

## Research Papers

# A novel enhanced SOC estimation method for lithium-ion battery cells using cluster-based LSTM models and centroid proximity selection

Mohammed Khalifa Al-Alawi<sup>\*</sup>, Ali Jaddoa, James Cugley, Hany Hassanin

School of Engineering, Technology and Design, Canterbury Christ Church University, Canterbury, UK



## ARTICLE INFO

## Keywords:

State of charge (SOC) estimation  
Lithium-ion batteries  
LSTM  
Battery management systems (BMS)  
Dynamic SOC estimation  
Cluster-based learning model

## ABSTRACT

In line with the global mission in achieving the net zero target through deployment of renewable energy technologies and electrifying the transportation sector; precise and adaptable State of Charge (SOC) estimation for Lithium-ion batteries has emerged as a critical need. The paper introduces a novel Cluster-Based Learning Model (CBLM) framework that integrates the strengths of K-Means and Fuzzy C-Means clustering with the predictive power of Long Short-Term Memory (LSTM) networks. This approach aims to enhance the precision and reliability of battery SOC estimations, adapting to the dynamic and complex operational conditions characteristic of Li-ion batteries. The key contributions of this study are the development and validation of the CBLM framework, which was proven to outperform state-of-art standalone deep learning techniques particularly under diverse operational conditions. Additionally, the introduction of a centroid proximity selection mechanism within the CBLM framework, which dynamically selects the most appropriate cluster model in real-time based on the proximity of the operational data to the cluster centroids. The performance of the proposed CBLM approach is evaluated using a Tesla Model 3, 170 Li-ion battery dataset. Results demonstrate the model's enhanced performance, with reductions in Root Mean Square Error (RMSE) to as low as 0.65 % and Mean Absolute Error (MAE) to 0.51 %, reducing state-of-art benchmark model errors by margins of 61.8 % and 68.5 % respectively. Additionally, the maximum error using CBLM was lower than benchmark, emphasising the model's reliability in worst-case-scenarios. The study also conducted comprehensive ablation tests on the proposed novel framework to further optimize its performance.

## 1. Introduction

With the growing electrification of various sectors, including transportation, there is a rising demand for Lithium-ion (Li-ion) batteries. This was reflected by the International Energy Association's 2023 report which documented a 65 % increase in Li-ion battery demand within the automotive sector in 2022 compared to the previous year [1]. This surge is a result to the widespread adoption of electric and hybrid vehicles. In 2022, Europe emerged as a significant player in this market, accounting for about 23 % of the global demand for automotive Li-ion batteries, making it the second-largest market after China, which held nearly 57 % of the worldwide demand [1].

Furthermore, as the global shift towards renewable energy sources for electricity generation continues, with the goal of achieving net-zero emission targets set by governments, Battery Energy Storage Systems (BESS) have become increasingly crucial. As a key component of this shift, lithium-ion batteries are being extensively deployed in major

markets for a variety of applications. These include managing the intermittency of wind and solar photovoltaic (PV) energy generation, offering ancillary services to maintain grid stability, and other related uses [2]. This broad adoption is driven by the need to effectively integrate clean energy sources into the grid mix, ensuring a consistent and reliable energy supply while advancing towards environmentally sustainable goals.

The increasing reliance on lithium-ion batteries, while crucial for the transition to cleaner energy, also raise concerns regarding their safety and reliability. Ensuring these batteries operate within safe limits is essential to prevent issues like overcharging, thermal runaway, and degraded performance [3], which do not only pose safety risks but also affect the longevity and efficiency of the batteries. In addition to SOC's importance in ensuring longevity and reliable performance of batteries, misestimating SOC could lead to financial implications. BESS operators who provide services to the grid including frequency regulation, backup power or load levelling which are critical services for maintaining the

<sup>\*</sup> Corresponding author.

E-mail address: [ma867@canterbury.ac.uk](mailto:ma867@canterbury.ac.uk) (M.K. Al-Alawi).

<https://doi.org/10.1016/j.est.2024.112866>

Received 12 March 2024; Received in revised form 12 June 2024; Accepted 7 July 2024

Available online 16 July 2024

2352-152X/© 2024 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

stability of the electrical grid, are one of the key entities that could be impacted by the consequences of SOC misestimation. Entities in this market are contracted with the grid to provide predefined level of power output; therefore, underestimating or overestimating the batteries' SOC could result in underperformance leading to financial penalties and in extreme scenarios, it could lead to exclusion from market [4]. Similarly, incorrect estimation of SOC enforces faulty trading decisions which results in less optimal usage of battery assets and hence reduced revenue.

This necessitates precise monitoring and management of the battery's SOC, a key indicator of its health and operational status. SOC can be defined as the ratio of the current available capacity in relation to the battery's maximum capacity [5], as demonstrated in Eq. (1),

$$SOC(\%) = \left( \frac{Q_{current}}{Q_{maximum}} \right) \times 100 \quad (1)$$

where  $Q_{current}$  is the current charge of the battery in ampere-hours (Ah),  $Q_{maximum}$  is the maximum charge capacity of the battery, also in Ah and  $SOC(\%)$  is expressed as a percentage, indicating the current charge level relative to the battery's maximum capacity.

Unlike battery parameters like voltage, current and temperature, measuring the current charge of the battery while its operating is challenging, which necessitates utilising estimation methods instead [6]. Measurable battery parameters are of great help in terms of estimating SOC; yet understanding the relationship among battery parameters, including current, voltage, temperature, and SOC becomes a complex task due to the uncertainties in electrochemical and thermodynamic reactions lithium-ion batteries and their nonlinear dynamics [7]. M. A. Hannan et al. have discussed the challenges of SOC estimation in terms of operational conditions and states that since batteries in various applications are not always performing in the same charge/discharge rate which significantly affects SOC estimation [8].

This paper introduces a novel SOC estimation approach. The primary contribution of this study lies in the development and validation of a novel Cluster-Based Learning Model (CBLM) framework that integrates the strengths of K-Means and Fuzzy C-Means clustering with the predictive power of LSTM. This innovative approach enhances the precision and reliability of battery SOC estimations, adapting to the dynamic and complex operational conditions characteristic of lithium-ion batteries. The novelty of this research is further underscored by the introduction of a centroid proximity selection mechanism within the CBLM framework. This mechanism is instrumental in assigning the corresponding cluster model for real-time SOC estimation.

This paper introduces a novel SOC estimation approach. The primary contribution of this study lies in the development and validation of CBLM framework that integrates K-Means clustering with LSTM networks to reduce the margin of estimation errors resulting of the traditional SOC estimation techniques and the advanced data driven methods, particularly under diverse operational conditions. Although clustering combined with neural networks has been investigated in other research domains, such as medical imaging, these implementations were not designed for real-time, regression-based estimation tasks. In fields like medical imaging, clustering is used for classification or segmentation tasks, which differ fundamentally from the continuous, real-time prediction required for SOC estimation of lithium-ion batteries. Additionally, a dynamic model assignment mechanism that selects the most appropriate cluster model in real-time based on the proximity to pre-defined directory of cluster centroids. This mechanism dynamically selects the most appropriate cluster model in real-time based on the proximity of the operational data to the cluster centroids. By continuously adapting to the closest cluster, this approach significantly improves the adaptability and precision of SOC estimation, especially under rapidly changing operational conditions.

The remainder of this paper is structured as follows: Section 2 reviews existing SOC estimation methods and the use of CBLM in various applications. Section 3 outlines the methodology, including dataset

preparation, clustering, and LSTM model development. Section 4 presents the experimental results, such as optimal cluster identification and performance evaluation of the proposed CBLM. Section 5 discusses the findings, advantages, limitations and results of robustness tests (ablation test) of the CBLM approach. The conclusion summarises the contributions and proposes future research directions.

## 2. Literature review

A wide range of methodologies have been explored for estimating the SOC of batteries, including approaches grounded in characteristic parameters, ampere-hour integration, data-driven models, and other model-based techniques. The precision of SOC estimation is pivotal, as it directly influences the optimization of battery performance and plays a critical role in ensuring the reliability of energy storage systems. According to [9], SOC estimation methods can be classified broadly to the following: Coulomb Counting [10], Look Up Table [11], Model Based [3,12], Data Driven [13] and Hybrid Models [14]. Coulomb counting and direct measurement techniques have been commonly used because of their simplicity. However, these approaches can accumulate errors and struggle to estimate with high precision under various operating conditions. The Coulomb counting method is widely used in industry however it is liable to significant errors due to initial misestimation of SOC and sensor inaccuracies [15]. Equivalent circuit models, such as Thevenin's model, provide a more detailed view of battery dynamics but face challenges in identifying parameters and are sensitive to noise [16]. In addition, the Extended Kalman Filter has been widely used for SOC estimation because it can deal with nonlinearities and uncertainties. However, EKF-based approaches need precise battery models and assume Gaussian noise, which might not be realistic in real-life situations [17]. This could result in notable estimation errors. Likewise, the combination of Unscented Kalman Filter with machine learning models has been employed to improve SOC estimation accuracy, but these methods may have limitations when applied in real-time scenarios due to their complexity [18].

Particularly, data-driven methods have earned significant attention due to their ability to adapt to the complex, non-linear behaviour of batteries, combined with the advantage of avoiding complex modelling and initial parameter identification, offering enhanced accuracy in SOC estimation [19]. Among the spectrum of data-driven techniques, neural network models have emerged as particularly promising in SOC estimation. According to a comprehensive review conducted by [20], which highlighted the benefits and drawbacks of data-driven SOC methods; neural network methods emerge as particularly promising for SOC estimation in lithium-ion batteries due to their advantages in dealing with the dynamic and non-linear nature of battery behaviour and successful operation under long-term dependencies, making them more suitable for SOC estimation compared to other methods like Decision Forests (DF) or Support Vector Machines (SVM). Numerous SOC estimation methods for lithium-ion batteries leveraging neural networks have been proposed [21–24], reflecting a growing trend towards utilising these advanced computational models for their capability of handling the complexities of battery behaviour.

The need for advanced SOC estimation methods that can adapt to diverse and real-world operational conditions is underscored by recent literature. [7,25] have emphasised the limitations of existing approaches that often rely on fixed charging and discharging currents, which do not accurately reflect actual battery usage. These observations suggest the need for models that can account for time-varying currents to enhance the applicability and accuracy of neural network methods in SOC estimation.

[13] conducted a study to estimate SOC using a Long Short-Term Memory (LSTM) neural network and compared its performance against other neural networks architecture; the findings showed that the LSTM model exhibited superior accuracy in SOC estimation, achieving <0.62 % of maximum standard error (MSE), outperforming Deep-Feed-

Forward Neural Network (DFNN) and Feed-Forward Neural Network (FFNN), which recorded MSEs of 5.37 % and 4.03 %, respectively. This leads to a conclusion that LSTM models are highly effective for SOC estimation; however, authors also recommended further investigation into the performance of the model across challenging operational conditions. A study by [26] introduced a Bi-LSTM neural network for accurate SOC estimation in lithium-ion batteries. Using a dataset at 0 °C, 10 °C, and 25 °C, the Bi-LSTM model demonstrated superior accuracy, with MAEs of 0.498 %, 0.411 %, and 0.738 %, and an overall MAE of 0.616 % across temperatures. The results, including a maximum error of 3.8 %, highlight the model's robustness and its significant improvement over existing methods like GRU networks in SOC estimation. Furthermore, [27] proposed a method that combines Temporal Convolutional Network (TCN) and LSTM referred to as the TCN-LSTM model, that is aimed to extract spatial features and long-term dependencies from battery data. The model demonstrated superior performance across various conditions, with an MAE of 0.48 %, RMSE of 0.60 %, and maximum error of 2.3 % under multiple temperature scenarios, outperforming standalone LSTM, TCN, and CNN-LSTM networks. SOC estimation has significantly improved with the use of machine learning and hybrid methods, but limitations related to real-time applicability, robustness, and adaptability under diverse conditions remain. According to [15], despite the major advances, SOC data-driven estimation models still face limitations in generalisability across different operational conditions.

Clustering-based learning models have demonstrated promising results in a diverse array of research fields, highlighting their adaptability and effectiveness. These models combine clustering techniques with advanced learning algorithms and have been applied to solve complex problems. This approach has led to significant improvements in accuracy and efficiency, as it tailors the learning process to the specific characteristics of each cluster. Such advancements in clustering-based learning models offer promising methodologies that can be adapted to enhance SOC estimation in lithium-ion batteries particularly with challenges of SOC estimation in diverse operational conditions.

[28] presents a method for short-term power load forecasting for larger consumers, leveraging LSTM deep learning framework K-Means clustering algorithm. The approach involves clustering users based on electrical attributes to create "load curves" for each cluster, representing different types of consumers. Compared with the traditional LSTM algorithm, K-Means LSTM model has significantly improved load forecasting, reducing the MAE by 5.12 %. Another significant study by [29] introduced an Improved multi-stage LSTM network with Clustering (ILSTMC) model which integrates the benefits of K-Means clustering and LSTM networks to enhance RUL prediction of aero-engines with higher accuracy. The proposed model integrates the benefits of clustering analysis (K-Means clustering) and LSTM networks to enhance Remaining Useful Life (RUL) prediction. The results show that, in the last stage of prediction, the ILSTMC model achieves a 0.85 % reduction in RMSE compared to LSTM. On average, across all stages, the RMSE reduction is 1.87 %, and the accuracy of the life cycle prediction is improved by 0.59 % over LSTM, with an average improvement of 1.84 % at each stage.

Fusion of clustering algorithms and advanced learning algorithms have not only enhanced regression tasks as discussed but also improved the performance of classification capabilities. [30] applied K-Means clustering and an improved deep learning model for diagnosing common corn leaf diseases. The pre-processing of images with clustering markedly improved model performance in disease classification and diagnosis. Extending to medical applications, [31] proposed a novel method for brain tumor segmentation and classification, combining K-Means clustering with deep learning. Their method, using a finetuned VGG19 model, achieved improved accuracy in tumor classification, showcasing the potential of clustering and deep learning in medical image analysis. Collectively, these studies showcased significant potential of clustering-based learning models in diverse domains of research, ranging from power forecasting to medical image analysis, in dealing with complex

data-driven tasks with high effectiveness and adaptability. The proven ability of these models in identifying patterns within segmented clusters and apply targeted learning strategies to each could potentially solve a key challenge in SOC estimation which is dealing with diverse operational conditions and usage patterns and the non-linear behaviour of Li-ion batteries.

This literature review confirms the increasing importance of accurate SOC estimation in lithium-ion batteries which is considered a critical element for the reliability of modern energy storage systems. The evidence shows that data-driven approaches, especially neural networks like the LSTM models, are particularly effective in understanding battery behaviour with no to minimal feature engineering. Combining these models with clustering techniques is a promising development, providing a clearer insight into different battery operational states. This study builds on this knowledge by introducing a CBLM that leverages the detailed analytics of LSTM and the data segmentation of clustering algorithms.

### 3. Methodology

This section describes the sequence of procedures employed in this study, as illustrated in Fig. 1. Initially, data cleaning eliminates overlapping timestamps, ensuring data integrity. Subsequently, data partitioning segments the dataset for training and validation, ensuring robust model training and effective validation. Data clustering then cluster the dataset into distinct groups, followed by data pre-processing, which prepares the data for the LSTM model. The model development phase constructs the LSTM network to process the sequential data and capture the battery's SOC behaviours. The final stage, dynamic estimation, applies a Centroid Proximity Mechanism using the trained models for each cluster for accurate SOC estimation. Further details of each step are explored in Sections 3.1 to 3.6.

#### 3.1. Dataset

This study utilised data obtained from the 4.5 Ah 'm80' cell of a Tesla Model 3, 170 lithium-ion battery, Nickel Cobalt Aluminium (NCA) chemistry, collected under 0 °C ambient conditions at McMaster University [32]. Testing was conducted in a 16 ft<sup>3</sup> Envirotronics SH16C thermal chamber with a temperature control accuracy of  $\pm 0.3$  °C, using an eight-channel, 60 A/channel Arbin cell cycler that ensures precise control over voltage and current.

The dataset encompasses a wide array of characterisation tests, including controlled discharge rates and Hybrid Pulse Power Characterisation (HPPC), aimed at capturing the fundamental properties of the battery. Additionally, various driving cycles were executed, including well-established ones such as UDDS, HWFET, LA92, and US06, as well as custom-designed cycles. Additionally, randomized cycles were introduced to simulate real-world driving patterns. This comprehensive set of both standard and specialized tests offers a comprehensive view of battery performance, facilitating the development of an accurate SOC estimation model capable of estimating SOC under a range of diverse operational conditions and usage patterns.

Table 1 provides a comparative overview of studies that used machine learning algorithms for battery SOC estimation that have specifically used the same dataset as our study. This includes various data-driven techniques employing neural that are directly comparable to our proposed CBLM framework. The table contrasts these studies based on their employed estimation techniques and key innovative contributions.

#### 3.2. Data cleaning

This stage involved the transformation of MATLAB (.mat) files into structured DataFrames using Python's Pandas library. Each DataFrame included key parameters such as time, voltage, current, SOC, and battery

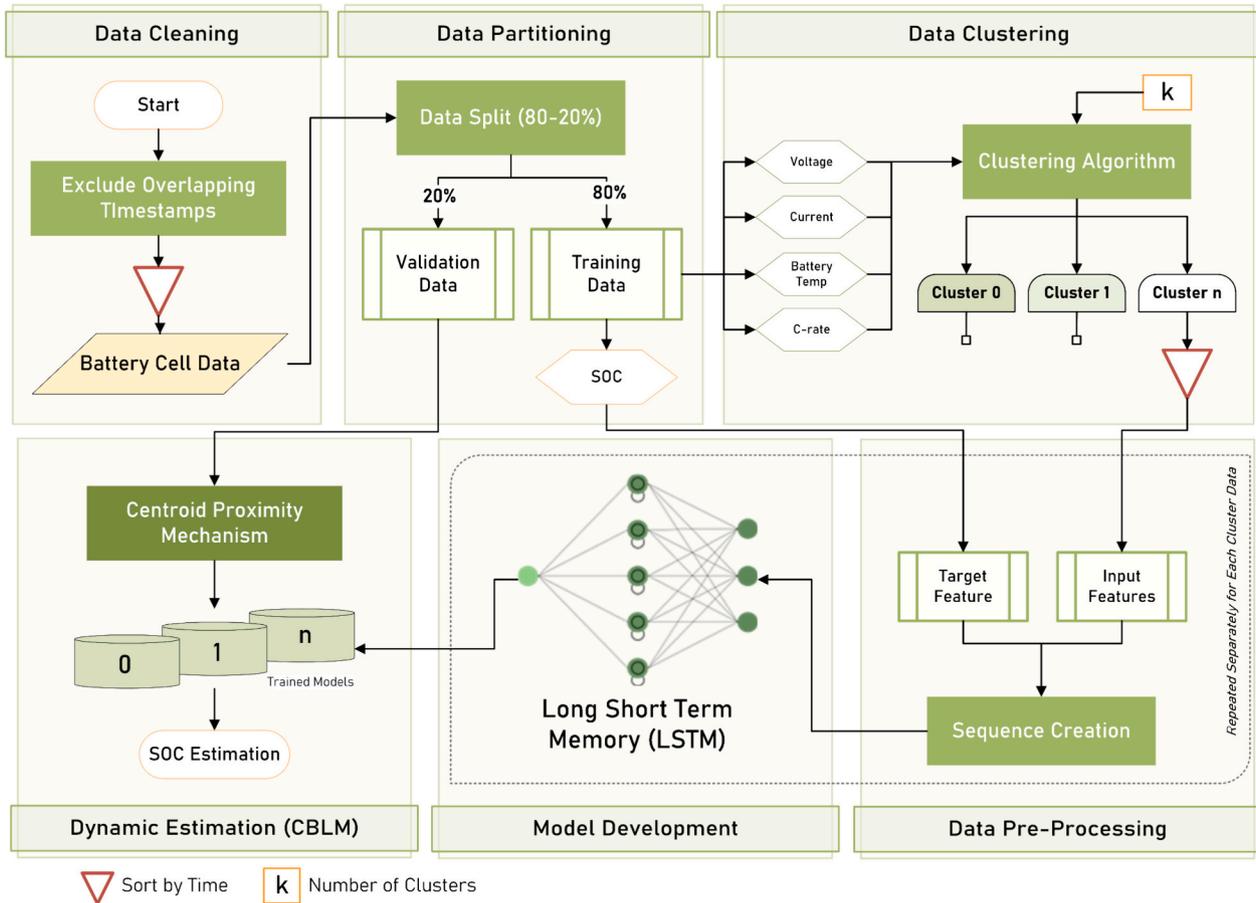


Fig. 1. Flowchart of the CBLM methodology for SOC using LSTM and centroid proximity.

Table 1

Comparative overview of machine learning approaches for battery SOC estimation using the same dataset.

Study	Data driven using neural networks	Algorithm	Clustering	Contribution
[33]	✓	FFNN and Nonlinear Autoregressive Exogenous (NARX) network	×	Demonstrated the benefit of using NARX for its robustness against current sensor errors.
[34]	✓	Multi-layer perceptron (MLP)	×	Highlighted the effectiveness of MLP in SOC Estimation
This Study	✓	LSTM	✓	Novel integration of clustering with LSTM and dynamic centroid proximity mechanism

temperature, in addition to the calculation of C-rate from the current data, which is It is a measure that indicates how fast a battery is being charged or discharged in relation to its maximum capacity. During the data processing phase, the presence of overlapping timestamps across different test files necessitated the selective exclusion of certain files to maintain data integrity and consistency. Specific filtering conditions were applied to address these overlaps. For the file named ‘8\_C\_20\_Discharge\_Charge\_10-03-21\_13.24’, data points were retained only if the ‘Time’ value was less than or equal to 330,507. Conversely, for the file ‘9\_CC\_CV\_charge\_10-11-21\_03.21’, data points were kept only if the ‘Time’ value was >330,507. These measures were important in ensuring that only non-overlapping, relevant data are included. Despite this, the retained data provides a robust foundation for the analysis of the SOC behaviour of the m80 cell under the specified ambient temperature condition, covering diverse operational scenarios. Fig. 2 demonstrates that thorough data cleaning was implemented ensuring continuity and validates that the C-rate distribution is variable and represents diverse operational conditions.

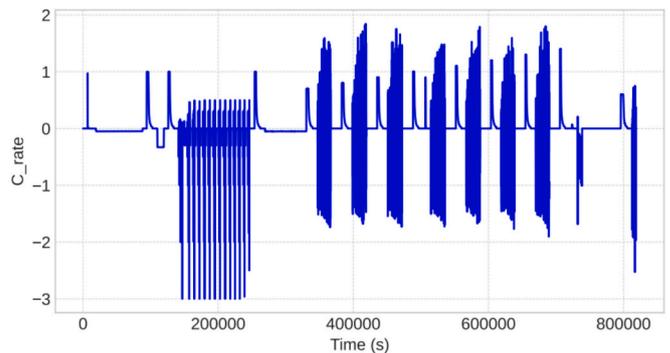


Fig. 2. C-rate over time for the battery data used.

### 3.3. Data partitioning

In order to establish a robust validation framework for the proposed cluster-based learning model, a discontinuous mixed sampling (DMS) method was employed for partitioning the cleaned dataset into training and validation sets. The rationale behind partitioning the data is to use training data for the development of the model, while the validation data are kept aside for the final stage, to assess the performance of the dynamic estimation using the cluster-based learning model. DMS reflects a more rigorous test environment, presenting the model with non-sequential data to evaluate its robustness against sudden and unexpected changes in working conditions as demonstrated in Fig. 3, where validation data is represented in red and training data is represented in blue.

### 3.4. Clustering

Clustering is an essential component in the proposed cluster-based learning model for SOC estimation to segment battery data into insightful clusters for specialized learning. The choice of K-Means clustering algorithm as a candidate was based on several factors; mainly the widespread use of the algorithm and its effectiveness in providing meaningful segments that resulted in enhanced estimation for both regression and classification tasks in diverse research domains as discussed in the literature. Furthermore, K-Means is well-known for its simplicity and low computational complexity requirements, making it a suitable algorithm to implement the proposed approach [35]. The algorithm assigns each datapoint into one distinct cluster, based on its proximity to the centroid of the cluster, ensuring clear and non-overlapping groups of battery data. For comparative analysis and to address the nature of battery data complexity, Fuzzy C-Means (FCM) clustering was also implemented, which is a soft clustering algorithm. Implementing the novel SOC estimation approach using the mentioned clustering algorithms offer a more detailed analysis of battery behaviour.

A fundamental step was determining the optimal number of clusters, 'k', for both algorithms as this directly influences the quality and effectiveness of the clustering process. Therefore, specific cluster quality assessment techniques were employed tailored for each algorithm.

#### a. K-Means Cluster Quality Metrics [36]:

- Within-Cluster Sum of Squares (WCSS): measures cluster compactness, by quantifying the variance within each cluster with lower WCSS values indicating tighter clustering.

- Davies-Bouldin (DB): evaluates the average similarity between clusters, where a lower value suggests better separation between clusters.
- Calinski-Harabasz indices (CH): assesses the ratio of between-cluster variance to within-cluster variance, where higher values are indicative of well-separated and distinct clusters.

For optimal choice of number of clusters, the average score is calculated using Eq. (2), where lowest average score indicates better clustering performance,

$$\text{Average Score} = \frac{\text{Normalised (WCSS + DB + Inverted CH)}}{3} \quad (2)$$

#### b. FCM Cluster Quality Metrics:

- Fuzzy Partition Coefficient (FPC): evaluates the clarity of cluster boundaries, with higher FPC values denoting more distinctly defined clusters [37].
- Fuzzy Entropy Coefficient (FPE): provides insight into the randomness in grouping data points, where lower FPE values signify well-structured and reliable clusters [37].
- Xie-Beni: as discussed by [38], it focuses on assessing cluster separation and compactness, striving for the smallest XB value for optimal clustering.

For optimal choice of number of clusters, the average score is calculated using Eq. (3), where lowest average score indicates better clustering performance,

$$\text{Average Score} = \frac{\text{Normalised (Inverted FPC + FPE + XB)}}{3} \quad (3)$$

For clustering, raw battery features including Voltage, Current, Battery Temperature and C\_rate are used. Following successful clustering, cluster labels are added to the original data frame and data are segmented to the respective cluster accordingly.

### 3.5. Learning model development

LSTM networks are a type of recurrent neural network (RNN) designed to capture and learn long-term dependencies within sequential input data [39]. These networks have demonstrated their effectiveness in a variety of applications including 3D human action recognition [40], rainfall-runoff modelling [39], and stock price prediction [41]. The LSTM model was chosen to implement the proposed SOC estimator. In

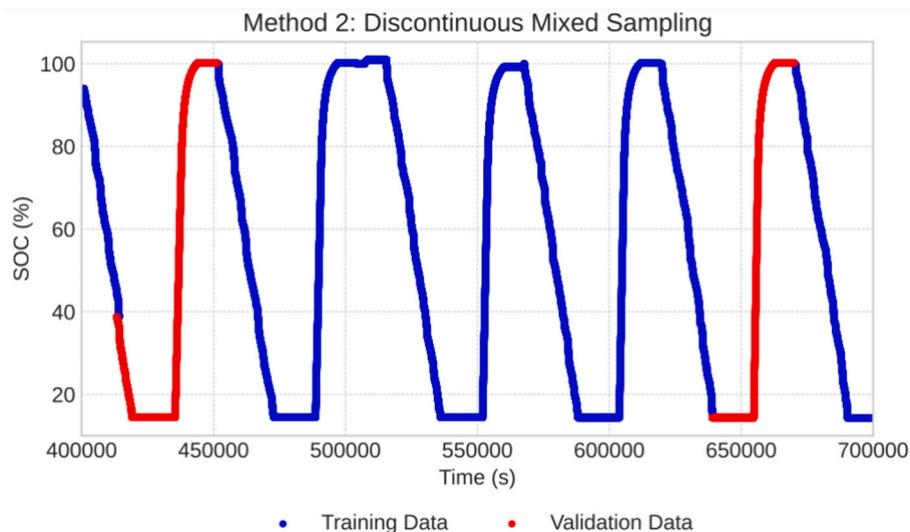


Fig. 3. SOC variation over time for a selected window highlighting training and validation data.

contrast to traditional RNNs, LSTM networks are equipped with gating mechanisms that control the flow of data, allowing them to selectively remember or forget specific information over long sequences [42]. This functionality underpins the LSTM's notable accuracy in SOC estimation tasks, as corroborated by the literature. This decision was based on the LSTM's proven high accuracy in SOC estimation, as evidenced in the literature. Additionally, the LSTM architecture is capable of handling the temporal dependencies and complex, non-linear nature of battery behaviour, making it a suitable choice for accurately modelling SOC dynamics [43].

The hyperparameters used to build and train the LSTM models are presented in Table 2. Research has proved that neural networks with a single hidden layer can achieve robust performance in SOC estimation [44]. Additionally, findings by [45] demonstrated that increasing the number of LSTM hidden layers led to reduced SOC estimation accuracy, highlighting the trade-off between model complexity and performance. Therefore, the model developed in this study has one LSTM hidden layer with 50 neurons, providing a balance between model complexity and computational demand. To avoid model overfitting, 0.2 was set as the dropout rate. Furthermore, the training parameters were selected to align with the available computational resources without compromising on the training performance. These hyperparameters were consistent across all LSTM models for the CBLM and the Benchmark models.

### 3.6. Dynamic estimation with cluster-based learning model

#### 3.6.1. Cluster assignment mechanism

The Cluster Assignment Mechanism is an essential component of the proposed SOC estimation model, facilitating dynamic and real-time SOC prediction. This mechanism operates based on a directory that stores the centroids for each cluster, which represent the average values of the features (Voltage, Current, Temperature, and C-rate) for the data points within each cluster. The selection of the appropriate cluster model for a given test datapoint is determined through the computation of the Euclidean distance between the datapoint and each cluster's centroid using Eq. (4)

$$\text{Euclidean Distance} = \sqrt{(V_i - V_c)^2 + (I_i - I_c)^2 + (T_i - T_c)^2 + (\text{Crate}_i - \text{Crate}_c)^2} \quad (4)$$

where,  $V_i$  and  $V_c$  are the voltage of the test datapoint and the voltage centroid of the cluster, respectively.  $I_i$  and  $I_c$  are the current of the test datapoint and the current centroid of the cluster, respectively.  $T_i$  and  $T_c$  are the temperature of the test datapoint and the temperature centroid of the cluster, respectively.  $\text{Crate}_i$  and  $\text{Crate}_c$  are the C-rate of the test datapoint and the C-rate centroid of the cluster, respectively.

This distance metric quantifies the similarity between a given test datapoint and the characteristic centre of each cluster. The test datapoint is then assigned to the cluster whose centroid is closest, i.e., the cluster with the minimum Euclidean Distance to the datapoint. This nearest centroid effectively determines the most similar operational condition represented by the cluster, thereby selecting the most appropriate cluster model for SOC estimation. The dynamic nature of this mechanism allows the model to adapt to changing operational

**Table 2**  
Hyperparameter settings of the LSTM models.

Type	Parameter	Description	Setting
Network structure	Hidden layers	Number of LSTM layers	1
	Neurons	Number of neurons in the LSTM layer	50
	Dropout rate	Rate of dropout to prevent overfitting	0.2
Training process	Optimizer	Algorithm for optimization	'adam'
	Epochs	Number of complete iterations through the dataset	20
	Batch size	The batch size for training	32

conditions, ensuring high accuracy and reliability in SOC estimation.

#### 3.6.2. Performance evaluation metrics

To rigorously evaluate the performance of the proposed CBLM for SOC estimation, three key metrics were employed: Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and Maximum Error (Max Error). Each metric provides a unique perspective on the model's accuracy and reliability, essential for validating the effectiveness of the proposed approach. RMSE is a standard metric used to measure the model's accuracy by calculating the square root of the average squared differences between the actual and predicted SOC values as defined in Eq. (5),

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\text{SOC}_i - \widehat{\text{SOC}}_i)^2} \quad (5)$$

where,  $n$  represents the number of observations,  $\text{SOC}_i$  is the actual SOC value and  $\widehat{\text{SOC}}_i$  is the estimated SOC value for the  $i$ th observation. RMSE is a useful metric in this case as it gives a relatively high weight to large errors, ensuring the model's precision across the entire battery dataset. MAE quantifies the average magnitude of SOC estimation errors as defined in Eq. (6),

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |\text{SOC}_i - \widehat{\text{SOC}}_i| \quad (6)$$

It is a key metric in evaluating the performance of the SOC estimation model as it provides insights on the consistency of the estimation performance without being overly influenced by few large errors. The Max Error identifies the largest single error in SOC estimation as defined in Eq. (7),

$$\text{Max Error} = \text{Max}(|\text{SOC}_i - \widehat{\text{SOC}}_i|) \quad (7)$$

It is another key metric as it assesses the reliability of the estimation model, particularly under extreme or worst-case conditions, ensuring the model's capability in maintaining accuracy even in challenging instances.

## 4. Experimental results and settings

### 4.1. Settings

#### 4.1.1. Optimal number of clusters ( $k$ ) identification

The WCSS metric shows a decreasing trend as the number of clusters increases, which is indicative of improved compactness within clusters. The significant drop from  $k = 2$  to  $k = 3$  suggests a substantial improvement in cluster compactness, making these cluster sizes of particular interest. The DB Score, which assesses the separation between clusters, is lowest (indicating better separation) at  $k = 3$ , suggesting that three clusters achieve a good balance between separation and compactness. The gradual increase in the DB Score beyond indicates diminishing returns in terms of cluster separation quality. The CH Score, which evaluates the validity of clustering by comparing within-cluster dispersion to between-cluster dispersion, shows an improvement as the number of clusters increases from  $k = 2$  to  $k = 4$ . This suggests that four clusters might offer an optimal balance between compactness and separation. The transition from  $k = 5$  to  $k = 6$  resulted in a reduction in the CH score, indicating better clustering validity (Fig. 4).

The FPC Scores remain high at  $k = 2$  and  $k = 3$ , indicating a strong degree of cluster membership certainty for these cluster sizes. The decrease in FPC Scores as the number of clusters increases beyond  $k = 4$  suggests that the clarity of cluster membership diminishes, making  $k \leq 4$  of particular interest. The FPE Scores, which assess the overall performance of the clustering, are lowest (indicating better performance) at  $k \leq 4$ . This aligns with the FPC Scores and suggests that these cluster sizes are optimal in terms of performance. The Xie-Beni Score, which measure

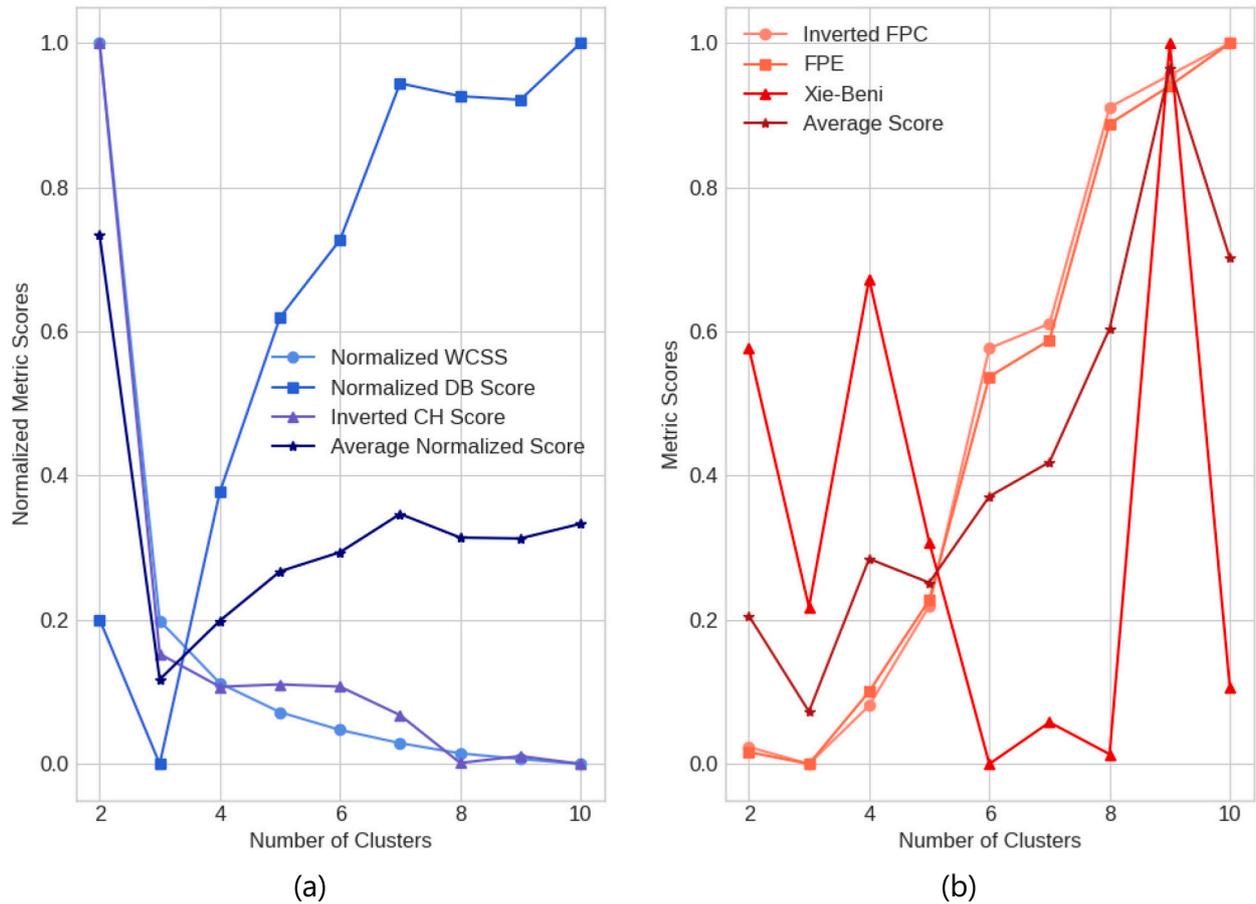


Fig. 4. Normalized clustering metrics for optimal number of clusters selection for (a) K-Means and (b) FCM, ranging from k = 2 to k = 10. The figure compares cluster quality metrics for each clustering algorithm, each with a merged average score.

the compactness and separation of fuzzy clusters, was the lowest at k = 6 suggesting better-defined clusters.

Therefore, based on the quality metrics findings of both clustering algorithms, the study selects k = 2, k = 3, k = 4 and k = 6 for further investigation as potential optimal clusters. For K-Means, the selection is driven by the balance between cluster compactness and separation, with k = 3 and k = 4 showing optimal cluster validity. For FCM, the focus on k = 2 and k = 3 is justified by the high degree of cluster membership certainty and performance, with k = 4 also considered due to its relatively strong performance across metrics. The inclusion of k = 6 for further investigation acknowledges the potential for insights into the limitations and challenges of clustering at higher cluster sizes, particularly for the FCM algorithm.

#### 4.1.2. Cluster distribution visualization

This section presents a detailed visualization of cluster distributions derived from the K-Means clustering algorithm applied to our SOC estimation model. Each cluster represents a unique operational condition of the battery, differentiated by characteristics such as charging behaviour, temperature variability, and C-rate variability. Through the use of box plots for each cluster across different values of k = 2, 3, 4, and 6, it is aimed to provide a visual interpretation of how the battery's operational states are segmented by the clustering process. This visualization serves as a pivotal step towards understanding the relationship between cluster characteristics and SOC estimation accuracy, ultimately guiding the selection of an optimal k value for precise SOC prediction.

Fig. 5 illustrates the distribution of C-rate and temperature across different clusters defined by the K-Means clustering algorithm for

various k-values. The box plots in this figure provide a visual representation of the operational conditions of battery usage, such as charging rates and temperature ranges of each cluster as segmented by the clustering algorithm. These plots display the mean ( $\mu$ ), standard deviation ( $\sigma$ ), and range of conditions within each cluster, emphasising the diversity and similarity of operational patterns across different clusters. Table 1 synthesises the analysis of central tendencies ( $\mu$ ) and variability ( $\sigma$ ) clusters offering a clearer understanding of the specific battery behaviours in each cluster. These interpretations correlated along with the SOC estimation performance of the CBLM models would provide valuable insights on the impact of operational state segmentation on the accuracy and reliability of the estimation models.

## 4.2. Results

This section presents the SOC estimation results using the proposed cluster-based LSTM models for various number of clusters, k = 2,3,4 and 6, for both clustering algorithms, namely, K-MEANS and FCM. The performance of each cluster-based LSTM model is evaluated using key metrics including RMSE, MAE and Max Error and compared against the benchmark model.

### 4.2.1. K-Means

Fig. 6 shows the estimated SOC against actual SOC over time (s) for K-MEANS-LSTM CBLM model with the corresponding cluster assignment subplot for different values of k. It is evident that the variation in value of k during cluster impacts the estimation accuracy; for instance, the performance of the model k = 4 was able to maintain a better SOC

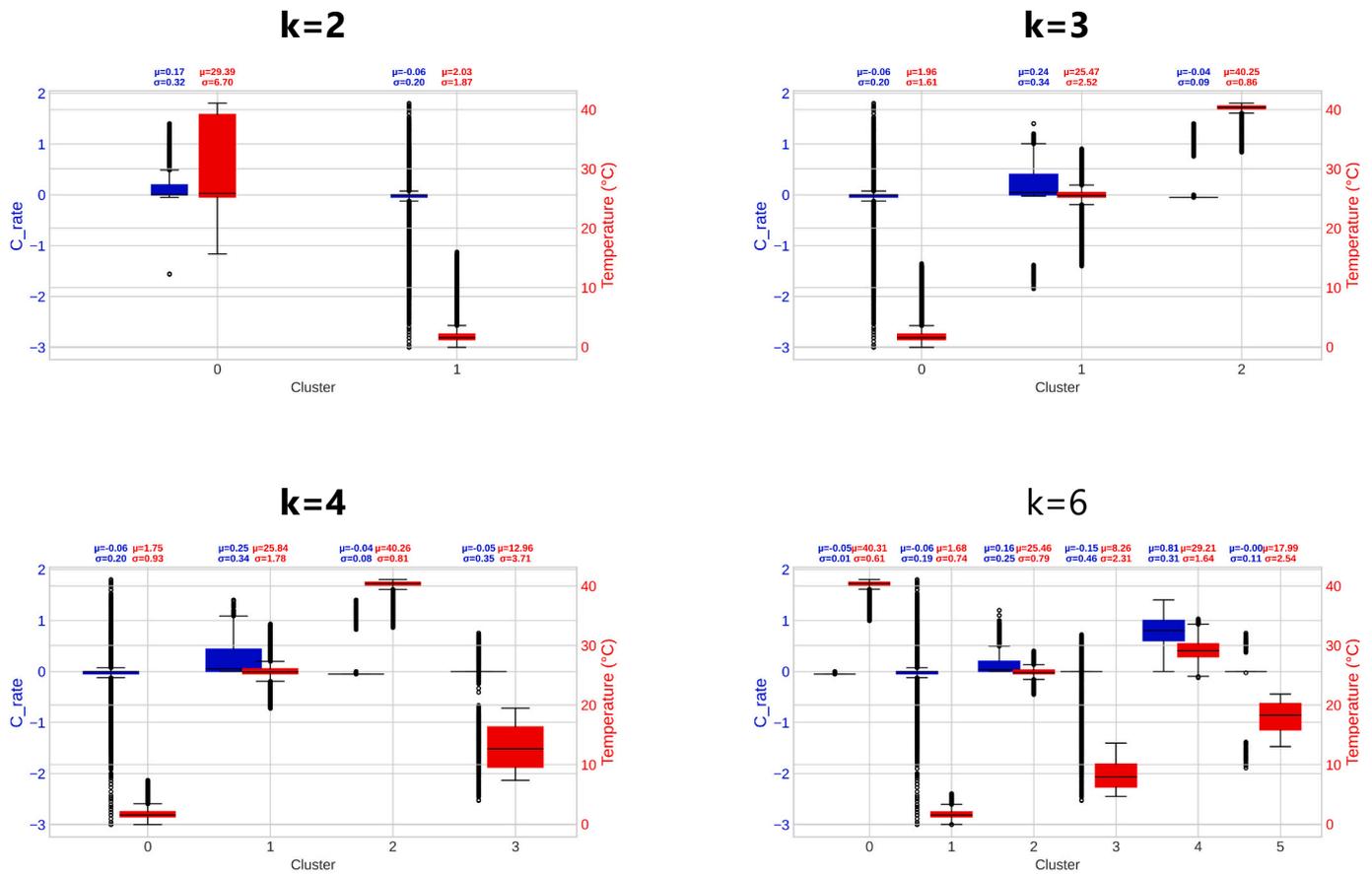


Fig. 5. Cluster distribution box plots for varying k values.

estimation during the steady state (>6000 s) as in Fig. 6.iii.a in comparison to other models. Additionally, the performance of CLBM with  $k = 6$  is significantly better in estimating SOC with less estimation errors in comparison to CLBM with  $k = 3$ . For a better understanding of the models accuracy and reliability in SOC estimation, Table 4 shows the values of errors for each CBLM and compared against benchmark LSTM model.

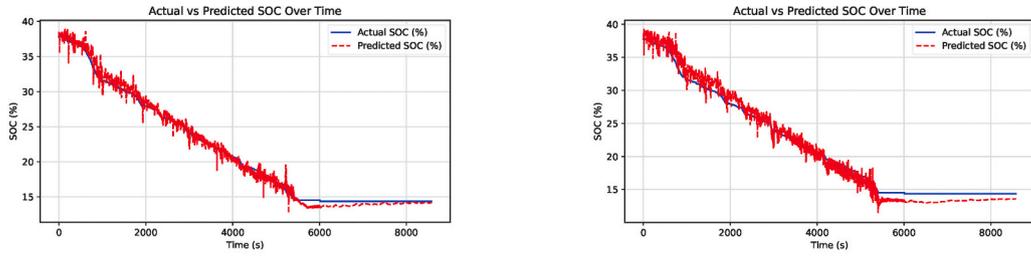
Notably, the RMSE values for K-Means CBLM show a general increase with the number of clusters, starting from 0.65 % at  $k = 2$  and peaking at 0.75 % at  $k = 6$ , except for a slight increase at  $k = 3$  to 1.05 %. This suggests that the model’s predictive accuracy is generally stable across different cluster sizes, with a notable exception at  $k = 3$ , where the error increases. This superior performance is consistent, in which each CBLM model across all  $k$  values demonstrated lower RMSE value compared to the benchmark, confirming that clustering is capable of enhancing the precision of SOC estimation. Similar to RMSE, the MAE increases from 0.51 % at  $k = 2$  to 0.63 % at  $k = 6$ , with a peak at  $k = 3$  (0.93 %). This indicates a consistent estimation across all cluster sizes, with  $k = 3$  showing the largest deviation from actual values. Additionally, CBLM model with  $k = 4$  also performed well in average with an MAE of 0.53 %, outperforming benchmark’s MAE of 1.62 %. To further evaluate the robustness of the approach, it is necessary to assess the model’s ability in maintaining accuracy in extreme working conditions. The CBLM with  $k = 4$  has the lowest maximum error of 3.65 % which is substantially lower than the benchmark’s maximum error of 6.59 %.

Although the CLBM with all  $k$  values outperformed the benchmark performance in each evaluation metric; the models were only using only one cluster model for the estimation as demonstrated in the cluster assignment subplots in Fig. 6 and this is not sufficient for a comprehensive assessment of the proposed approach and necessities extended tests.

#### 4.2.2. FCM

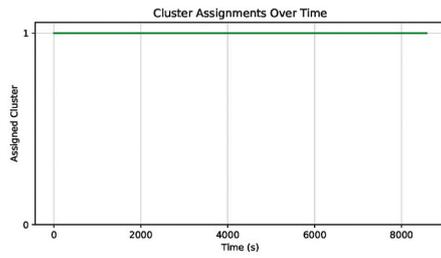
Fig. 7 shows the estimated SOC against actual SOC over time (s) for FCM-LSTM CBLM model with the corresponding cluster assignment subplot for different values of  $k$ . It is evident that the choice of the number of clusters ( $k$ ) markedly influences the accuracy of SOC estimation. Models with lower  $k$  values exhibit fewer spikes in SOC estimation, indicating a smoother estimation curve, whereas CBLM with higher  $k$  values tend to show increased fluctuation in SOC estimation. Although CBLM with  $k = 6$  performance showed increased number of spikes in terms of SOC estimation; it was able to maintain better estimation accuracy in steady state (>6000 s) compared to models with lower  $k$  values. Table 3 presents error metrics for each CBLM, offering insights into the models’ precision and consistency in SOC estimation, with a comparative analysis against the Benchmark LSTM model.

The RMSE for FCM CBLM is lowest at  $k = 3$  (0.64 %), significantly outperforming the benchmark’s RMSE of 1.70 %, and highest at  $k = 6$  (4.71 %), indicating a significant loss in estimation accuracy at higher cluster numbers. This suggests that the model performs best with a low to moderate number of clusters, with performance degrading significantly as clusters increase. MAE results convey that  $k = 3$  and  $k = 2$  models performed with superior consistency, that is confirmed by the lowest MAE figures of 0.47 % and 0.51 % respectively, compared to the benchmark’s 1.62 %. In terms of the maximum error the lower  $k$  models consistently exhibit a lower error margin compared to the benchmark, reinforcing the robustness of lower  $k$  FCM CBLM in extreme scenarios. On the other hand, the  $k = 6$  model’s maximum error at 79.06 % exposes its considerable limitations in maintaining prediction accuracy under extreme operational scenarios and this is a result of the rapid transition between cluster models as shown in Fig. 7.iv, as models with lower  $k$  were only using one cluster model for estimation. Overall, these findings highlight that FCM clustering can improve SOC estimation, yet selecting

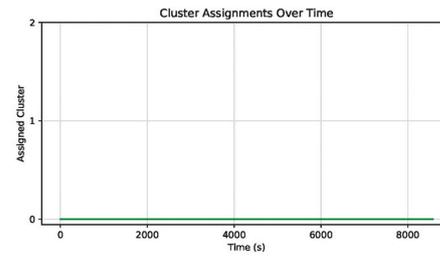


(a)

(a)



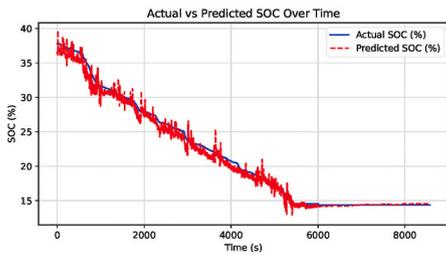
(b)



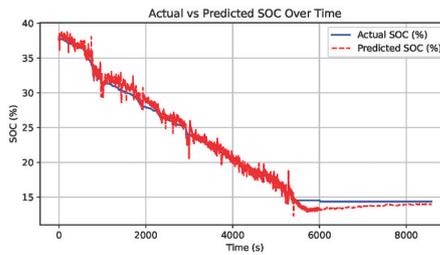
(b)

*i: K-means CBLM with k=2 (a) Estimated SOC vs Actual and (b) Cluster assignments*

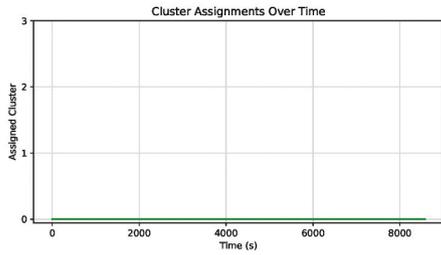
*ii: K-means CBLM with k=3 (a) Estimated SOC vs Actual and (b) Cluster assignments*



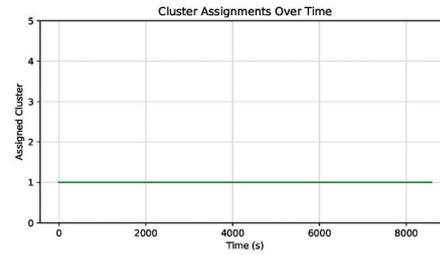
(a)



(a)



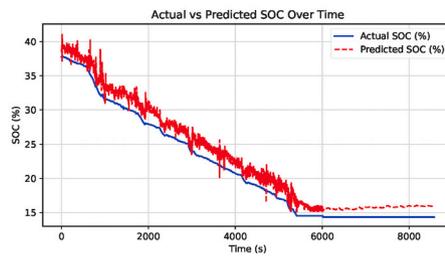
(b)



(b)

*iii: K-means CBLM with k=4 (a) Estimated SOC vs Actual and (b) Cluster assignments*

*iv: K-means CBLM with k=6 (a) Estimated SOC vs Actual and (b) Cluster assignments*



*v: Benchmark model*

**Fig. 6.** Actual vs. estimated SOC over time and corresponding cluster assignments for CLBM  $k = 2$ ,  $k = 3$ ,  $k = 4$ , and  $k = 6$  using K-Means algorithm compared against Benchmark model performance.

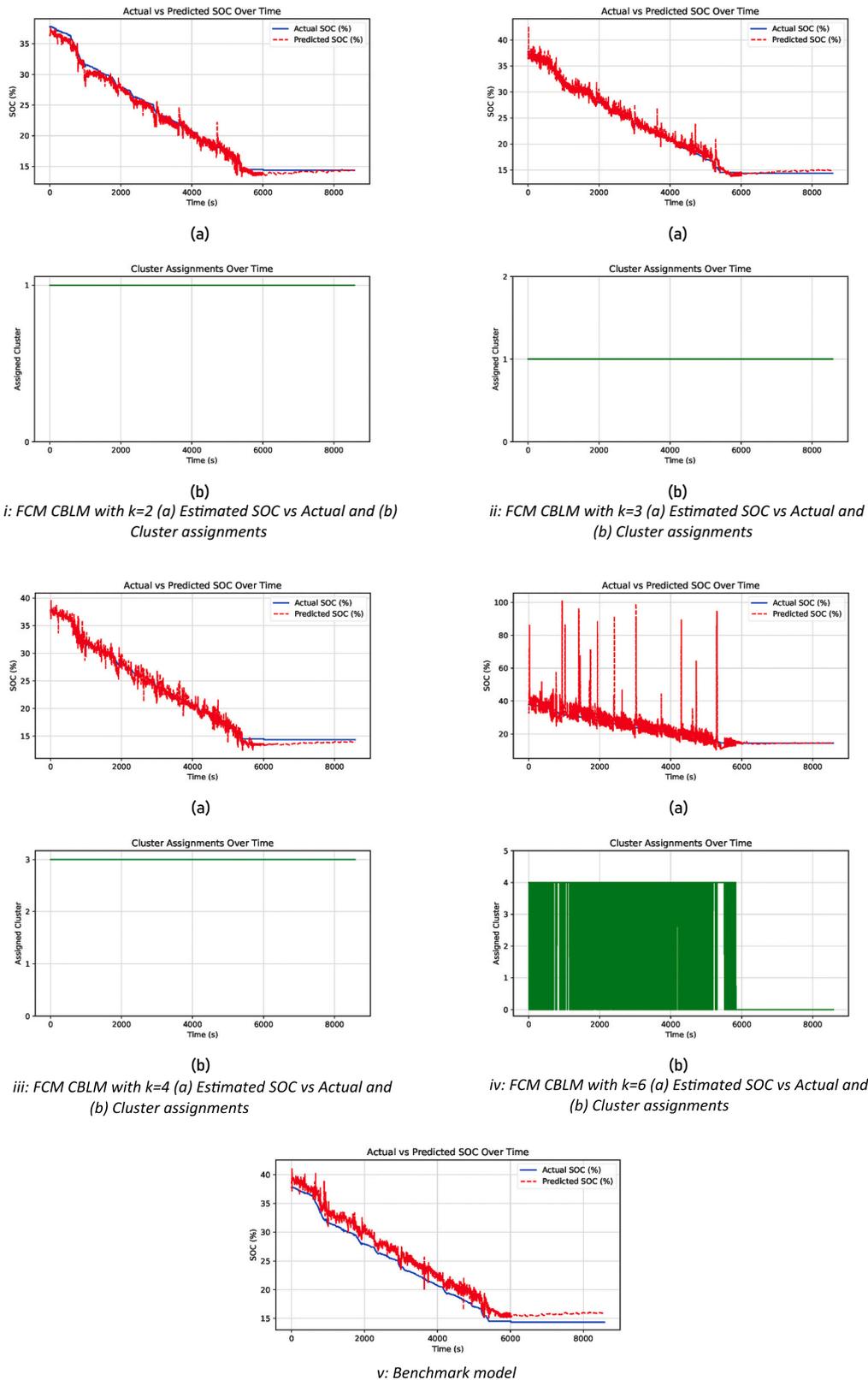


Fig. 7. Actual vs. estimated SOC over time and corresponding cluster assignments for CLBM  $k = 2$ ,  $k = 3$ ,  $k = 4$ , and  $k = 6$  using FCM algorithm compared against Benchmark model performance.

**Table 3**  
Interpretation of cluster distribution plots in terms of C-rate and battery temperature.

k	Cluster	K-Means	FCM
k = 2	0	Moderate charging, higher temps.	Moderate charging, higher temps.
	1	Low variability discharging, near ambient temps.	Low variability discharging, significantly higher temps.
k = 3	0	Low variability discharging, near ambient temps.	Low variability discharging, moderate temps.
	1	Moderate variability charging, warmer temps.	High variability discharging, wide temp. range.
	2	Low to moderate variability discharging, near ambient temps.	Mixed charging/discharging, elevated temps.
k = 4	0	Low variability discharging, just above ambient temps.	Charging, slightly above ambient temps.
	1	Moderate variability charging, warmer temps.	High variability discharging, broad temp. range.
	2	Very low C-rate, much higher than ambient temps.	Lower variability charging, high temps.
	3	Moderate variability discharging, moderately above ambient temps.	Low variability, likely standby, lower temps.
k = 6	0	Very low variability discharging, significantly higher temps.	Minimal C-rate variability discharging, higher temps.
	1	Low variability discharging, close to ambient temps.	Balanced charge/discharge, moderate temps.
	2	Moderate variability charging, warmer temps.	Varied use discharging, broader temp. range.
	3	High variability discharging, moderately above ambient temps.	Active/high-rate charging, higher temps.
	4	Moderate variability charging, higher temps.	Intensive discharge rates, high temps.
	5	Balanced charge/discharge, above ambient temps.	Erratic/extreme conditions, high temps.

**Table 4**  
Comparative performance evaluation, K-Means.

Metric	Cluster based learning model				Benchmark
	k = 2	k = 3	k = 4	k = 6	
RMSE (%)	0.65	1.05	0.72	0.75	1.70
MAE (%)	0.51	0.93	0.53	0.63	1.62
Max error (%)	3.62	4.64	3.56	3.91	6.59

**Table 5**  
Comparative performance evaluation, FCM CBLM.

Metric	Cluster based learning model				Benchmark
	k = 2	k = 3	k = 4	k = 6	
RMSE (%)	0.66	0.64	0.70	4.71	1.70
MAE (%)	0.51	0.47	0.58	2.27	1.62
Max error (%)	3.98	5.60	4.52	79.06	6.59

an optimal number of clusters is crucial, as lower k values tend to yield more dependable and accurate results (Table 5).

**5. Discussion**

Both K-Means and FCM CBLM findings indicate that the model’s estimation accuracy varies with the number of clusters. However, FCM shows a more pronounced deterioration in performance at the highest number of clusters, suggesting that K-Means may provide more stable predictions across a wider range of cluster sizes.

The K-Means CBLM showcases a remarkable improvement in accuracy and consistency over the benchmark. With a RMSE reduced by 61.8 % and MAE by 68.5 % for the k = 2 configuration, it is evident that the K-Means CBLM is significantly superior in estimating SOC. This high level of precision, combined with a 45 % reduction in maximum error, underscores the K-Means CBLM’s ability to provide reliable estimations under varied conditions. The success of the K-Means model, particularly with fewer clusters, highlights its efficacy in capturing the non-linear behaviour of lithium-ion batteries, striking an optimal balance between simplicity and the complexity needed for precise SOC estimation.

Comparing K-Means and FCM models directly, the K-Means CBLM,

**Table 6**  
Comparative performance evaluation, K-MEANS CBLM–Extended test.

Metric	Cluster based learning model				Benchmark
	k = 2	k = 3	k = 4	k = 6	
RMSE (%)	1.24	1.39	1.31	1.16	1.64
MAE (%)	0.72	0.85	0.91	0.73	1.46
Max error (%)	7.79	9.82	16.18	5.42	6.60

especially with k = 2, emerges as the more robust model for SOC estimation. This is attributed to its greater accuracy, lower error rates, and consistency across estimation. While the FCM model with k = 3 comes close in terms of RMSE and MAE improvements over the benchmark, the K-Means model’s overall stability and reliability make it the preferable choice. The significant increase in maximum error observed in the FCM k = 6 model further highlights the K-Means model’s superiority in handling diverse and dynamic battery behaviours without compromising on estimation accuracy. Additionally, the FCM CBLM performance across all k values exhibited undesirable spikes as demonstrated in Fig. 7; although there was no transition in cluster models for the lower values of k which puts forward a huge concern on its performance under a variety of operational conditions where transitions are required.

Despite the strengths of the K-Means CBLM, the reliance on a single cluster model for estimation highlights a potential limitation in its adaptability to rapidly changing operational states. This underscores the importance of further testing under a variety of conditions to fully assess the dynamic SOC estimation capabilities of the K-Means CBLM SOC estimation model, which is critical for real-time SOC estimation.

**5.1. Extended test for K-Means CBLM**

To comprehensively evaluate the robustness of the K-Means CBLM in SOC estimation and its ability to dynamically transition between cluster models, extended testing was conducted under a variety of operational conditions. The performance of the models for different values of k is demonstrated in Fig. 8.

For k = 2, the model demonstrates a high degree of accuracy, with the predicted SOC closely following the actual SOC throughout the duration, including the steady state and dynamic transitions. The cluster assignment graph indicates that the model is relying on both cluster models which shows successful implementation of the proposed dynamic approach. For k = 3, the model continues to perform well, but introduces slightly more variability in the SOC estimation, as evidenced

**Table 7**  
CBLM performance evaluation under the test: Removing one feature at a time.

Scenario	MAE (%)	RMSE (%)	MaxError (%)	Finding
All features (original)	0.37	0.47	3.05	Lowest error
Remove current	0.58	0.71	5.57	Increased error
Remove voltage	20.93	27.21	80.51	Increased error
Remove battery temperature	0.81	1.88	22.68	Increased error
Remove C-rate	0.36	0.46	4.16	Similar performance to original scenario, increased max error

**Table 8**  
CBLM performance evaluation under the test: Changing model hyperparameters.

Case 1: Changing learning rate					
Configuration	Learning rate	MAE (%)	RMSE (%)	MaxError (%)	Finding
1	0.0001	0.55	0.71	6.30	Increased error
2 (original)	0.001	0.37	0.47	3.05	Lowest error
3	0.01	0.92	1.38	14.22	Increased error
Case 2: Changing optimizer					
Configuration	Optimizer	MAE (%)	RMSE (%)	MaxError (%)	Finding
1	SGD	2.56	4.18	12.51	Increased error
2	NAdam	0.42	0.53	4.46	Increased error
3 (original)	Adam	0.37	0.47	3.05	Lowest error
4	RMSprop	1.42	1.57	4.91	Increased error
Case 3: Training epochs					
Configuration	Epochs	MAE (%)	RMSE (%)	MaxError (%)	Finding
1	10	0.41	0.57	4.46	Increased error
2 (original)	20	0.37	0.47	3.05	Lowest error
3	50	0.45	0.62	20.06	Increased error
Case 4: Number of hidden neurons					
Configuration	Neurons	MAE (%)	RMSE (%)	MaxError (%)	Finding
1	20	0.64	0.80	3.71	Increased error
2 (original)	50	0.37	0.47	3.05	Lowest error
3	100	0.60	0.77	5.64	Increased error
Case 5: Output layer activation function					
Configuration	Activation function	MAE (%)	RMSE (%)	MaxError (%)	Finding
1	Sigmoid	0.44	0.55	3.75	Increased error
2 (original)	Linear	0.37	0.47	3.05	Original performance
3	ReLU	0.32	0.42	2.39	Lowest error
4	PReLU	0.38	0.49	4.00	Increased error
5	MISH	0.32	0.41	3.44	Reduced error

by the mild fluctuations. These fluctuations are mainly present at the periods of transitions between cluster models, which introduces complexity but also allows for a finer distinction between different operating conditions. With  $k = 4$ , the SOC estimation remains closely aligned with the actual SOC, even as the model navigates through more frequent transitions between clusters. This suggests that while the model can capture a greater range of operating conditions, it may also become more vulnerable to fluctuations due to rapid changes in cluster assignments. The  $k = 6$  model displays a significant increase in accuracy compared to  $k = 4$ , which was evident in the shorter spikes which could

be a result of the non-presence of rapid change to different cluster models.

Evidently, there is an increase in the number of spikes and estimation errors during the dynamic transitions between clusters. This may indicate that while the model can finely distinguish between numerous conditions, the frequent switching between clusters can lead to instability in SOC predictions. In summary, these results suggest that while increasing the number of clusters can provide a more detailed representation of battery behaviour, it also raises the complexity of the model's estimation task, which can lead to less stable predictions during periods of rapid operational changes. This analysis underscores the need to balance the granularity of clustering with the model's ability to maintain accurate and stable SOC predictions across a range of conditions. Table 6 provides focused insights on the error margins of the models compared to the benchmark LSTM model.

RMSE percentages suggest that the CBLM with  $k = 6$  achieved the lowest RMSE of 1.16 %, which is an improvement over the benchmark's RMSE of 1.64 %. MAE percentages further reinforce the superiority of the  $k = 6$  model with the lowest MAE of 0.73 %, compared to the benchmark's MAE of 1.46 %. This lower MAE indicates more consistent SOC estimation from the  $k = 6$  model. Interestingly, the  $k = 2, 3$  and  $4$  models also outperform the benchmark with an MAE of 0.72 %, 0.85 % and 0.91 % respectively, which is closely aligned with the  $k = 2$  model's MAE, suggesting that despite the higher complexity and more frequent cluster transitions, the  $k = 6$  model can still maintain a reliable level of precision in average. Notably, when examining the Maximum Error, all CBLM configurations demonstrate a higher maximum error than the benchmark, except for CBLM with  $k = 6$ .

From this analysis, it is evident that the  $k = 6$  model has the best overall performance across all metrics suggesting an optimal balance between accuracy and consistency, with lower maximum error in comparison to benchmark model. Reflecting on the cluster data distribution interpretations presented in Table 6, the superior performance of the  $k = 6$  configuration can be justified by its detailed segmentation of battery behaviour K-means algorithm with  $k = 6$  was able to clearly separate minor differences in charging and discharging behaviours, along with a detailed awareness of temperature changes and operational conditions. Such granularity allowed the model to accurately estimate battery SOC across a range of working conditions, enhancing the estimator's precision robustness as indicated by the lowest errors presented in Table 6.

Although all CBLM configurations performed significantly better in terms of MAE in comparison to Benchmark; CBLM with  $k = 2, 3$  and  $4$  exhibited increased maximum error which is majorly due to the transition between different cluster models. These insights should be taken into consideration when choosing a CBLM for practical applications, as the operational context might dictate the preference for either consistent average performance or resilience against maximum prediction errors.

Additionally, the robustness of the centroid proximity-based cluster assignment process presented in this research for dynamic SOC estimation must be emphasised. The findings have shown that this process is successful in assigning cluster models for estimation under a range of operational conditions. However, certain limitations have become apparent, particularly the undesired estimation spikes caused by rapid transitions between cluster models. This issue points to the necessity for enhancing the current method to better manage these transitions, ensuring smoother SOC estimations and overall model reliability. A potential advancement of the proposed method is including a threshold of transition between one cluster model to another to reduce the chances of potential rapid transitions, particularly instances of momentarily switching to an alternative cluster model, only to revert to the previous one in short time period as noticed in the cluster assignment sublots of  $k = 3$  and  $k = 4$  in Fig. 8.

The benchmark's model performance with the highest RMSE and MAE and relatively higher maximum error highlights the critical role of granularity provided by the CBLM approach in understanding complex



develop high performance models using deep learning algorithms as in [46–48]. The ablation tests are aimed to evaluate the impact of individual input features on the estimation accuracy, and the hyperparameter tuning experiments explore the effect of various model configurations on the overall performance. All experiments presented in this section were fully carried out on the  $k = 2$  CBLM configuration for simplicity. To challenge the models, the test set included both a Constant Current Constant Voltage (CCCV) charge and a drive cycle to present diverse operational scenarios. The testing data used for this robustness analysis was different from the data employed for the rest of the experiments in the manuscript.

### 5.2.1. Remove one feature at a time

This experiment was carried out by systematically removing one input feature at a time from the original feature set that included current, voltage, battery temperature, and C-rate. The aim is to identify importance and impact of individual features on the SOC estimation using CBLM to highlight any redundant features that do not significantly contribute to enhanced performance. The results of this ablation test are presented in Table 7.

Results in Table 7 clearly highlight the critical importance of the voltage feature for accurate SOC estimation using the CBLM. Removing the voltage feature led to a significant increase in all error metrics (MAE, RMSE, and MaxError), indicating that voltage is an essential feature in this context. The battery temperature feature also played an important role, as its removal resulted in an increase in all metrics, particularly leading to higher errors. While the current feature's removal had a slight increase in the errors, removing the C-rate feature had a minimal impact on MAE and RMSE, although it slightly increased the MaxError. This suggests that the C-rate feature may be less crucial than the others for SOC Estimation using the proposed approach.

### 5.2.2. Changing model hyperparameters

This study was carried out to optimize the performance of CBLM which involved varying key hyperparameters of the model and assessing their effects on the model's performance. The hyperparameters investigated in this study include Learning Rate, Optimizer, Training Epochs, Hidden Neurons and Output Layer Activation Function. The results of this hyperparameter tuning study are presented in Table 8.

The hyperparameter tuning study revealed the optimal configurations for the CBLM's performance. The original learning rate of 0.001 and the Adam optimizer outperformed other settings, indicating their suitability for this task. Regarding the number of training epochs, the original setting of 20 epochs achieved the lowest errors, suggesting that either fewer or more epochs could lead to underfitting or overfitting, respectively. The original configuration with 50 hidden neurons achieved the best performance, indicating an appropriate level of complexity. Fewer neurons limited the model's capacity, while more neurons resulted in increased errors that could be due to overfitting. Finally, the ReLU activation function in the output layer achieved the lowest errors, outperforming the original linear activation, resulting in reduced errors overall with particular emphasis on the reduction in maximum error which is highly important in the context of Lithium-ion battery SOC estimation.

## 6. Conclusion

In conclusion, the primary contribution of this study is the development and validation of a novel CBLM framework that integrates the strengths of K-Means and Fuzzy C-Means clustering with the predictive power of LSTM networks. This innovative approach enhances the precision and reliability of battery SOC estimations, adapting to the dynamic and complex operational conditions characteristic of lithium-ion batteries. The integration of a dynamic SOC estimation process in the proposed CBLM framework, facilitated by a centroid proximity, The paper also introduces a centroid proximity cluster model selection

mechanism within the CBLM framework, which is instrumental in assigning the corresponding cluster model for real-time SOC estimation of batteries.

The results demonstrated that the K-Means CBLM, particularly with  $k = 2$  and  $k = 6$  clusters, outperformed the benchmark standalone LSTM model in terms of RMSE, MAE, and maximum error, showcasing the effectiveness of the proposed approach. The findings also highlighted the importance of selecting an optimal number of clusters, as the performance can vary depending on the granularity of the operational condition representation. The study conducted comprehensive ablation tests to assess the robustness of the proposed CBLM framework and evaluated the impact of individual input features and the effect of various model hyperparameters on the estimation accuracy.

The study has made significant contribution to SOC estimation, but it acknowledges certain limitations. Rapid transitions between cluster models, particularly with FCM CBLM, could lead to undesirable estimation spikes and indicate a need for refinement in the cluster assignment mechanism. To address these limitations and enhance the SOC estimation process further, future research should focus on refining the CBLM SOC estimation by improving the centroid proximity cluster assignment mechanism for smoother transitions between cluster models, exploring dynamic switching mechanisms to enhance model accuracy during different operational. Additionally, we plan to further assess the generalisability of the proposed CBLM framework by applying it to datasets involving different battery chemistries. This will allow us to evaluate the robustness and adaptability of the approach across a wider range of battery technologies. Building upon the optimal cluster identification investigation conducted in this study, future research could focus on developing a more automated mechanism for selecting the optimal number of clusters.

## CRedit authorship contribution statement

**Mohammed Khalifa Al-Alawi:** Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Methodology, Investigation, Formal analysis, Conceptualization. **Ali Jaddoa:** Writing – review & editing, Supervision, Resources, Methodology, Investigation. **James Cugley:** Writing – review & editing, Supervision. **Hany Hassanin:** Writing – review & editing, Supervision.

## 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.

## Data availability

Used open access data

## References

- [1] Global EV outlook 2023: trends in batteries, Available: <https://www.iea.org/reports/global-ev-outlook-2023/trends-in-batteries>.
- [2] M.K. Al-Alawi, J. Cugley, H. Hassanin, Techno-economic feasibility of retired electric-vehicle batteries repurpose/reuse in second-life applications: a systematic review, Energy and Climate Change 3 (2022) 100086, <https://doi.org/10.1016/j.egycc.2022.100086>. Available: <https://www.sciencedirect.com/science/article/pii/S2666278722000162>.
- [3] L. Wu, et al., Physics-based battery SOC estimation methods: recent advances and future perspectives, Journal of Energy Chemistry 89 (2024) 27–40, <https://doi.org/10.1016/j.jechem.2023.09.045>.
- [4] The financial implications of inaccurate SOC in LFP batteries, Available: <https://www.accure.net/battery-knowledge/lfp-soc-estimation-challenges> (Dec 14).
- [5] S. Wang, et al., Multidimensional Lithium-ion Battery Status Monitoring, First edition, CRC Press, Boca Raton; London; New York, 2023.
- [6] M. Naguib, P. Kollmeyer, A. Emadi, Lithium-ion battery pack robust state of charge estimation, cell inconsistency, and balancing: review, Access 9 (2021) 50570–50582, <https://doi.org/10.1109/ACCESS.2021.3068776>. Available: <https://ieeexplore.ieee.org/document/9386065>.

- [7] Z. Cui, et al., A comprehensive review on the state of charge estimation for lithium-ion battery based on neural network, *Int. J. Energy Res.* 46 (5) (2022) 5423–5440, <https://doi.org/10.1002/er.7545>. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/er.7545>.
- [8] M.A. Hannan, et al., A review of lithium-ion battery state of charge estimation and management system in electric vehicle applications: challenges and recommendations, *Renew. Sustain. Energy Rev.* 78 (2017) 834–854, <https://doi.org/10.1016/j.rser.2017.05.001>.
- [9] D.V.S.R.S., C. Badachi, R.C. Green II, A review on data-driven SOC estimation with Li-Ion batteries: implementation methods & future aspirations, *Journal of Energy Storage* 72 (2023) 108420, <https://doi.org/10.1016/j.est.2023.108420>.
- [10] K.S. Ng, et al., Enhanced coulomb counting method for estimating state-of-charge and state-of-health of lithium-ion batteries, *Appl. Energy* 86 (9) (2009) 1506–1511, <https://doi.org/10.1016/j.apenergy.2008.11.021>.
- [11] H. He, et al., Online model-based estimation of state-of-charge and open-circuit voltage of lithium-ion batteries in electric vehicles, *Energy (Oxford)* 39 (1) (2012) 310–318, <https://doi.org/10.1016/j.energy.2012.01.009>.
- [12] W. S. Rui Xiong, *Advanced Battery Management Technologies for Electric Vehicles*. (1st ed.) Newark: John Wiley & Sons, Ltd, 2019.
- [13] E. Almaita, et al., State of charge estimation for a group of lithium-ion batteries using long short-term memory neural network, *Journal of Energy Storage* 52 (2022) 104761, <https://doi.org/10.1016/j.est.2022.104761>.
- [14] F. Feng, et al., Co-estimation of lithium-ion battery state of charge and state of temperature based on a hybrid electrochemical-thermal-neural-network model, *J. Power Sources* 455 (2020) 227935, <https://doi.org/10.1016/j.jpowsour.2020.227935>.
- [15] I.B. Espedal, et al., Current trends for state-of-charge (SoC) estimation in lithium-ion battery electric vehicles, *Energies (Basel)* 14 (11) (2021) 3284, <https://doi.org/10.3390/en14113284>. Available: <https://search.proquest.com/docview/2539699219>.
- [16] S.V. Kishore N, V.S. Sravan Kumar, *Comparative Analysis of Model-based Approaches for State-of-charge Estimation in Batteries*, Nov 24, 2022, pp. 1–6.
- [17] S. Sunil, B. Balasingam, K.R. Pattipati, *State-of-charge Estimation of Batteries Using the Extended Kalman Filter: Insights Into Performance Analysis and Filter Tuning*, Dec 9, 2022, pp. 1–6.
- [18] Y. Zeng, Y. Li, T. Yang, State of charge estimation for lithium-ion battery based on unscented Kalman filter and long short-term memory neural network, *Batteries (Basel)* 9 (7) (2023) 358, <https://doi.org/10.3390/batteries9070358>. Available: <https://doaj.org/article/33466da870d64b1f8996e2764d055080>.
- [19] W. Zhou, et al., Review on the battery model and SOC estimation method, *Processes* 9 (9) (2021) 1685, <https://doi.org/10.3390/pr9091685>. Available: <https://search.proquest.com/docview/2576497406>.
- [20] M.S. Hossain Lipu, et al., Data-driven state of charge estimation of lithium-ion batteries: algorithms, implementation factors, limitations and future trends, *J. Clean. Prod.* 277 (2020) 124110, <https://doi.org/10.1016/j.jclepro.2020.124110>.
- [21] F. Yang, et al., State-of-charge estimation of lithium-ion batteries based on gated recurrent neural network, *Energy (Oxford)* 175 (2019) 66–75, <https://doi.org/10.1016/j.energy.2019.03.059>.
- [22] Q. Gong, P. Wang, Z. Cheng, A novel deep neural network model for estimating the state of charge of lithium-ion battery, *Journal of Energy Storage* 54 (2022) 105308, <https://doi.org/10.1016/j.est.2022.105308>.
- [23] X. Fan, et al., SOC estimation of Li-ion battery using convolutional neural network with U-Net architecture, *Energy (Oxford)* 256 (2022) 124612, <https://doi.org/10.1016/j.energy.2022.124612>.
- [24] B. Fu, et al., An improved neural network model for battery smarter state-of-charge estimation of energy-transportation system, *Green Energy and Intelligent Transportation* 2 (2) (2023) 100067, <https://doi.org/10.1016/j.geits.2023.100067>.
- [25] O. Demirci, et al., Review of battery state estimation methods for electric vehicles - Part I: SOC estimation, *Journal of Energy Storage* 87 (2024) 111435, <https://doi.org/10.1016/j.est.2024.111435>.
- [26] Z. Zhang, et al., *A State-of-charge Estimation Method Based on Bidirectional LSTM Networks for Lithium-ion Batteries*, Dec 13, 2020, pp. 211–216.
- [27] C. Hu, et al., State of charge estimation for lithium-ion batteries based on TCN-LSTM neural networks, *Jes* 169 (3) (2022) 30544, <https://doi.org/10.1149/1945-7111/ac5cf2>. Available: <https://iopscience.iop.org/article/10.1149/1945-7111/ac5cf2>.
- [28] J. Hao et al., "Short-term Power Load Forecasting for Larger Consumer Based on TensorFlow Deep Learning Framework and Clustering-regression Model," in Oct 2018, pp. 1–6.
- [29] J. Liu, et al., Prediction of remaining useful life of multi-stage aero-engine based on clustering and LSTM fusion, *Reliab. Eng. Syst. Saf.* 214 (2021) 107807, <https://doi.org/10.1016/j.res.2021.107807>.
- [30] H. Yu, et al., Corn leaf diseases diagnosis based on K-means clustering and deep learning, *Access* 9 (2021) 143824–143835, <https://doi.org/10.1109/ACCESS.2021.3120379>. Available: <https://ieeexplore.ieee.org/document/9576102>.
- [31] A.R. Khan, et al., Brain tumor segmentation using K-means clustering and deep learning with synthetic data augmentation for classification, *Microsc. Res. Tech.* 84 (7) (2021) 1389–1399, <https://doi.org/10.1002/jemt.23694>. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/jemt.23694>.
- [32] P. J. Kollmeyer et al. Tesla Model 3 2170 Li-ion Cell Dataset and Battery SOC Estimation Blind Modeling Tool. <https://doi.org/10.5683/SP3/ZVTR4B>.
- [33] R. N. Vieira et al., "Feedforward and NARX Neural Network Battery State of Charge Estimation With Robustness to Current Sensor Error," in Jun 21, 2023, pp. 1–6.
- [34] V. Pendyala, F. Nishanth, *Development of a Machine Learning Technique to Accurately Estimate Battery State of Charge*, Dec 9, 2022, pp. 1–6.
- [35] A.E. Ezugwu, et al., A comprehensive survey of clustering algorithms: state-of-the-art machine learning applications, taxonomy, challenges, and future research prospects, *Eng. Appl. Artif. Intell.* 110 (2022) 104743, <https://doi.org/10.1016/j.engappai.2022.104743>.
- [36] E. B. Osunwoke et al., "A Machine Learning-enabled Clustering Approach for Large-scale Classification of Solar Data," in Nov 14, 2021, pp. 1.
- [37] C. Sánchez-Rebollo, et al., *Detection of jihadism in social networks using big data techniques supported by graphs and fuzzy clustering*, *Complexity* 2019 (2019).
- [38] O. Ozdemir, A. Kaya, Effect of parameter selection on fuzzy clustering, *Mehmet Akif Ersoy Üniversitesi Uygulamalı Bilimler Dergisi* 2 (1) (2018) 22–33, <https://doi.org/10.31200/makuubd.348688>.
- [39] F. Kratzert, et al., Rainfall-runoff modelling using long short-term memory (Lstm) networks, *Hydrol. Earth Syst. Sci.* 22 (11) (2018) 6005–6022, <https://doi.org/10.5194/Hess-22-6005-2018>. Available: <https://search.proquest.com/docview/2136490956>.
- [40] B. Leibe, et al., *Spatio-temporal LSTM with trust gates for 3D human action recognition*, in: *Computer Vision - ECCV 2016* Anonymous, Springer International Publishing AG, Switzerland, 2016, pp. 816–833.
- [41] C. Han, X. Fu, Challenge and opportunity: deep learning-based stock price prediction by using bi-directional LSTM model, *Frontiers in Business, Economics and Management* 8 (2) (2023) 51–54, <https://doi.org/10.54097/fbem.v8i2.6616>. Available: <https://explore.openaire.eu/search/result?id=doi::ef993f05447450f19f358a6de8764d72>.
- [42] H. Jiang, et al., Construction and analysis of emotion computing model based on LSTM, *Complexity (New York, N.Y.)* 2021 (2021) 1–12, <https://doi.org/10.1155/2021/8897105>.
- [43] F. Yang, et al., State-of-charge estimation of lithium-ion batteries using LSTM and UKF, *Energy (Oxford)* 201 (2020) 117664, <https://doi.org/10.1016/j.energy.2020.117664>.
- [44] J. Chen, et al., SOC estimation for lithium-ion battery using the LSTM-RNN with extended input and constrained output, *Energy (Oxford)* 262 (2023) 125375, <https://doi.org/10.1016/j.energy.2022.125375>.
- [45] Z. Chen, et al., Synthetic state of charge estimation for lithium-ion batteries based on long short-term memory network modeling and adaptive H-infinity filter, *Energy (Oxford)* 228 (2021) 120630, <https://doi.org/10.1016/j.energy.2021.120630>.
- [46] S. Montaha, et al., Time distributed-CNN-LSTM: a hybrid approach combining CNN and LSTM to classify brain tumor on 3D MRI scans performing ablation study, *Access* 10 (2022) 60039–60059, <https://doi.org/10.1109/ACCESS.2022.3179577>. Available: <https://ieeexplore.ieee.org/document/9786658>.
- [47] F. Karim, S. Majumdar, H. Darabi, Insights into LSTM fully convolutional networks for time series classification, *Access* 7 (2019) 67718–67725, <https://doi.org/10.1109/ACCESS.2019.2916828>. Available: <https://ieeexplore.ieee.org/document/8713870>.
- [48] N.S. Ranawat, et al., Performance evaluation of LSTM and Bi-LSTM using non-convolutional features for blockage detection in centrifugal pump, *Eng. Appl. Artif. Intell.* 122 (2023) 106092, <https://doi.org/10.1016/j.engappai.2023.106092>.