discharge protocol (rate)       | variance  | CC (1C)   | CC (1C)                           |
| Cut-off voltage (V)             | charge to 4.2<br>discharge to 2.7               | charge to 4.2<br>discharge to 2.7/<br>2.5/2.2/2.5 | charge to 4.2<br>discharge to 2.7 |
| Cut-off current (A)             | –   | charge to 0.02                                    | charge to 0.05                    |
| Experiment temperature (°C)     | 40  | 24  | 24                                |
| Cycle numbers                   | 8200/7700/8100/<br>5100/5000/5000/<br>8100/8100 | 168/168/168/<br>132                               | 881/935/<br>971/995               |

of Excellence Data Set Repository-NASA. Four Lithium-ion 18650 cylindrical batteries (#5, 6, 7, and 18) from LG Chem, each with 2Ah nominal capacity, undergo three operational profiles at room temperature: charging at 1.5A (0.75C) in CC mode to 4.2V, followed by CV mode until the current fell to 20 mA; discharging at 2A in CC mode down to 2.7V, 2.5V, 2.2V and 2.5V, respectively. The aging process is accelerated through repeated cycles, with EOL defined as a 30 % capacity loss to 1.4Ahr. The positive electrode material is lithium nickel cobalt aluminum oxide (NCA).

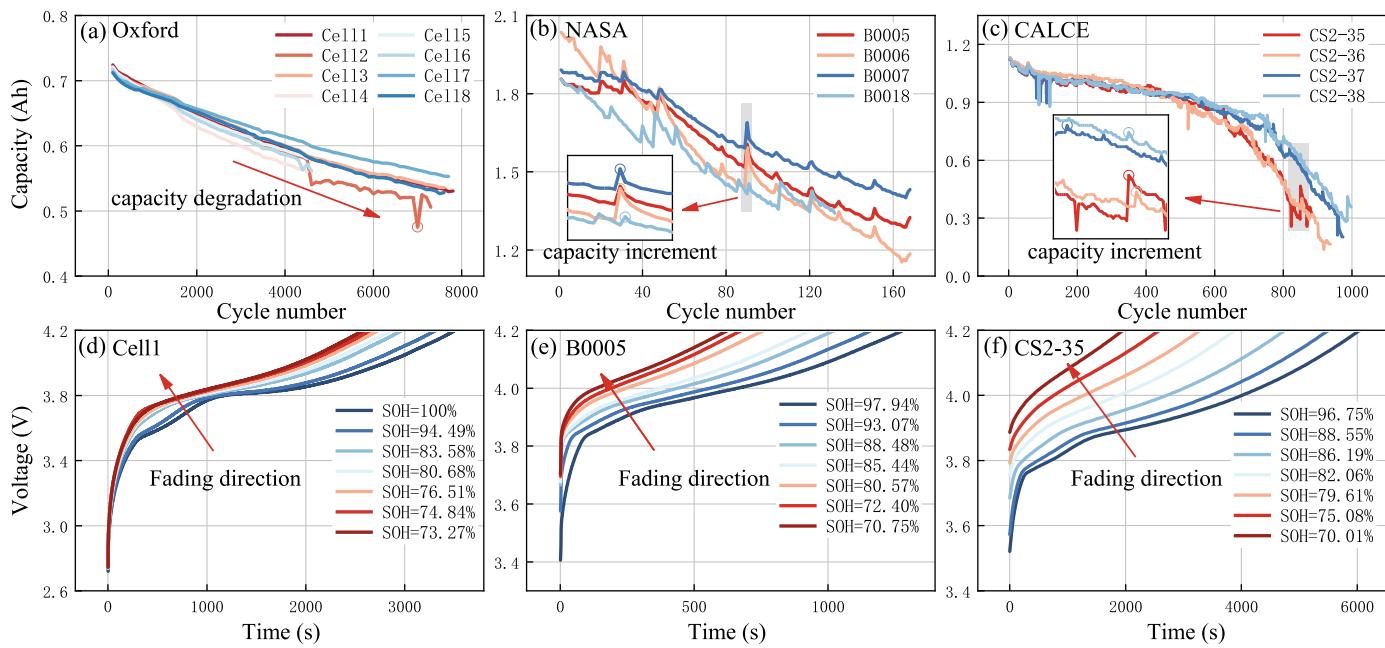
### 2.2.3. CALCE lithium battery dataset

The CALCE lithium-ion battery dataset originates from the Center for Advanced Life Cycle Engineering (CALCE) at the University of Maryland. These batteries have a rated capacity of 1.1 Ah and utilize lithium cobalt oxide (LiCoO<sub>2</sub>) as the cathode material. The aging tests for battery life are performed at room temperature. The charge strategy is to charge at a constant-current rate of 0.5 C (0.55A) until the voltage reaches 4.2 V and then hold at 4.2 V until the charging current drops below 0.02 A, then the batteries are discharged with a constant current of 2 A until the voltage decreases to 2.7 V.

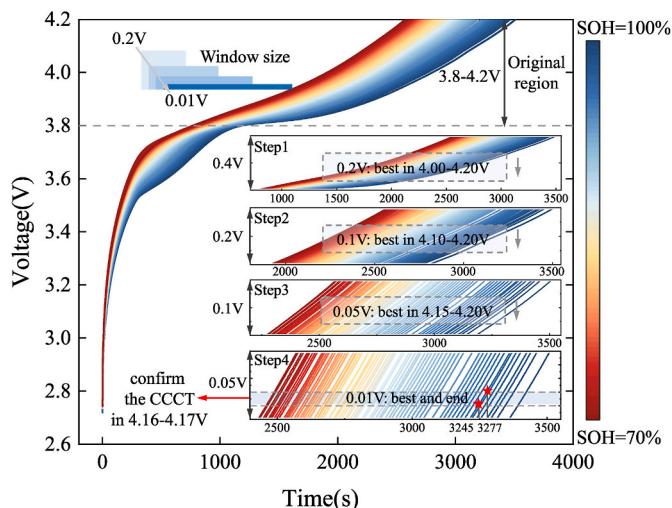
## 2.3. Health indicators extraction

The selected health indicators (HIs) are inspired by a range of references [51–58] serving various research purposes. Considering the high frequency of the intermediate SOC during daily battery operation, charge segment from 3.8 V to 4.2 V and discharge segment from 3.8 V to 3.4 V are selected for HI extraction.

As shown in Fig. 3, in the first step of our HIs extraction procedures, a 0.2V window size with a 0.2V moving step size is utilized to divide the selected charge or discharge segment into shorter segments: 3.8–4.0V, 3.9–4.1V, and 4.0–4.2V in charge segment (3.8–3.6V, 3.7–3.5V, 3.6–3.4V in discharge segment). Then the minimum Pearson correlation coefficient (PCC) value across all batteries of each dataset is taken as the baseline, and the segment with the highest baseline value is selected as the search segment for the next step. In the 2nd, 3rd, and 4th steps, the search procedures are repeated with progressively smaller windows and steps until a 0.01V window is reached or no higher baseline is identified than in the previous step. Finally, the constant current charge time (CCCT) series from 4.16V to 4.17V and the constant current discharge time (CCDT) series from 3.8V to 3.4V are respectively confirmed as the



**Fig. 2.** Three different kinds of battery datasets: (a)–(c) capacity degradation and temporary increment curves of the Oxford, NASA and CALCE cells, respectively. (d)–(f) constant current charge voltage curves of Oxford Cell1, NASA B0005 and CALCE CS2-35, respectively.



**Fig. 3.** The extraction procedures of the selected HI in the selected charge segment.

selected HIs.

The PCC is a widely utilized statistical metric for evaluating the strength of linear association between HIs and true SOH [59,60]. The value of the PCC ranges from  $-1$  to  $1$ , with the absolute value closer to  $1$

indicating a stronger relationship. The PCC formula is shown as follows:

$$\text{PCC} = \frac{\sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})}{\sqrt{\sum_{k=1}^n (x_k - \bar{x})^2 \sum_{k=1}^n (y_k - \bar{y})^2}} \quad (2)$$

where  $x, y, \bar{x}, \bar{y}$  and  $n$  represent the value of the selected HI sequence and the battery SOH sequence, the average value of  $x$  and  $y$ , the number of samples, respectively.

As shown in Table 2, the selected HI, which is the constant current charging time from  $4.16$  to  $4.17$ V, has the highest PCC compared to others of recent methods. Similarly, as shown in Table 3, the selected HI, which is the constant current discharge time from  $3.8$  to  $3.4$ V, has the highest PCC compared to others.

### 3. The developed BMSFormer

In this section, the overall architecture and workflow of BMSFormer are first introduced. Then, the designed Multi-scale Depthwise Separable Convolution modules, including DSConv-L and DSConv-S, are described. Following that, the proposed Local-Global Fusion Attention (LGFA) module is detailed. Finally, several foundational modules are briefly introduced.

**Table 2**  
PCC comparison of different HIs on Oxford datasets.

| Battery series | selected HI | CCCT from ICA [55] | CCCT from DTA [56] | CCQ [57] | CCDQ [55] | IC peak [58] | sample entropy [55] | Discharge voltage curve slope [58] |
|----------------|-------------|--------------------|--------------------|----------|-----------|--------------|---------------------|------------------------------------|
| Cell1          | 0.99996     | 0.9999             | 0.9989             | 0.99947  | 0.9998    | 0.9444       | -0.9853             | 0.9999                             |
| Cell2          | 0.99999     | 0.9999             | 0.9959             | 0.98677  | 0.9999    | -            | -0.9839             | -                                  |
| Cell3          | 0.99996     | 0.9999             | 0.9993             | 0.99948  | 0.9989    | 0.9375       | -0.9876             | 0.9998                             |
| Cell4          | 0.99998     | 0.9104             | -                  | 0.99981  | 0.9999    | -            | -0.9885             | -                                  |
| Cell5          | 0.99992     | 0.9999             | -                  | 0.99960  | 0.9999    | -            | -0.9133             | -                                  |
| Cell6          | 0.99998     | 0.9999             | -                  | 0.99961  | 0.9991    | -            | -0.9848             | -                                  |
| Cell7          | 0.99999     | 0.9999             | -                  | 0.99973  | 0.9999    | 0.9296       | -0.9896             | 0.9998                             |
| Cell8          | 0.99999     | 0.9999             | -                  | 0.99970  | 0.9999    | 0.9378       | -0.9908             | 0.9998                             |

**Table 3**

PCC comparison of different HIs on NASA and CALCE datasets.

| Battery series | selected HI   | CCCT [61] | CCCQ [57] | Max projection distance of start and end charge voltage [62] | IC peak [58] | Relax time between charge and discharge cycle [63] |
|----------------|---------------|-----------|-----------|--|--------------|--|
| B0005          | <b>0.9934</b> | 0.9703    | 0.94268   | 0.95   | 0.9910       | -0.8941  |
| B0006          | <b>0.9858</b> | 0.9366    | 0.94518   | 0.97   | 0.9800       | -0.9373  |
| B0007          | <b>0.9983</b> | 0.9745    | 0.92037   | 0.99   | 0.9251       | -0.9049  |
| B0018          | <b>0.9987</b> | -         | 0.81248   | -  | -            | -  |
| CS2-35         | <b>0.9960</b> | 0.9799    | 0.90290   | 0.95   | -            | -0.9843  |
| CS2-36         | <b>0.9968</b> | 0.9827    | 0.93370   | 0.99   | -            | -0.9829  |
| CS2-37         | <b>0.9943</b> | 0.9789    | 0.93920   | 0.98   | -            | -  |
| CS2-38         | <b>0.9952</b> | 0.9748    | -         | 0.96   | -            | -  |

### 3.1. Architecture overview

In this research, to solve the problem that the accuracy and efficiency of the current models often cannot be optimal together, a lightweight network termed BMSFormer with both global features and local information capturing capabilities is proposed.

The entire framework is illustrated in Fig. 4 and described as follows: HIs are taken as input, segmented into fragments via window splitting, embedded in a high-dimensional space, and then fed into a BMSFormer block, which includes the LGFA module and DSConv-L module. The output from LGFA module is transposed, passed through a DSConv-L module, and transposed back. Finally, the output from all BMSFormer blocks is fed into a Multilayer Perceptron (MLP) layer before outputting the result. During the training process, the true SOH value right following each window segment is utilized as a label, which guides the model to adjust the gradient descent optimization.

Mathematically, this means, for the input  $x_i$  of BMSFormer block  $i$ , the corresponding output  $y_i$  of the block is calculated using

$$x'_i = W_a \cdot \text{LGFA}(x_i) + \text{LN}(x_i), \quad (3)$$

$$x''_i = \text{Back}(\text{DSConv-L}(\text{LN}(\text{Turn}(x'_i))) + x'_i), \quad (4)$$

$$y_i = \text{MLP}(\text{LN}(x''_i)) + x'_i, \quad (5)$$

where LGFA, LN, DSConv-L, MLP respectively represent the Local-Global Fusion Attention, layer normalization, depthwise separable

convolution with large kernel size, and multilayer perceptron. Turn and Back refer to the transposing operation of tensor, and  $W_a$  is a learnable weight parameter.

### 3.2. The designed multi-scale depthwise separable convolutional modules

In this section, the fundamental structure of DSConv is introduced and compared to standard convolution, then the multi-scale properties of DSConv-S and DSConv-L are elaborated subsequently.

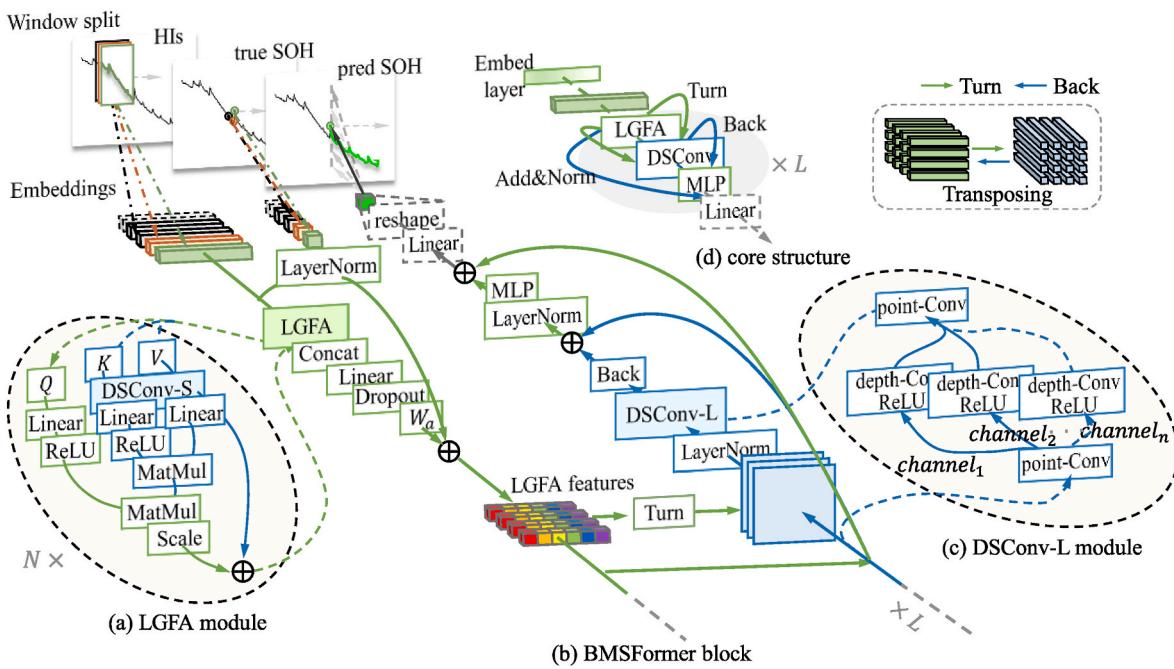
#### 3.2.1. Fundamental structure of DSConv-S and DSConv-L modules

Convolutions operations are adept in extracting local information through sliding filters [30,46,64]. Standard convolution operates across all input channels, with each convolutional kernel generating a feature map as one output channel. Therefore, the consumption of computational resources and training time will significantly increase when dealing with large time-series datasets. The standard convolutions have the computational cost of:

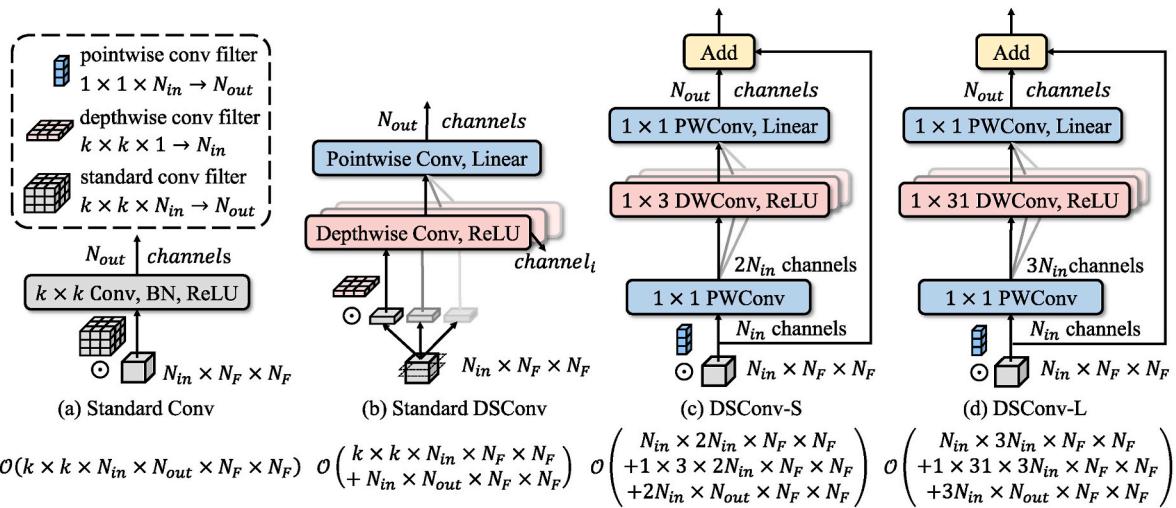
$$k \times k \times N_{in} \times N_{out} \times N_F \times N_F \quad (6)$$

where the  $k, N_{in}, N_{out}, N_F$  represent the kernel size, input channels, output channels, and the feature map size, respectively.

Compared to standard convolution, firstly, in general, depthwise separable convolution utilizes smaller filter than standard convolution [15,65–68], as illustrated in the upper left area of Fig. 5. Secondly, depthwise convolution separates the features map into multiple



**Fig. 4.** Framework of BMSFormer. (a) LGFA module, (b) BMSFormer block, (c) DSConv-L module, (d) BMSFormer core structure.



**Fig. 5.** The fundamental structure of: (a) Standard convolution; (b) Standard Depthwise separable convolution; (c) DSConv-S module; (d) DSConv-L module.

channels, and each channel is convolved with its own set of filters. These two operations not only allow for the extraction of a variety of features but also significantly decreases the computational load. Thirdly, the outputs from the depthwise step are then combined using  $1 \times 1$  filter of pointwise convolution. This step mixes the features across all channels, allowing for learning interactions between them. Finally, the residual connection helps in mitigating the vanishing gradient problem and allows the network to learn more complex functions.

DSConv separates the spatial and channel-wise operations, reducing computational complexity by first applying a depthwise convolution across spatial dimensions, followed by a pointwise convolution to adjust the number of output channels:

$$k \times k \times N_{in} \times N_F \times N_F + N_{in} \times N_{out} \times N_F \times N_F \quad (7)$$

where  $k, N_{in}, N_{out}, N_F$  represent the kernel size, input channels, output channels, and the feature map size, respectively.

By breaking convolution operation into a two-step process of depthwise filtering followed by pointwise combination, a reduction is achieved in computation of:

$$\frac{k \times k \times N_{in} \times N_F \times N_F + N_{in} \times N_{out} \times N_F \times N_F}{\frac{1}{N_{out}} + \frac{1}{k^2}} \quad (8)$$

### 3.2.2. Computational analysis of DSConv-S and DSConv-L modules

The depthwise separable convolution with small kernel size (DSConv-S) and large kernel size (DSConv-L) are designed to extract multi-scale and multi-channel feature, aiming to addresses the expressiveness limitations of ReLU attention, which is constrained by its linear attention mechanism, while maintaining high computational efficiency.

As illustrated in Fig. 5, with a double expansion factor for input channels, the DSConv-S module applies a  $1 \times 3$  depthwise convolutional kernel size to enhance the quality of input information and the sensitivity of attention for local variety. This module further incorporates two layers of  $1 \times 1$  pointwise convolutions to integrate multi-channel features, further improving feature representation. The computational cost of DSConv-S can be expressed as:

$$N_{in} \times 2N_{in} \times N_F \times N_F + 1 \times 3 \times 2N_{in} \times N_F \times N_F + 2N_{in} \times N_{out} \times N_F \times N_F \quad (9)$$

where  $2N_{in}$  and  $N_F \times N_F$  represents the input channels after expanding and the feature map size, respectively.

With a triple expansion factor for input channels, the DSConv-L

module applies a  $1 \times 31$  depthwise convolutional kernel size to extract long-term features, aiming to further enhance the feature diversity and model generalization. This module also incorporates two layers of  $1 \times 1$  pointwise convolutions. Therefore, the computational cost of DSConv-L can be expressed as:

$$N_{in} \times 3N_{in} \times N_F \times N_F + 1 \times 31 \times 3N_{in} \times N_F \times N_F + 3N_{in} \times N_{out} \times N_F \times N_F \times N_F \quad (10)$$

where  $3N_{in}$  and  $N_F \times N_F$  represents the input channels after expanding and the feature map size, respectively.

### 3.3. The proposed Local-Global Fusion Attention module

In this section, the general form of self-attention mechanism in Transformers is described, the advantages and disadvantages of traditional Softmax and linear attention are briefly analyzed, and two main superiorities of proposed Local-Global Fusion Attention (LGFA) module, including global receptive field and local information sensitivity are introduced subsequently.

Most existing SOH estimation methods still rely on traditional Softmax attention [45–47], suffering from excessive computational complexity. Meanwhile, although the linear attention reduces the complexity, it continues to struggle with insufficient model expressiveness. Therefore, the Local-Global Fusion Attention (LGFA) module is proposed to integrate linear complexity and high expressiveness, which are potentially beneficial for efficient and accurate SOH estimation.

#### 3.3.1. General form of multi-head self-attention

To improve the model's representational capacity, multi-head self-attention processes the input through several parallel "heads", each representing an independent self-attention mechanism. Each head calculates its own attention weights and generates a unique attention representation. The self-attention mechanism computes the relationships (attention weights) between each element in the input and all other elements. It then produces a weighted sum of these relationships to create new representations. This allows the model to focus on important information in the sequence, especially long-range dependencies.

In  $i$ -th head, given an input  $X \in \mathbb{R}^{N \times C}$ , the general form of self-attention can be formulated as follows:

$$Q = XW_Q, \quad (11)$$

$$K = XW_K, \quad (12)$$

$$V = XW_V, \quad (13)$$

$$O_i = \sum_{j=1}^N \frac{\text{Sim}(Q_i, K_j)}{\sum_{j=1}^N \text{Sim}(Q_i, K_j)} V_j, \quad (14)$$

where  $W_Q$ ,  $W_K$ , and  $W_V \in \mathbb{R}^{C \times d}$  are the learnable linear projection matrices,  $C$  and  $d$  are the channel dimensions of module and each head.  $\text{Sim}(\cdot, \cdot)$  is the similarity function of self-attention mechanism.

### 3.3.2. Softmax attention mechanism and linear attention mechanism

Transformers mainly utilize Softmax Attention [38] with the  $\text{Sim}(Q, K) = \exp(QK^T / \sqrt{d})$  in Eq. (14). However, the attention map is obtained by computing the similarity between all query-key pairs, leading to the computation complexity of  $\mathcal{O}(N^2)$  as illustrated in Fig. 6(a). In contrast, linear attention [69], as shown in Fig. 6(b), efficiently reduces the computation complexity by carefully designed similarity functions with  $\text{Sim}(Q, K) = \phi(Q)\phi(K)^T$ , which changes the computation order from  $(\phi(Q)\phi(K)^T)V$  to  $\phi(Q)(\phi(K)^T V)$  base on the associative property of matrix multiplication. The computation complexity with respect to the number of tokens in the sequence is reduced to  $\mathcal{O}(N)$ , where  $N$  represents the sequence length. This linear complexity allows the model to efficiently handle long sequences, which is particularly beneficial for time series analysis.

Traditional Softmax attention amplifies the differences between all query-key pairs using the exponential function, significantly enhancing larger similarity values when calculating attention weights. Although Softmax attention provides stronger expressive power, its computational cost is a bit unbeneficial to its application in large-scale sequences. In contrast, linear attention uses different kernel functions  $\phi(\cdot)$  to map  $Q$  and  $K$  into a new space, thereby enhancing the differences. While this may result in some loss of expressiveness, the new calculation order of the weight matrix in linear attention reduces computational complexity, which offers an advantage in processing long time-series sequences like battery SOH estimation.

### 3.3.3. The proposed attention module

As illustrated in Fig. 6(c), the proposed Local-Global Fusion Attention (LGFA) module utilizes the Rectified Linear Unit attention mechanism and enhances feature diversity through the DSConv-S module introduced in Section 3.2.1. The similarity function Eq. (14) of this

attention mechanism is defined as:

$$\text{Sim}(Q, K) = \text{ReLU}(Q)\text{ReLU}(\text{DSConv-S}(K)^T) \quad (15)$$

where the self-attention module can be rewritten as:

$$O_i = \frac{\sum_{j=1}^N [\text{ReLU}(Q_i)\text{ReLU}(\text{DSConv-S}(K_j)^T)] \text{DSConv-S}(V_j)}{\text{ReLU}(Q_i) \sum_{j=1}^N \text{ReLU}(\text{DSConv-S}(K_j)^T)} \quad (16)$$

Then, by using the associative property of matrix multiplication, the computational complexity can be reduced from quadratic to linear while maintaining the same functionality:

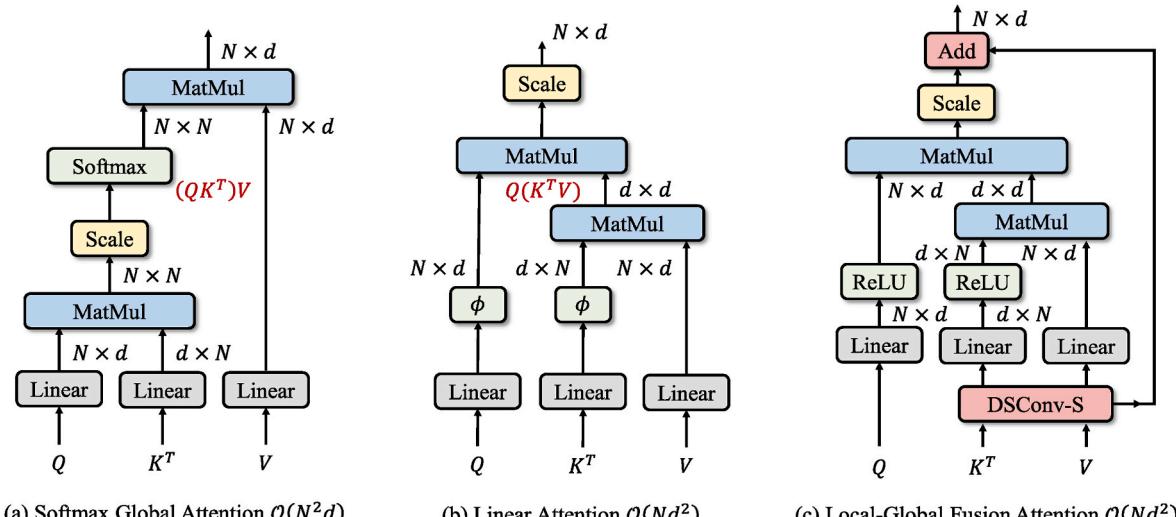
$$O_i = \frac{\sum_{j=1}^N \{\text{ReLU}(Q_i)\text{ReLU}[\text{DSConv-S}(K_j)^T]\} \text{DSConv-S}(V_j)}{\text{ReLU}(Q_i) \sum_{j=1}^N \text{ReLU}[\text{DSConv-S}(K_j)^T]} \\ = \frac{\text{ReLU}(Q_i) \left\{ \sum_{j=1}^N \text{ReLU}[\text{DSConv-S}(K_j)^T] \text{DSConv-S}(V_j) \right\}}{\text{ReLU}(Q_i) \sum_{j=1}^N \text{ReLU}[\text{DSConv-S}(K_j)^T]} \quad (17)$$

As shown in Eq. (17), the computation of  $\sum_{j=1}^N \text{ReLU}(K_j)^T V_j \in \mathbb{R}^{d \times d}$  and  $\sum_{j=1}^N \text{ReLU}(K_j)^T \in \mathbb{R}^{d \times 1}$  are performed only once, after which these results can be reused for each query. This approach thereby requires only  $\mathcal{O}(N)$  computational cost and  $\mathcal{O}(N)$  storage memory.

Another key advantage of ReLU-based global attention is that it prevents the vanishing gradient problem and involves simple operation  $\max(0, x)$ , making it faster to compute than Softmax, which requires exponentiation and normalization. For instance, some previous latency comparisons [70] show that, with similar computational demands, ReLU-based linear attention outperforms Softmax attention, particularly on mobile devices.

## 4. Experiments and analysis

This section presents a comprehensive evaluation of five models (BMSFormer, CNN-Transformer, CNN-LSTM, Transformer and LSTM) across three public battery datasets. The Cell1, B0005, and CS2-35 are employed as training sets, with the first 30 % of the data utilized for training across 384 hyperparameter combinations, including variations



**Fig. 6.** Difference between (a) traditional Softmax attention mechanism, (b) traditional Linear Attention mechanism, and (3) the proposed Local-Global Fusion attention module.

in epochs, learning rate, number of blocks and layers, and the dimensions of the dense and embedding layers. The remaining 70 % of the data is reserved for validation and comparison to identify the best-performing model parameters. The selected model is then directly tested on the entire datasets of other batteries to assess its generalization capability and to compare the performance of BMSFormer against the alternative models. Moreover, the robustness of the five models is compared through a detailed analysis of the 384 training results, without any further hyperparameter refinement.

The training process is conducted on an ASUS TUF A15 computer, which is equipped with an AMD Ryzen R7-4800H (2.90 GHz), an NVIDIA GeForce RTX 2060 and 24 GB of RAM. The operation system on this machine is Windows 11 Professional. The prediction models are built using PyTorch 2.0.0, utilizing Python 3.9.16 as the programming language.

#### 4.1. The evaluation criteria

Several widely adopted evaluation metrics are selected to assess the models' performance: Mean Absolute Error (MAE), Mean Absolute Percentage Error (MAPE), Root Mean Squared Error (RMSE), the coefficient of determination  $R^2$ , and the Average Root Mean Squared Error (ARMSE). MAE measures the average magnitude of the errors in predictions without considering their direction, which refers to whether the predicted values are above or below the actual values, providing a straightforward measure of prediction accuracy. MAPE expresses accuracy as a percentage, making it useful for comparing performance across different datasets. RMSE is the square root of the average squared differences between predicted and actual values, giving larger weight to larger errors. The  $R^2$  score indicates the proportion of variance in the dependent variable predictable from the independent variables, with larger values showing better performance. ARMSE, the average of RMSE values over several subsets or folds, provides a stable measure of model performance across different data partitions. The formulas for these metrics are as follows:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (18)$$

$$\text{MAPE} = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \quad (19)$$

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (20)$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (21)$$

$$\text{ARMSE} = \frac{1}{m} \sum_{j=1}^m \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (22)$$

where  $y$ ,  $\hat{y}$ ,  $\bar{y}$ ,  $n$  and  $m$  are the actual SOH, predicted SOH, the average value of  $y$ , the number of testing samples, and the number of batteries, respectively [71].

#### 4.2. Comparison on Oxford dataset

##### 4.2.1. Accuracy and generalization comparison

Comparative experiments are conducted between BMSFormer and four alternative models. The layer configurations of these models are presented in Table 4.

Fig. 7 and Table 5 show the estimation results and errors of BMSFormer, CNN-Transformer, CNN-LSTM, Transformer and LSTM on

**Table 4**  
The layer configurations of different deep learning models.

| BMSFormer              | CNN-Transformer  | Transformer            | CNN-LSTM                        | LSTM            |
|------------------------|--|------------------------|---------------------------------|-----------------|
| Input Embed layer (32) | Input Embed layer (32)                                     | Input Embed layer (32) | Input Conv1D (1 × 3, 1)<br>ReLU | Input Linear    |
| LGFA module            | Positional Encoder   | Positional Encoder     | MaxPooling1D (2 × 2, 2)         | LSTM layer (32) |
| DSConv-L module        | Softmax Attention  | Softmax Attention      | Conv1D (1 × 3, 1)<br>ReLU       | LSTM layer (32) |
| Add&Norm               | Conv1D (1 × 3, 1)<br>ReLU                                  | Add&Norm               | LSTM layer (32)                 | LSTM layer (32) |
| MLP                    | MaxPooling1D (2, 2)  | Softmax Attention      | LSTM layer (32)                 | Linear          |
| Add&Norm               | Conv1D (1 × 3, 1)<br>ReLU,<br>Add&Norm                     | Add&Norm               | LSTM layer (32)                 |                 |
| Linear                 | Softmax Attention<br>Add&Norm<br>MLP<br>Add&Norm<br>Linear | MLP                    | Linear                          |                 |

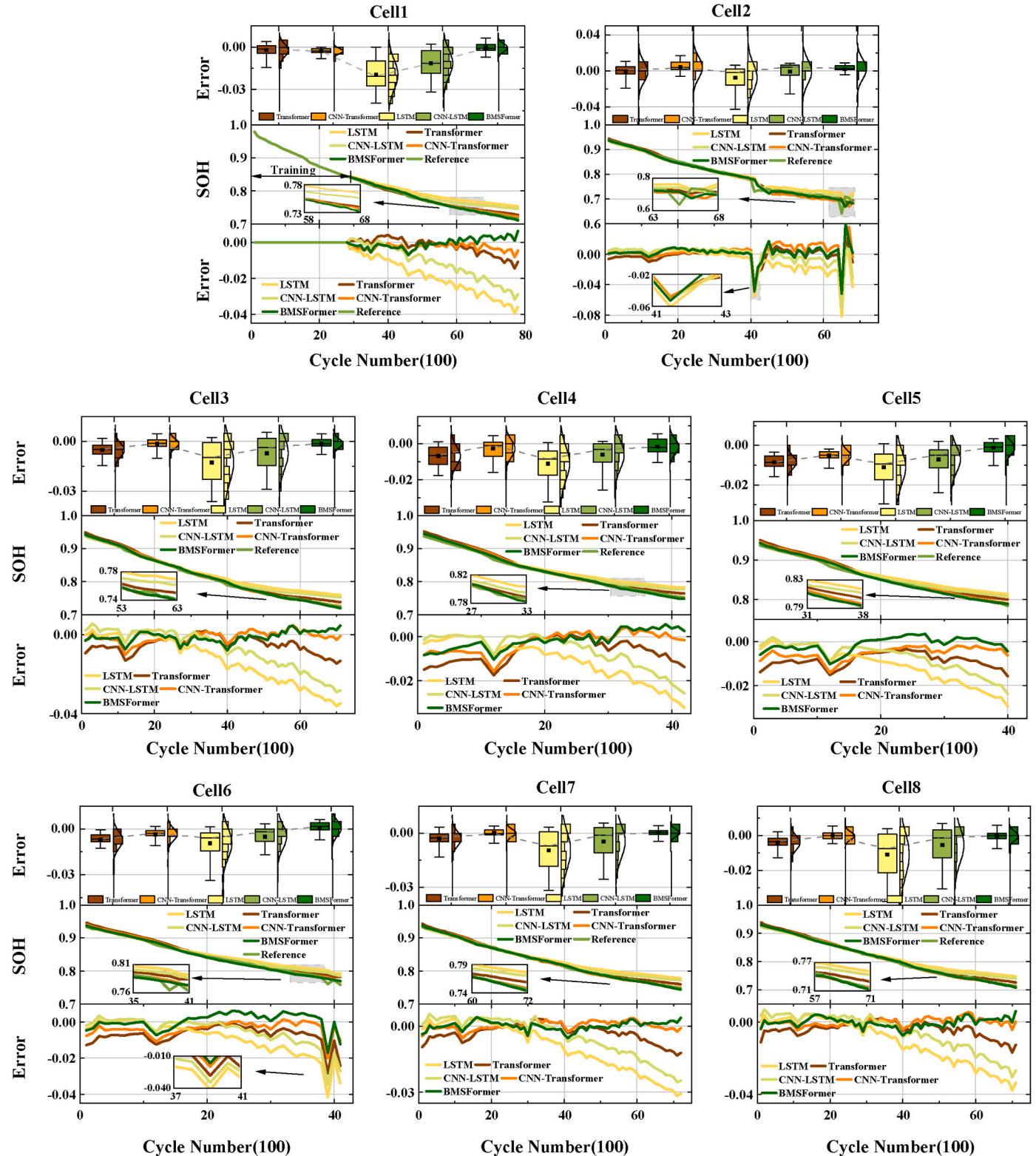
the Oxford dataset. As shown in the figure, in the validation results for Cell1, CNN-Transformer's estimation error IQR is close to zero. However, the asymmetry in the whiskers of the boxplot indicates a greater number of low outliers in the estimation results. This trend is also observed from Cell2 to Cell8, suggesting that the model's estimation stability may be insufficient. The boxplot of prediction errors and the normal distribution plot across the eight batteries clearly show that BMSFormer has smaller and more concentrated prediction errors, illustrating superior prediction accuracy and robustness.

Notably, in Cells 2 and 6, where there are several sudden drops, CNN-Transformer and BMSFormer were able to track the true values most accurately, with BMSFormer showing the smallest prediction errors in response to sudden changes in Cell6. For the other six batteries with relatively stable SOH changes (excluding Cell3 and Cell8), BMSFormer performed the best and ranked first in the mean prediction results across all eight batteries, showing improvements in the average MAE of 14.81 %, 59.65 %, 46.51 %, 73.56 %, the average RMSE of 16.21 %, 62.19 %, 44.64 %, and 73.27 %, and the average MAPE of 20.00 %, 63.63 %, 50.00 %, 76.32 % compared to CNN-Transformer, Transformer, CNN-LSTM, and LSTM, respectively.

##### 4.2.2. Efficiency comparison

In practical application, it is important to consider training cost (floating-point operations, training time, parameters) and hardware cost (storage size) besides the accuracy [72]. To ensure a fair comparison, the same training and model hyperparameters are set to evaluate the comprehensive performance of different models. The definitions of these evaluation metrics and the settings of hyperparameters are described as follows.

- Evaluation Metrics:** (1) floating-point operations (FLOPs), which are measured using the profile function from the THOP library, quantify the numbers of FLOPs required for a single forward pass. (2) Training time, measured in seconds using Python's time module, represents the duration required to train a model on the training dataset. (3) Parameters, which are counted using the summary function from the torchsummary library, typically measured in millions (Million), indicate the total number of trainable parameters in the model. (4) Storage size, which is counted using the os.path.



**Fig. 7.** The experimental results of different models on the Oxford dataset.

- getsize function of Python, measured in kilobytes (KB), indicates the storage requirements of the model parameters on a device.
2. **Training Hyperparameters:** For the experiments, the same training hyperparameters are used across all models. Each model is trained with a batch size of 128 for 1000 epochs. The learning rate is uniformly set to 0.01, and a dropout rate of 0.1 is applied.
  3. **Model Hyperparameters:** To compare the computational complexity of different model structures under the same input data and training hyperparameters, the model hyperparameters are set as follows: For models based on the Transformer architecture, the parameters are set uniformly: an embed dimension of 16, a dense dimension (MLP) of 16, four attention heads, and both the encoder and decoder comprise one layer each. Four LSTM layers and hidden

**Table 5**  
SOH estimation errors of different models on the Oxford dataset.

| Cells   | Methods         | MAE           | MAPE          | RMSE          | R <sup>2</sup> |
|---------|-----------------|---------------|---------------|---------------|----------------|
| Cell1*  | BMSFormer       | 0.0032        | <b>0.0022</b> | <b>0.0018</b> | <b>0.9937</b>  |
|         | CNN-Transformer | <b>0.0022</b> | 0.0038        | 0.0025        | 0.9913         |
|         | CNN-LSTM        | 0.0085        | 0.0154        | 0.0109        | 0.8641         |
|         | Transformer     | 0.0027        | 0.0049        | 0.0037        | 0.9814         |
|         | LSTM            | 0.0140        | 0.0252        | 0.0160        | 0.7044         |
| Cell2   | BMSFormer       | <b>0.0040</b> | <b>0.0074</b> | <b>0.0083</b> | <b>0.9802</b>  |
|         | CNN-Transformer | 0.0058        | 0.0107        | 0.0091        | 0.9752         |
|         | CNN-LSTM        | 0.0059        | 0.0105        | 0.0095        | 0.9755         |
|         | Transformer     | 0.0047        | 0.0084        | 0.0082        | 0.9798         |
|         | LSTM            | 0.0075        | 0.0139        | 0.0124        | 0.9765         |
| Cell3   | BMSFormer       | 0.0019        | 0.0031        | 0.0023        | 0.9977         |
|         | CNN-Transformer | <b>0.0017</b> | <b>0.0028</b> | <b>0.0022</b> | <b>0.9978</b>  |
|         | CNN-LSTM        | 0.0062        | 0.0110        | 0.0087        | 0.9710         |
|         | Transformer     | 0.0040        | 0.0068        | 0.0049        | 0.9898         |
|         | LSTM            | 0.0093        | 0.0165        | 0.0126        | 0.9395         |
| Cell4   | BMSFormer       | <b>0.0027</b> | <b>0.0044</b> | <b>0.0032</b> | <b>0.9936</b>  |
|         | CNN-Transformer | 0.0029        | 0.0046        | 0.0041        | 0.9901         |
|         | CNN-LSTM        | 0.0045        | 0.0077        | 0.0069        | 0.9765         |
|         | Transformer     | 0.0052        | 0.0083        | 0.0064        | 0.9753         |
|         | LSTM            | 0.0077        | 0.0132        | 0.0102        | 0.9487         |
| Cell5   | BMSFormer       | <b>0.0020</b> | <b>0.0031</b> | <b>0.0026</b> | <b>0.9940</b>  |
|         | CNN-Transformer | 0.0039        | 0.0061        | 0.0043        | 0.9829         |
|         | CNN-LSTM        | 0.0051        | 0.0085        | 0.0069        | 0.9649         |
|         | Transformer     | 0.0064        | 0.0101        | 0.0068        | 0.9591         |
|         | LSTM            | 0.0078        | 0.0128        | 0.0098        | 0.9286         |
| Cell6   | BMSFormer       | <b>0.0026</b> | <b>0.0043</b> | <b>0.0035</b> | <b>0.9905</b>  |
|         | CNN-Transformer | 0.0028        | 0.0045        | 0.0039        | 0.9884         |
|         | CNN-LSTM        | 0.0041        | 0.0069        | 0.0067        | 0.9720         |
|         | Transformer     | 0.0054        | 0.0087        | 0.0067        | 0.9673         |
|         | LSTM            | 0.0068        | 0.0115        | 0.0095        | 0.9429         |
| Cell7   | BMSFormer       | <b>0.0011</b> | 0.0018        | <b>0.0014</b> | <b>0.9987</b>  |
|         | CNN-Transformer | 0.0012        | <b>0.0015</b> | 0.0019        | <b>0.9987</b>  |
|         | CNN-LSTM        | 0.0053        | 0.0091        | 0.0072        | 0.9719         |
|         | Transformer     | 0.0029        | 0.0047        | 0.0037        | 0.9917         |
|         | LSTM            | 0.0076        | 0.0132        | 0.0104        | 0.9422         |
| Cell8   | BMSFormer       | 0.0015        | 0.0026        | 0.0020        | 0.9983         |
|         | CNN-Transformer | <b>0.0012</b> | <b>0.0021</b> | <b>0.0016</b> | <b>0.9989</b>  |
|         | CNN-LSTM        | 0.0060        | 0.0106        | 0.0082        | 0.9739         |
|         | Transformer     | 0.0033        | 0.0057        | 0.0043        | 0.9924         |
|         | LSTM            | 0.0087        | 0.0156        | 0.0119        | 0.9455         |
| Average | BMSFormer       | <b>0.0023</b> | <b>0.0036</b> | <b>0.0031</b> | <b>0.9934</b>  |
|         | CNN-Transformer | 0.0027        | 0.0045        | 0.0037        | 0.9904         |
|         | CNN-LSTM        | 0.0057        | 0.0099        | 0.0082        | 0.9587         |
|         | Transformer     | 0.0043        | 0.0072        | 0.0056        | 0.9796         |
|         | LSTM            | 0.0087        | 0.0152        | 0.0116        | 0.9160         |

Cell1\* denotes a training set of BMSFormer with 80 % data, and its values are averaged among ten experiments. The best results are marked in bold and underlined.

dimension of 16 are also employed for models based on the LSTM architecture. Additionally, for the CNN-Transformer models, the CNN module corresponds to our depthwise separable convolutional module and is also placed after the multi-head self-attention module. The basic hyperparameters for the two convolutional layers are set as follows: the first convolutional layer has a kernel size of 3 and a stride of 1, the second pooling layer has a kernel size of 2 and a stride of 2,

and the third convolutional layer has a kernel size of 3 and a stride of 1. The activation function used is ReLU.

Under the almost same configuration and Oxford dataset, as illustrated in Table 6, the LSTM-based networks results in lower training time than other three networks due to the inherently sequential structure; however, this architecture limits their ability to capture complex degradation patterns, leading to lower estimation accuracy. In contrast, compared to CNN-Transformer and Transformer, BMSFormer achieves the shortest training time. Specifically, despite having 63.04% more FLOPs and 12.03% more parameters than Transformer, BMSFormer still reduces training time and storage capacity by 21.37% and 7.88%, respectively.

When adapting to different battery types or tasks, models usually require adjustments in hyperparameter combinations to match data patterns efficiently. Among the five adjustable hyperparameters (including  $e$ ,  $lr$ ,  $b$ ,  $d_e$  and  $n$ ), the LSTM network lacks embed dimension ( $d_e$ ), and storage size depends only on model hyperparameters (including  $d_d$  and  $n$ ). Therefore, to comprehensively evaluate storage memory consumption across different hyperparameter combinations, three hidden layer dimensions (16, 32, 64) and four model layers (1, 2, 3, 4) are configured, resulting in a total of 12 combinations applied to five models, as shown in Table 7.

As shown in Fig. 8(a), the bar chart for the 12 combinations indicates that BMSFormer consistently exhibits a lower storage size compared to CNN-Transformer and Transformer across the G1-G12 configurations, making it more suitable for deployment on resource-constrained devices. Moreover, LSTM and CNN-LSTM models experience a significant increase in memory usage when dimensions are raised to 64 (as shown in G3, G6, G9, G12), far exceeding the memory usage of the other three models. Although lower dimensions and fewer layers reduce the storage size of LSTM and CNN-LSTM, these models may suffer from a decline in estimation accuracy, failing to achieve an optimal balance between precision and efficiency.

As illustrated in Fig. 8(b), the box plot shows that the proposed model has the shortest box and whiskers, indicating the lowest average

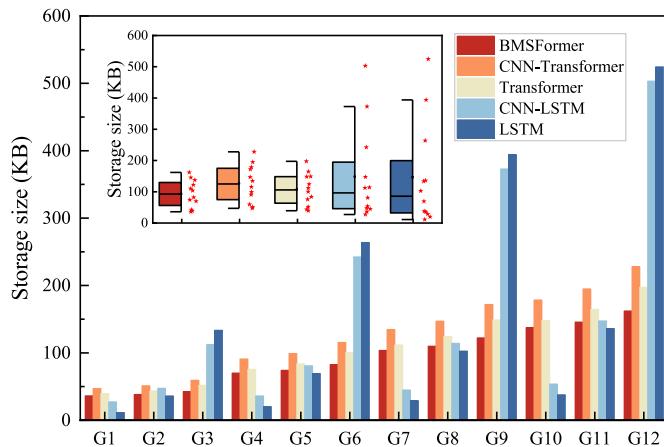
**Table 7**  
Grouping of models hyperparameter.

| Group number | Dense dimension | Layer number |
|--------------|-----------------|--------------|
| G1           | 16              | 1            |
| G2           | 32              | 1            |
| G3           | 64              | 1            |
| G4           | 16              | 2            |
| G5           | 32              | 2            |
| G6           | 64              | 2            |
| G7           | 16              | 3            |
| G8           | 32              | 3            |
| G9           | 64              | 3            |
| G10          | 16              | 4            |
| G11          | 32              | 4            |
| G12          | 64              | 4            |

**Table 6**  
The comprehensive performance of models under the same and optimal configurations.

| Models          | Training hyperparameters |      |     | Model hyperparameters |       |     | Performance indicator |                   |             |                  |
|-----------------|--------------------------|------|-----|-----------------------|-------|-----|-----------------------|-------------------|-------------|------------------|
|                 | $e$                      | $lr$ | $b$ | $d_e$                 | $d_d$ | $n$ | FLOPs (Million)       | Training time (s) | Parameters  | Storage size(KB) |
| BMSFormer       | 1000                     | 0.01 | 128 | 16                    | 16    | 4   | 0.46                  | 19.83             | 5330        | <b>36.37</b>     |
| CNN-Transformer | 1000                     | 0.01 | 128 | 16                    | 16    | 4   | 0.27                  | 26.57             | 6257        | 47.10            |
| Transformer     | 1000                     | 0.01 | 128 | 16                    | 16    | 4   | <b>0.17</b>           | 25.22             | <b>4689</b> | 39.48            |
| CNN-LSTM        | 1000                     | 0.01 | 128 | -                     | 16    | 4   | 0.48                  | 11.84             | 12719       | 53.72            |
| LSTM            | 1000                     | 0.01 | 128 | -                     | 16    | 4   | 0.83                  | <b>8.07</b>       | 8753        | 37.85            |

Here, the  $e$ ,  $lr$ ,  $b$ ,  $d_e$ ,  $d_d$  and  $n$  represent the epochs, learning rate, batch size, embedding dimension of MLP or hidden dimension of LSTM, and layer number, respectively. The best results are marked in bold and underlined.



**Fig. 8.** Storage size of five models across G1-G12 hyperparameter combinations: (a) Bar chart comparing storage sizes for each group; (b) Box plot showing overall distribution.

memory consumption with minimal fluctuation across the 12 hyperparameter configurations. This characteristic is helpful when handling more complex data patterns, as the model maintains stable storage size while increasing dimensions or depth to enhance estimation accuracy.

In summary, the experimental results indicate that the proposed model exhibits low and stable memory usage across 384 hyperparameter combinations, making it well-suited for usage on resource-limited devices.

#### 4.3. Comparison on NASA and CALCE datasets

This section discusses the performance of BMSFormer and four mainstream deep learning models on the NASA and CALCE datasets, which are different from the Oxford dataset in battery materials and charge/discharge protocols. For both datasets, 30 % of the data from the first cell (B0005 for NASA and CS2-35 for CALCE) is used as the training set, while the remaining 70 % is reserved for validation. The full datasets from the remaining six cells (B0006, B0007, B0018 for NASA and CS2-36, CS2-37, CS2-38 for CALCE) are used as the test set.

The five models' evaluation process is described as follows. First, the five models are respectively trained on their training sets using 384 distinct hyperparameter combinations and validated. All the validation results are then used to evaluate the stability of each model and the best-performance models, with highest validation accuracy, are selected as the final trained models. Second, the selected five final models are then directly tested on the unseen test sets to further evaluate the model generalization.

##### 4.3.1. Stability evaluation across 384 hyperparameter combinations

Table 8 presents the comprehensive set of hyperparameter combinations. Each column represents an adjustable hyperparameter category, including the train epochs, the learning rate, the number of attention block, the number of LSTM layer (or attention head), the dense layer dimension of MLP (or the hidden dimension of LSTM), and the dimension of embedding layer. This can help identify the sets of hyperparameter combinations that yield the high validation accuracies and reflects the model's stability across different combinations.

To ensure a comprehensive exploration of training hyperparameters, "epochs" are chosen as 1000 and 200, and "learning rates" are set to

0.001 and 0.01, resulting in a total of  $2 \times 2 = 4$  basic combinations, which correspond to the four segments on the x-axis of Fig. 9. Meanwhile, various model hyperparameter settings are also considered: BMSFormer, CNN-Transformer, and Transformer are based on attention mechanism and each have  $2 \times 3 \times 4 \times 4 = 96$  possible hyperparameter groups, in which 2, 3, 4 and 4 represent the "block number", "head number", "dense dimension" of MLP layer, and "embed dimension" of embedding layer, respectively. In contrast, LSTM and CNN-LSTM models just have two additional adjustable hyperparameters for the "layer number" of LSTM and the "dense dimension" of hidden layer, resulting in  $3 \times 4 = 12$  combinations.

Consequently, the total number of hyperparameter combinations of the models represented on the y-axis of the figure is  $4 \times 12 = 48$  for the first two models and  $4 \times 96 = 384$  for the latter three models. Each color bar corresponds to a fitting score, with redder colors indicating higher accuracy and bluer colors indicating worse. The red vertical line on the far right represents the number of combinations with a fitting score greater than 0.98, with a higher number of groups indicating better performance.

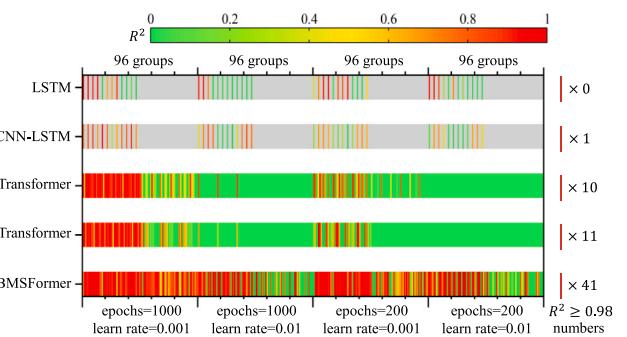
Overall, compared to recent mainstream models, BMSFormer shows robust performance and accuracy across different hyperparameter combinations. Specifically, BMSFormer has 41 combinations where the prediction model's  $R^2$  value exceeds 0.98, representing 10.67 % of the total combinations. Additionally, 43.23 % of the combinations have an  $R^2$  value greater than 0.9. BMSFormer also achieves the highest  $R^2$  of 0.9892 than others.

##### 4.3.2. Generalization evaluation

This section presents and compares the testing results obtained through the process of directly testing on the unseen test sets using the five optimal models with highest validation accuracies.

As illustrated in Fig. 10 (a)-(c), the estimated SOH of BMSFormer is generally closer to the true SOH compared to other models, and the estimation error is also closer to zero, both indicating higher accuracy. The localized magnifications in the 'SOH' plots further illustrate the better stability of BMSFormer when encountering sudden SOH changes in the degradation process of ten batteries.

As detailed in Table 9, BMSFormer achieves improvements of 37.15 %, 42.65 %, 23.67 %, and 47.90 % in average MAE, 27.30 %, 48.20 %, 20.42 %, and 53.71 % in average MAPE, and 34.29 %, 40.69 %, 18.84 %, and 46.29 % in average RMSE, compared to CNN-Transformer, CNN-LSTM, Transformer, and LSTM, respectively, across eight battery cells from NASA and CALCE datasets. Notably, the zoomed-in prediction

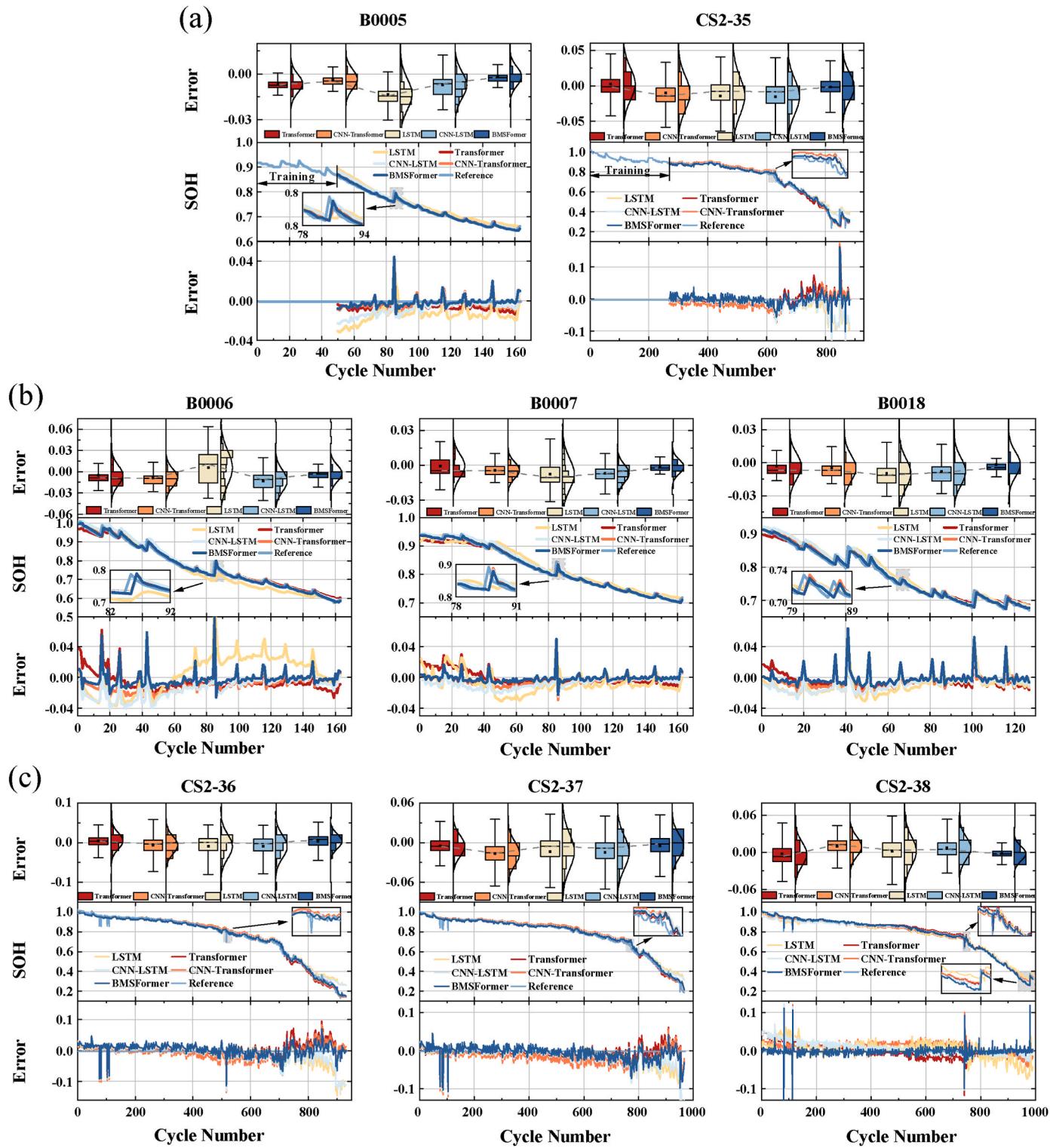


**Fig. 9.** The whole validating results when using different models and hyperparameter combinations.

**Table 8**

The combination of hyperparameter.

| Hyperparameter | Epochs      | Learning rate | Block number | Layer number | Dense dimension   | Embedding dimension |
|----------------|-------------|---------------|--------------|--------------|-------------------|---------------------|
| Values         | [200, 1000] | [0.001, 0.01] | [1, 2]       | [1, 2, 4]    | [16, 32, 64, 128] | [16, 32, 64, 128]   |



**Fig. 10.** (a) The validation results of five models on B0005 and CS2-35. (b) The test results using five weighted models trained on B0005. (c) The test results using five weighted models trained on CS2-35.

curves reveal that BMSFormer performs exceptionally well even in scenarios with the most severe SOH fluctuations, such as those observed in B0018 and CS2-38.

LSTM still exhibits the lowest prediction accuracy due to insufficient sensitivity to local features. Although CNN-LSTM partially addresses this issue, it still suffers from structural limitations, resulting in the loss of important information after convolution processing. Transformer also

struggles with local feature recognition, leading to larger fluctuations in error curves and particularly noticeable deviations in B0005. CNN-Transformer and BMSFormer show similar performance, but CNN-Transformer's small convolution kernel size of  $1 \times 3$  and the lack of feature diversity constrain its ability to identify complex features pattern, resulting in slightly lower prediction accuracy compared to BMSFormer. Increasing the convolution kernel size in CNN-Transformer

**Table 9**  
Estimation results on the NASA and CALCE datasets.

| Cell    | Methods         | MAE           | MAPE          | RMSE          | R <sup>2</sup> |
|---------|-----------------|---------------|---------------|---------------|----------------|
| B0005*  | BMSFormer       | <b>0.0042</b> | <b>0.0053</b> | <b>0.0065</b> | <b>0.9950</b>  |
|         | CNN-Transformer | 0.0071        | 0.0089        | 0.0089        | 0.9907         |
|         | CNN-LSTM        | 0.0104        | 0.0129        | 0.0121        | 0.9830         |
|         | Transformer     | 0.0091        | 0.0115        | 0.0109        | 0.9861         |
|         | LSTM            | 0.0140        | 0.0181        | 0.0157        | 0.9716         |
| B0006   | BMSFormer       | <b>0.0073</b> | <b>0.0117</b> | <b>0.0092</b> | <b>0.9905</b>  |
|         | CNN-Transformer | 0.0123        | 0.0151        | 0.0155        | 0.9843         |
|         | CNN-LSTM        | 0.0158        | 0.0194        | 0.0197        | 0.9730         |
|         | Transformer     | 0.0118        | 0.0156        | 0.0148        | 0.9848         |
|         | LSTM            | 0.0209        | 0.0274        | 0.0246        | 0.9581         |
| B0007   | BMSFormer       | <b>0.0037</b> | <b>0.0045</b> | <b>0.0061</b> | <b>0.9939</b>  |
|         | CNN-Transformer | 0.0062        | 0.0074        | 0.0081        | 0.9893         |
|         | CNN-LSTM        | 0.0086        | 0.0102        | 0.0105        | 0.9821         |
|         | Transformer     | 0.0079        | 0.0095        | 0.0098        | 0.9841         |
|         | LSTM            | 0.0131        | 0.0159        | 0.0152        | 0.9626         |
| B0018   | BMSFormer       | <b>0.0068</b> | <b>0.0111</b> | <b>0.0087</b> | <b>0.9771</b>  |
|         | CNN-Transformer | 0.0093        | 0.0119        | 0.0123        | 0.9713         |
|         | CNN-LSTM        | 0.0122        | 0.0153        | 0.0148        | 0.9591         |
|         | Transformer     | 0.0094        | 0.0122        | 0.0120        | 0.9729         |
|         | LSTM            | 0.0139        | 0.0178        | 0.0163        | 0.9503         |
| CS2-35* | BMSFormer       | <b>0.0116</b> | <b>0.0208</b> | <b>0.0197</b> | <b>0.9892</b>  |
|         | CNN-Transformer | 0.0198        | 0.0262        | 0.0324        | 0.9808         |
|         | CNN-LSTM        | 0.0191        | 0.0416        | 0.0289        | 0.9767         |
|         | Transformer     | 0.0140        | 0.0259        | 0.0227        | 0.9856         |
|         | LSTM            | 0.0193        | 0.0439        | 0.0304        | 0.9742         |
| CS2-36  | BMSFormer       | <b>0.0144</b> | <b>0.0242</b> | <b>0.0188</b> | <b>0.9936</b>  |
|         | CNN-Transformer | 0.0169        | 0.0306        | 0.0233        | 0.9902         |
|         | CNN-LSTM        | 0.0197        | 0.0536        | 0.0306        | 0.9831         |
|         | Transformer     | 0.0138        | 0.0278        | 0.0209        | 0.9921         |
|         | LSTM            | 0.0204        | 0.0582        | 0.0320        | 0.9815         |
| CS2-37  | BMSFormer       | 0.0125        | 0.0221        | 0.0184        | 0.9904         |
|         | CNN-Transformer | 0.0196        | 0.0245        | 0.0304        | 0.9831         |
|         | CNN-LSTM        | 0.0185        | 0.0386        | 0.0286        | 0.9770         |
|         | Transformer     | <b>0.0117</b> | <b>0.0202</b> | <b>0.0168</b> | <b>0.9920</b>  |
|         | LSTM            | 0.0179        | 0.0402        | 0.0299        | 0.9749         |
| CS2-38  | BMSFormer       | <b>0.0214</b> | <b>0.0356</b> | <b>0.0276</b> | <b>0.9778</b>  |
|         | CNN-Transformer | 0.0391        | 0.0615        | 0.0441        | 0.9432         |
|         | CNN-LSTM        | 0.0385        | 0.0696        | 0.0487        | 0.9309         |
|         | Transformer     | 0.0296        | 0.0473        | 0.0338        | 0.9667         |
|         | LSTM            | 0.0377        | 0.0708        | 0.0500        | 0.9271         |
| Average | BMSFormer       | <b>0.0102</b> | <b>0.0169</b> | <b>0.0143</b> | <b>0.9884</b>  |
|         | CNN-Transformer | 0.0163        | 0.0233        | 0.0219        | 0.9791         |
|         | CNN-LSTM        | 0.0179        | 0.0327        | 0.0242        | 0.9706         |
|         | Transformer     | 0.0134        | 0.0213        | 0.0177        | 0.9830         |
|         | LSTM            | 0.0197        | 0.0365        | 0.0268        | 0.9625         |

B0005\* and CS2-35\* denote that they are the training sets with 30 % of the early data. The best results are marked in bold and underlined.

might improve performance but would also increase the computational cost.

## 5. Conclusion

This paper presents BMSFormer, an innovative and efficient deep learning model for online state-of-health (SOH) estimation of lithium-ion batteries, alleviating the limitations of traditional approaches that cannot keep well balance between accuracy and efficiency. By incorporating the constructed Local-Global Fusion Attention module, BMSFormer can not only capture both short-term and long-term features, but also reduce computational complexity compared to traditional Softmax-based attention methods. Meanwhile, the integration of multi-scale Depthwise Separable Convolutions enhances feature extraction across multiple scales and channels, contributing to a more diverse representation of battery degradation characteristics. These advancements collectively provide technical advantages in tackling challenges related to battery behaviors, such as random capacity increments and gradual performance degradation.

The comprehensive experiments on the three different kinds of public battery datasets have shown the advantages of BMSFormer in estimation accuracy, efficiency, and generalization. The validation

results also illustrate the stability of BMSFormer across 384 hyper-parameter combinations. These advancements make BMSFormer a potentially robust solution for SOH estimation, particularly in some resource-constrained environments.

Finally, there is still room for further optimization of the model architecture, particularly in enhancing computational efficiency and reducing model complexity. Future research will focus on refining the structure of BMSFormer, expanding its applicability to a wider range of battery types and use cases, and integrating the model into Battery Management Systems (BMS) for real-world deployment and evaluation. These efforts aim to extend the practical impact of BMSFormer in advancing the field of battery health management.

## CRediT authorship contribution statement

**Xiaopeng Li:** Writing – original draft, Visualization, Validation, Software, Methodology. **Minghang Zhao:** Writing – review & editing, Methodology, Formal analysis. **Shisheng Zhong:** Project administration, Funding acquisition. **Junfu Li:** Writing – review & editing, Methodology. **Song Fu:** Writing – review & editing. **Zhiqi Yan:** Supervision, Project administration.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgments

This work is supported by National Key R&D Program of China (2023YFB4302400).

## Data availability

Data will be made available on request.

## References

- Krause FC, Ruiz JP, Jones SC, Brandon EJ, Darcy EC, Iannello CJ, Bugga RV. Performance of commercial Li-ion cells for future NASA missions and aerospace applications. *J Electrochim Soc* 2021;168:040504. <https://doi.org/10.1149/1945-7111/abf05f>.
- Barré A, Deguilhem B, Grolleau S, Gérard M, Suard F, Riu D. A review on lithium-ion battery ageing mechanisms and estimations for automotive applications. *J Power Sources* 2013;241:680–9. <https://doi.org/10.1016/j.jpowsour.2013.05.040>.
- Ma S, Jiang M, Tao P, Song C, Wu J, Wang J, Deng T, Shang W. Temperature effect and thermal impact in lithium-ion batteries: a review. *Prog Nat Sci: Mater Int* 2018; 28:653–66. <https://doi.org/10.1016/j.pnsc.2018.11.002>.
- Xia Y, Chen G, Zhou Q, Shi X, Shi F. Failure behaviours of 100% SOC lithium-ion battery modules under different impact loading conditions. *Eng Fail Anal* 2017;82. <https://doi.org/10.1016/j.englfailanal.2017.09.003>.
- Plett GL. Sigma-point Kalman filtering for battery management systems of LiPB-based HEV battery packs: Part 2: simultaneous state and parameter estimation. *J Power Sources* 2006;161:1369–84. <https://doi.org/10.1016/j.jpowsour.2006.06.004>.
- Plett GL. Extended Kalman filtering for battery management systems of LiPB-based HEV battery packs: Part 1. Background. *J Power Sources* 2004;134:252–61. <https://doi.org/10.1016/j.jpowsour.2004.02.031>.
- Peng S, Zhang D, Dai G, Wang L, Jiang Y, Zhou F. State of charge estimation for LiFePO4 batteries joint by PID observer and improved EKF in various OCV ranges. *Appl Energy* 2024;377:124435. <https://doi.org/10.1016/j.apenergy.2024.124435>.
- Chaturvedi NA, Klein R, Christensen J, Ahmed J, Kojic A. Algorithms for advanced battery-management systems. *IEEE Control Syst Mag* 2010;30:49–68. <https://doi.org/10.1109/MCS.2010.936293>.
- Zhang C, Allafi W, Dinh Q, Ascencio P, Marco J. Online estimation of battery equivalent circuit model parameters and state of charge using decoupled least squares technique. *Energy* 2018;142:678–88. <https://doi.org/10.1016/j.energy.2017.10.043>.
- Hearst MA, Dumais ST, Osuna E, Platt J, Scholkopf B. Support vector machines. *IEEE Intell Syst Their Appl* 1998;13:18–28. <https://doi.org/10.1109/5254.708428>.

- [11] Zhang Q-J, Gupta KC, Devabhaktuni VK. Artificial neural networks for RF and microwave design - from theory to practice. *IEEE Trans Microw Theor Tech* 2003; 51:1339–50. <https://doi.org/10.1109/TMTT.2003.809179>.
- [12] Kim J, Oh J, Lee H. Review on battery thermal management system for electric vehicles. *Appl Therm Eng* 2019;149:192–212. <https://doi.org/10.1016/j.applthermeng.2018.12.020>.
- [13] Sista S, Sista A. Intelligent BMS solution using AI and prognostic SPA. In: Proceedings of the FISITA 2012 world automotive congress. Berlin, Heidelberg: Springer; 2013. p. 755–64. [https://doi.org/10.1007/978-3-642-33741-3\\_4](https://doi.org/10.1007/978-3-642-33741-3_4).
- [14] Krizhevsky A, Sutskever I, Hinton GE. ImageNet classification with deep convolutional neural networks. *Commun ACM* 2017;60:84–90. <https://doi.org/10.1145/3065386>.
- [15] Howard AG, Zhu M, Chen B, Kalenichenko D, Wang W, Weyand T, Andreetto M, Adam H. MobileNets: efficient convolutional neural networks for mobile vision applications. <https://arxiv.org/abs/1704.04861>. [Accessed 22 August 2024].
- [16] Duan W, Song S, Xiao F, Chen Y, Peng S, Song C. Battery SOH estimation and RUL prediction framework based on variable forgetting factor online sequential extreme learning machine and particle filter. *J Energy Storage* 2023;65:107322. <https://doi.org/10.1016/j.est.2023.107322>.
- [17] Yang P, Yang HD, Meng XB, Song CR, He TL, Cai JY, Xie YY, Xu KK. Joint evaluation and prediction of SOH and RUL for lithium batteries based on a GBLS booster multi-task model. *J Energy Storage* 2024;75:109741. <https://doi.org/10.1016/j.est.2023.109741>.
- [18] Guo F, Wu X, Liu L, Ye J, Wang T, Fu L, Wu Y. Prediction of remaining useful life and state of health of lithium batteries based on time series feature and Savitzky-Golay filter combined with gated recurrent unit neural network. *Energy* 2023;270: 126880. <https://doi.org/10.1016/j.energy.2023.126880>.
- [19] Ge M-F, Liu Y, Jiang X, Liu J. A review on state of health estimations and remaining useful life prognostics of lithium-ion batteries. *Measurement* 2021;174:109057. <https://doi.org/10.1016/j.measurement.2021.109057>.
- [20] Degradation diagnostics for lithium ion cells. *J Power Sources* 2017;341:373–86. <https://doi.org/10.1016/j.jpowsour.2016.12.011>.
- [21] Zhang L, Sun C, He T, Jiang Y, Wei J, Huang Y, Yuan M. High-performance quasi-2D perovskite light-emitting diodes: from materials to devices. *Light Sci Appl* 2021; 10:61. <https://doi.org/10.1038/s41377-021-00501-0>.
- [22] Rechkemmer SK, Zang X, Zhang W, Sawodny O. Empirical Li-ion aging model derived from single particle model. *J Energy Storage* 2019;21:773–86. <https://doi.org/10.1016/j.est.2019.01.005>.
- [23] Luo K, Chen X, Zheng H, Shi Z. A review of deep learning approach to predicting the state of health and state of charge of lithium-ion batteries. *J Energy Chem* 2022;74:159–73. <https://doi.org/10.1016/j.jecchem.2022.06.049>.
- [24] Fei Z, Yang F, Tsui K-L, Li L, Zhang Z. Early prediction of battery lifetime via a machine learning based framework. *Energy* 2021;225:120205. <https://doi.org/10.1016/j.energy.2021.120205>.
- [25] Richardson RR, Osborne MA, Howey DA. Gaussian process regression for forecasting battery state of health. *J Power Sources* 2017;357:209–19. <https://doi.org/10.1016/j.jpowsour.2017.05.004>.
- [26] Li Y, Zou C, Berecibar M, Nanini-Maury E, Chan JC-W, van den Bossche P, Van Mierlo J, Omar N. Random forest regression for online capacity estimation of lithium-ion batteries. *Appl Energy* 2018;232:197–210. <https://doi.org/10.1016/j.apenergy.2018.09.182>.
- [27] Shi Y, Li J, Li Z. Gradient boosting with piece-wise linear regression trees. <https://doi.org/10.48550/arXiv.1802.05640>; 2019.
- [28] Ye Z, Kim MK. Predicting electricity consumption in a building using an optimized back-propagation and Levenberg–Marquardt back-propagation neural network: case study of a shopping mall in China. *Sustain Cities Soc* 2018;42:176–83. <https://doi.org/10.1016/j.scs.2018.05.050>.
- [29] Yang N, Song Z, Hofmann H, Sun J. Robust State of Health estimation of lithium-ion batteries using convolutional neural network and random forest. *J Energy Storage* 2022;48:103857. <https://doi.org/10.1016/j.est.2021.103857>.
- [30] Lee G, Kwon D, Lee C. A convolutional neural network model for SOH estimation of Li-ion batteries with physical interpretability. *Mech Syst Signal Process* 2023;188: 110004. <https://doi.org/10.1016/j.ymssp.2022.110004>.
- [31] Lu J, Xiong R, Tian J, Wang C, Sun F. Deep learning to estimate lithium-ion battery state of health without additional degradation experiments. *Nat Commun* 2023;14: 2760. <https://doi.org/10.1038/s41467-023-38458-w>.
- [32] Luo T, Liu M, Shi P, Duan G, Cao X. A hybrid data preprocessing-based hierarchical attention BiLSTM network for remaining useful life prediction of spacecraft lithium-ion batteries. *IEEE Trans. Neural Netw. Learning Syst* 2024;1–14. <https://doi.org/10.1109/TNNLS.2023.3311443>.
- [33] Peng S, Zhu J, Wu T, Yuan C, Cang J, Zhang K, Pecht M. Prediction of wind and PV power by fusing the multi-stage feature extraction and a PSO-BiLSTM model. *Energy* 2024;298:131345. <https://doi.org/10.1016/j.energy.2024.131345>.
- [34] Zhang Y, Wang Y, Zhang C, Qiao X, Ge Y, Li X, Peng T, Nazir MS. State-of-health estimation for lithium-ion battery via an evolutionary Stacking ensemble learning paradigm of random vector functional link and active-state-tracking long-short-term memory neural network. *Appl Energy* 2024;356:122417. <https://doi.org/10.1016/j.apenergy.2023.122417>.
- [35] Ma G, Zhang Y, Cheng C, Zhou B, Hu P, Yuan Y. Remaining useful life prediction of lithium-ion batteries based on false nearest neighbors and a hybrid neural network. *Appl Energy* 2019;253:113626. <https://doi.org/10.1016/j.apenergy.2019.113626>.
- [36] Xu H, Wu L, Xiong S, Li W, Garg A, Gao L. An improved CNN-LSTM model-based state-of-health estimation approach for lithium-ion batteries. *Energy* 2023;276: 127585. <https://doi.org/10.1016/j.energy.2023.127585>.
- [37] Zheng Y, Hu J, Chen J, Deng H, Hu W. State of health estimation for lithium battery random charging process based on CNN-GRU method. *Energy Rep* 2023;9:1–10. <https://doi.org/10.1016/j.egyr.2022.12.093>.
- [38] Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser L, Polosukhin I. Attention is all you need. <http://arxiv.org/abs/1706.03762>. [Accessed 22 August 2024].
- [39] Dong L, Xu S, Xu B. Speech-Transformer: a No-recurrence sequence-to-sequence model for speech recognition. In: 2018 IEEE International Conference on Acoustics, Speech and signal processing (ICASSP). Calgary, AB: IEEE; 2018. p. 5884–8. <https://doi.org/10.1109/ICASSP.2018.8462506>.
- [40] Devlin J, Chang M-W, Lee K, Toutanova K. BERT: pre-training of deep bidirectional Transformers for language understanding. <https://doi.org/10.48550/arXiv.1810.04805>; 2019.
- [41] Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, Dehghani M, Minderer M, Heigold G, Gelly S, Uszkoreit J, Houlsby N. An image is worth 16x16 words: Transformers for image recognition at scale. <http://arxiv.org/abs/2010.11299>. [Accessed 16 June 2024].
- [42] Hannan MA, How DNT, Lipu MSH, Mansor M, Ker PJ, Dong ZY, Sahari KSM, Tiong SK, Muttaqi KM, Mahlia TMI, Blaabjerg F. Deep learning approach towards accurate state of charge estimation for lithium-ion batteries using self-supervised transformer model. *Sci Rep* 2021;11:19541. <https://doi.org/10.1038/s41598-021-98915-8>.
- [43] Gomez W, Wang F-K, Chou J-H. Li-ion battery capacity prediction using improved temporal fusion transformer model. *Energy* 2024;296:131114. <https://doi.org/10.1016/j.energy.2024.131114>.
- [44] Jia C, Tian Y, Shi Y, Jia J, Wen J, Zeng J. State of health prediction of lithium-ion batteries based on bidirectional gated recurrent unit and transformer. *Energy* 2023; 285:129401. <https://doi.org/10.1016/j.energy.2023.129401>.
- [45] Gu X, See KW, Li P, Shan K, Wang Y, Zhao L, Lim KC, Zhang N. A novel state-of-health estimation for the lithium-ion battery using a convolutional neural network and transformer model. *Energy* 2023;262:125501. <https://doi.org/10.1016/j.energy.2022.125501>.
- [46] Bai T, Wang H. Convolutional transformer-based multiview information perception framework for lithium-ion battery state-of-health estimation. *IEEE Trans Instrum Meas* 2023;72:1–12. <https://doi.org/10.1109/TIM.2023.3300451>.
- [47] Chen L, Xie S, Lopes AM, Bao X. A vision transformer-based deep neural network for state of health estimation of lithium-ion batteries. *Int J Electr Power Energy Syst* 2023;152:109233. <https://doi.org/10.1016/j.ijepes.2023.109233>.
- [48] Birkil CR, Roberts MR, McTurk E, Bruce PG, Howey DA. Degradation diagnostics for lithium ion cells. *J Power Sources* 2017;341:373–86. <https://doi.org/10.1016/j.jpowsour.2016.12.011>.
- [49] Obisakun I, Ekeanyanwu CV. State of health estimation of lithium-ion batteries using support vector regression and long short-term memory. *Open J Appl Sci* 2022;12:1366–82. <https://doi.org/10.4236/ojapps.2022.128094>.
- [50] He W, Williard N, Osterman M, Pecht M. Prognostics of lithium-ion batteries based on Dempster-Shafer theory and the Bayesian Monte Carlo method. *J Power Sources* 2011;196:10314–21. <https://doi.org/10.1016/j.jpowsour.2011.08.040>.
- [51] Chen J, Hu Y, Zhu Q, Rashid H, Li H. A novel battery health indicator and PSO-LSSVR for LiFePO<sub>4</sub> battery SOH estimation during constant current charging. *Energy* 2023;282:128782. <https://doi.org/10.1016/j.energy.2023.128782>.
- [52] Wang F, Zhai Z, Zhao Z, Di Y, Chen X. Physics-informed neural network for lithium-ion battery degradation stable modeling and prognosis. *Nat Commun* 2024;15: 4332. <https://doi.org/10.1038/s41467-024-48779-z>.
- [53] Lin C, Xu J, Shi M, Mei X. Constant current charging time based fast state-of-health estimation for lithium-ion batteries. *Energy* 2022;247:123556. <https://doi.org/10.1016/j.energy.2022.123556>.
- [54] Peng S, Zhu J, Wu T, Tang A, Kan J, Pecht M. SOH early prediction of lithium-ion batteries based on voltage interval selection and features fusion. *Energy* 2024;308: 132993. <https://doi.org/10.1016/j.energy.2024.132993>.
- [55] Bai J, Huang J, Luo K, Yang F, Xian Y. A feature reuse based multi-model fusion method for state of health estimation of lithium-ion batteries. *J Energy Storage* 2023;70:107965. <https://doi.org/10.1016/j.est.2023.107965>.
- [56] Dai H, Wang J, Huang Y, Lai Y, Zhu L. Lightweight state-of-health estimation of lithium-ion batteries based on statistical feature optimization. *Renew Energy* 2024; 222:119907. <https://doi.org/10.1016/j.renene.2023.119907>.
- [57] Wu J, Liu Z, Zhang Y, Lei D, Zhang B, Cao W. Data-driven state of health estimation for lithium-ion battery based on voltage variation curves. *J Energy Storage* 2023; 73:109191. <https://doi.org/10.1016/j.est.2023.109191>.
- [58] Feng H, Zhang L. A heterogeneous learner fusion method with supplementary feature for lithium-ion batteries state of health estimation. *J Energy Storage* 2024; 92:111896. <https://doi.org/10.1016/j.est.2024.111896>.
- [59] Roman D, Saxena S, Robu V, Pecht M, Flynn D. Machine learning pipeline for battery state of health estimation. <http://arxiv.org/abs/2102.00837>. [Accessed 5 June 2024].
- [60] Lin M, Ke L, Wang W, Meng J, Guan Y, Wu J. Health prognosis via feature optimization and convolutional neural network for lithium-ion batteries. *Eng Appl Artif Intell* 2024;133:108666. <https://doi.org/10.1016/j.engappai.2024.108666>.
- [61] Tan Y, Zhao G. Transfer learning with long short-term memory network for state-of-health prediction of lithium-ion batteries. *IEEE Trans Ind Electron* 2020;67: 8723–31. <https://doi.org/10.1109/TIE.2019.2946551>.
- [62] Jin H, Cui N, Cai L, Meng J, Li J, Peng J, Zhao X. State-of-health estimation for lithium-ion batteries with hierarchical feature construction and auto-configurable Gaussian process regression. *Energy* 2023;262:125503. <https://doi.org/10.1016/j.energy.2022.125503>.
- [63] Huang K, Yao K, Guo Y, Lv Z. State of health estimation of lithium-ion batteries based on fine-tuning or rebuilding transfer learning strategies combined with new

- features mining. Energy 2023;282:128739. <https://doi.org/10.1016/j.energy.2023.128739>.
- [64] Kiranyaz S, Avci O, Abdeljaber O, Ince T, Gabbouj M, Inman DJ. 1D convolutional neural networks and applications: a survey. Mech Syst Signal Process 2021;151: 107398. <https://doi.org/10.1016/j.ymssp.2020.107398>.
- [65] Chollet F. Xception: deep learning with depthwise separable convolutions. In: 2017 IEEE Conference on computer vision and pattern recognition (CVPR); 2017. p. 1800–7. <https://doi.org/10.1109/CVPR.2017.195>.
- [66] Sandler M, Howard A, Zhu M, Zhmoginov A, Chen L-C. MobileNetV2: Inverted residuals and linear bottlenecks. In: 2018 IEEE/CVF Conference on computer vision and pattern recognition; 2018. p. 4510–20. <https://doi.org/10.1109/CVPR.2018.00474>.
- [67] Zhang X, Zhou X, Lin M, Sun J. ShuffleNet: an extremely efficient convolutional neural network for mobile devices. <https://doi.org/10.48550/arXiv.1707.01083>; 2017.
- [68] Ma N, Zhang X, Zheng H-T, Sun J. ShuffleNet V2: practical guidelines for efficient CNN architecture design. In: Ferrari V, Hebert M, Sminchisescu C, Weiss Y, editors. Computer vision – eccv 2018. Cham: Springer International Publishing; 2018. p. 122–38. [https://doi.org/10.1007/978-3-030-01264-9\\_8](https://doi.org/10.1007/978-3-030-01264-9_8).
- [69] Han D, Pan X, Han Y, Song S, Huang G. FLatten transformer: vision transformer using focused linear attention. In: 2023 IEEE/CVF International Conference on computer vision (ICCV). Paris, France: IEEE; 2023. p. 5938–48. <https://doi.org/10.1109/ICCV51070.2023.00548>.
- [70] Cai H, Li J, Hu M, Gan C, Han S. EfficientViT: lightweight multi-scale attention for high-Resolution dense prediction. In: 2023 IEEE/CVF International conference on computer vision (ICCV). Paris, France: IEEE; 2023. p. 17256–67. <https://doi.org/10.1109/ICCV51070.2023.01587>.
- [71] Chicco D, Warrens M, Jurman G. The coefficient of determination R-squared is more informative than SMAPE, MAE, MAPE, MSE and RMSE in regression analysis evaluation. PeerJ Computer Science 2021;7:e623. <https://doi.org/10.7717/peerj.cs.623>.
- [72] Li P, Zhang Z, Grosu R, Deng Z, Hou J, Rong Y, Wu R. An end-to-end neural network framework for state-of-health estimation and remaining useful life prediction of electric vehicle lithium batteries. Renew Sustain Energy Rev 2022; 156:111843. <https://doi.org/10.1016/j.rser.2021.111843>.