

# Attention towards chemistry agnostic and explainable battery lifetime prediction

Fuzhan Rahmanian<sup>1,2,4,5,6\*</sup>, Robert M. Lee<sup>3</sup>, Dominik Linzner<sup>3</sup>, Kathrin Michel<sup>3</sup>, Leon Merker<sup>1,2</sup>, Balazs B. Berkes<sup>3</sup>, Leah Nuss<sup>1,4,5,6</sup>, and Helge Sören Stein<sup>1,4,5,6\*</sup>

<sup>1</sup>Helmholtz Institute Ulm, Applied Electrochemistry, Helmholtzstr. 11, 89081 Ulm, Germany

<sup>2</sup>Karlsruhe Institute of Technology, Institute of Physical Chemistry, Fritz-Haber-Weg 2, 76131 Karlsruhe, Germany

<sup>3</sup>BASF SE, Ludwigshafen, Germany

<sup>4</sup>Technische Universität München, School of Natural Sciences, Department of Chemistry Lichtenbergstr 4, 85748 Garching, Germany

<sup>5</sup>Technische Universität München, Munich Data Science Institute, Walther-von-Dyck-Straße 10, 4, 85748 Garching, Germany

<sup>6</sup>Technische Universität München, Munich Institute for Robotic and Machine Intelligence, Georg-Brauchle-Ring 60-62, 80992 Munich, Germany

\*corresponding author(s): Helge Sören Stein (helge.stein@tum.de), Fuzhan Rahmanian (fuzhan.rahmanian@tum.de)

## Abstract

Predicting and monitoring battery life early and across chemistries is a significant challenge due to the plethora of degradation paths, form factors, and electrochemical testing protocols. Existing models typically translate poorly across different electrode, electrolyte, and additive materials, mostly require a fixed number of cycles, and are limited to a single discharge protocol. Here, an attention-based recurrent algorithm for neural analysis (ARCANA) architecture is developed and trained on a unique, ultra-large, proprietary dataset from BASF and a large Li-ion dataset gathered from literature across the globe. ARCANA generalizes well across this diverse set of chemistries, electrolyte formulations, battery designs, and cycling protocols and thus allows for universal extraction of data-driven knowledge of the degradation mechanisms. The model's adaptability is further demonstrated through fine-tuning on Na-ion batteries. ARCANA advances the frontier of large-scale time series models in analytical chemistry beyond textual data and holds the potential to significantly accelerate discovery-oriented battery research endeavors.

## 1 Introduction

Lithium-ion batteries (LIBs) enable the electrification of everything, yet there is a maze of challenges that must be navigated in order to optimize the batteries of the future <sup>5, 61, 51, 15</sup>. Critical to the advancement of battery research is the rapid understanding of why and how some batteries degrade and what needs to be changed to prevent premature capacity fade <sup>54</sup>. Material degradation can occur due to numerous factors, including unpreventable solid electrolyte interphase growth, loss of active material, and other electrochemical phenomena <sup>31</sup>. However, investigating battery degradation is a time-consuming task, as non-linear capacity loss can occur over hundreds or thousands of cycles <sup>6</sup>. Another challenge in early lifetime prediction is the diversity of battery chemistries in the anode, cathode, and electrolyte, along with various form factors and testing protocols.

Battery lifetime can be evaluated through various methods, such as conventional cycling until end of life (EOL) under constant current-constant voltage (CC-CV) conditions or cycling for a predetermined

42 number of cycles. From these data, measures such as coulombic efficiency (CE) can be calculated <sup>62</sup>  
43 and correlated to more in-depth techniques such as electrochemical impedance spectroscopy (EIS) <sup>46</sup> to  
44 fundamentally assess the underlying degradation mechanisms. Accurate measurement of CE <sup>18, 52</sup> does,  
45 however, require bespoke instrumentation and a considerable amount of time, i.e. cycling a battery  
46 for 1000 cycles at  $1C/1D$  takes approximately 11 weeks. Reducing the required number of cycles  
47 by a factor of 10, while maintaining a high level of fidelity, is therefore of great interest <sup>7</sup>. Machine  
48 Learning (ML) and Deep Learning (DL) can accelerate testing by lowering the number of cycles  
49 required to understand the underlying chemistries <sup>3</sup>. An example of predicting EOL of batteries using  
50 initial discharge capacity curves was demonstrated by Severson et al. <sup>51</sup>, who used regression models.  
51 They integrated data generation with data-driven models to forecast the lifetime of LFP/graphite cells  
52 based on  $\Delta Q(V)$  and classified their longevity. In further work, Attia et al. <sup>7</sup> employed a Bayesian  
53 algorithm to accelerate the optimization of fast-charging protocols. By using early-cycle data for low-  
54 fidelity predictions, the approach enabled the optimization of high-fidelity experimental outcomes, thus  
55 significantly reducing the experimental duration from 500 to 16 days.

56 The most reliable models do not, however, merely predict just predict a quantity but also allow as-  
57 sessment of the model's uncertainty. Emblematic of this is the work by Tong et al. <sup>56</sup>, who introduced  
58 ADLSTM-MC, a hybrid predictive model using adaptive dropout long short term memory (LSTM)  
59 with Monte Carlo simulations. This approach, which requires minimal training data, enhances robust-  
60 ness through Bayesian-optimized dropout rates and improves the remaining useful life of two types  
61 of LIBs. In a correlative study <sup>47</sup>, a recurrent autoregressive deep ensemble network with aleatoric  
62 and epistemic uncertainties was developed along with saliency analysis to assess the impact of input  
63 parameters on output prediction. This provided an intuitive understanding of feature importance.  
64 Another advantage of using DL algorithms is their ability to use raw data, which has gained interest in  
65 the estimation of battery State of Health (SOH). For instance, Yang et al. <sup>63</sup> developed a novel hybrid  
66 convolutional neural network architecture with parallel residual connections, which utilizes raw data  
67 across multiple dimensions. By incorporating attention mechanisms, their model achieves remarkable  
68 accuracy in predicting the early stages of degradation. Although these approaches are applied in bat-  
69 tery research <sup>27, 68</sup>, their prominence is not as widespread as in other scientific fields. However, this  
70 lesser emphasis provides an opportunity for further exploration and discovery. Beyond these early life-  
71 time prediction models, sequence-to-sequence (Seq-to-Seq) models have been used to monitor battery  
72 lifetime and SOH <sup>27, 34, 20</sup>. They leverage intrinsic temporal dependencies in degradation data, provid-  
73 ing high predictive accuracy and computational efficiency. Li et al. <sup>34</sup> developed a one-shot LSTM-based  
74 Seq-to-Seq framework which not only predicts future capacities, but also identifies knee points in the  
75 degradation curve, maintaining stability even in the face of stochastic disturbances. Although Seq-to-  
76 Seq models demonstrate robust predictions, they also exhibit limitations in generalization and require  
77 large and diverse datasets to enhance performance <sup>15</sup>.

78 Despite the promises made by ML and DL for lifetime predictions <sup>11, 24, 55</sup>, these models, while  
79 robust, face challenges of precision and trustworthiness <sup>36</sup>. Existing models often focus on single-task  
80 learning, neglecting the potential benefits of multi-objective learning for various predictive settings <sup>15</sup>.  
81 In particular, data-driven approaches <sup>40</sup> tend to overlook the inherent variations between, for ex-  
82 ample, production batches or individual cells <sup>9</sup>. Such discrepancies, originating from manufacturing  
83 processes or aging mechanisms, can profoundly impact lifetime predictions. Addressing these varia-  
84 tions for accurate forecasting remains a central yet unresolved research question. Furthermore, despite  
85 the assertions of recent studies that they are chemistry-agnostic <sup>47, 48</sup>, they often require enhanced  
86 explainability to optimize their effectiveness in various chemistry settings. Transfer learning offers a  
87 promising solution to the challenge of scarce data, but requires more investigation for transparency  
88 and interpretability <sup>30</sup>. The acquisition of extensive datasets, essential for DL algorithms <sup>64</sup>, remains a  
89 significant hurdle <sup>22, 40, 59</sup>. Nevertheless, innovative strategies, such as the use of common features in  
90 databases and the documentation of various chemistries and protocols <sup>35</sup>, establish the foundation for  
91 more in-depth research <sup>64</sup>. Our goal is to develop a model that is not only universally applicable, but  
92 also robust, with the capability to provide both uncertainty quantification and explainability. Such a  
93 model would be invaluable to the academic community and would find marketable applications in the  
94 real world <sup>64</sup>, accelerating battery design and data collection based on active learning.

## 2 Results

### 2.1 Data resources

Developing a model that generalizes well necessitates a diverse and large dataset<sup>40</sup> that ideally covers a spectrum of chemistries and formats given high-dimensional correlations and cell variations<sup>66, 30</sup>, obtained from various laboratories and measured under different operating conditions<sup>7</sup>. Data diversity not only ensures an accurate representation of different cycling behaviors, but also tames the irreducible uncertainty in the predictions while mitigating the risk of overfitting. However, the scarcity of large and comprehensive datasets<sup>36</sup> that include both high and low-performing cells creates a challenge for training generalized models, i.e., to overcome a positive bias<sup>45, 30</sup>. Available data often exhibit noise, discontinuities, and varying formats that require extensive curation, adding a layer of complexity. Initiatives such as Battery Archive<sup>19</sup> or other cloud services<sup>32</sup> are therefore commendable in promoting Findable, Accessible, Interoperable, and Reusable (FAIR) data<sup>23, 58</sup> handling in battery research<sup>22, 59</sup>.

In this study, we develop a model trained on ca. 17400 batteries from BASF research laboratories that covers a diverse range of LIBs chemistries and multiple cycling protocols. Exposure of our model to such a wide variety of data enables robust generalization. Utilizing our pre-trained model on a unique set of unseen data, we effectively predict the early degradation trajectory. The ultimate test of our model, therefore, is to apply it to data from cells produced in a different location and with varying chemistries. Due to intellectual property constraints that prevent the authors from making the model trained on the BASF dataset openly accessible, we have retrained our model by leveraging a diverse array of publicly available datasets from respected institutions and research groups, including the Toyota Research Institute (TRI) in partnership with MIT and Stanford<sup>2, 1</sup>, NASA<sup>50</sup>, the Center for Advanced Life Cycle Engineering (CALCE)<sup>29</sup>, Karlsruhe Institute of Technology (KIT)<sup>69</sup>, Hawaii Natural Energy Institute (HNEI)<sup>21</sup>, and Sandia National Laboratories (SNL)<sup>21</sup>. Furthermore, we have incorporated data from our in-house cycled cells<sup>67, 38, 39, 42</sup> with successful and failed experiments, to further enrich model training and reduce bias. In the Supplementary Section 1 we provide an overview of all datasets; we include a brief summary in table 1 with an indication of which datasets were used during training and which remained completely unseen for model testing. This approach ensures a thorough understanding of the data sources, thus improving the transparency and reproducibility of our research.

### 2.2 Architecture Overview

Central to this study is the Attention-based ReCurrent Algorithm for Neural Analysis with LSTM (ARCANA) model. This is an attention-based Seq-to-Seq architecture specifically engineered to assess early stage battery degradation and perform lifecycle monitoring. The model demonstrates superior multitasking capabilities, supported by its high modularity and dynamic adaptability. It is designed to utilize a flexible range of past battery cycle data, known as historical temporal segments, for input. In addition, the model includes predetermined parameters for future conditions, such as discharge rates and cycle numbers. These parameters are known in advance of the experiment, i.e. they are controlled by the measurement device and are referred to as encoded temporal segments. This dual capability offers multifaceted advantages, from cost and time savings to improved material selection and protocol optimization.

The ARCANA model is augmented with additional features such as the attention mechanism, which provides insight into the decision-making process of the model. This feature distinguishes between predictions based on underlying patterns and those arising from stochastic variability. Saliency analysis is additionally performed to emphasize the relative importance of each parameter through a computation of the absolute gradient of the model output relative to the input of the test set. It quantifies the sensitivity of the input parameters, revealing how minor variations significantly alter the output results<sup>47</sup>, thus aligning the internal logic of the model with domain-specific knowledge. Adding another layer of robustness is uncertainty quantification, which is valuable not only for understanding the reliability of cycling protocols, but also for assessing material performance across different battery chemistries.

As illustrated in the Unified Modeling Language (UML) diagram (Fig. 1), the ARCANA model consists of four principal classes, each performing a different function, and is designed to accept raw data, thus negating the need for preliminary feature engineering. This design versatility extends to its operational modes with Naive Training for initial experiments, Dynamic Tuning for real-time

Table 1: An overview of the collected cycling data utilized for training and testing. The model  $M(B)$ , was trained with data provided by BASF, and the model  $M(P)$  was trained with publicly available data. The model  $M(P)_f$  represents a fine-tuned version of  $M(P)$  for lithium-ion coin cell data.  $M(P)_{Na}$  and  $M(B)_{Na}$  models are fine-tuned  $M(B)$  and  $M(P)$ , respectively, adapted for sodium coin cells.

Location	Cell form	Cell chemistry	Protocol Charge\Discharge	No. Cell	Cycle range	Nominal capacity [Ah]	Usage
BASF	Coin	heterogenous	multimodal	17400	multimodal	multimodal	$M(B)$ Train\Val
TRI <sup>2</sup>	Cylindrical commercial	LFP\graphite	$CC1(Q1)CC2$ , $CC - CV@1C, 4.2V$ \CC@4C	124	169 – 2235	1.1	$M(P)$ Train\Val
TRI <sup>1</sup>	Cylindrical commercial	LFP\graphite	$CC1(20\%)CC2(40\%)$ $CC3(60\%)CC4(80\%)$ , $CC - CV@1C, 4.2V$ \CC - CV@4C, 2V	233	100 – 862	1.1	$M(P)$ Train\Val
CALCE <sup>29</sup>	Prismatic commercial CX2	LCO\graphite	$CC - CV@0.5C, 4.2V$ , 6 \CC@(0.5C, 1C)	6	781 – 1082	1.35	Testing
CALCE <sup>29</sup>	Prismatic commercial CS2	LCO\graphite	$CC - CV@0.5C, 4.2V$ , 6 \CC@0.5C	6	1701 – 2016	1.1	$M(P)$ Train\Val
KIT <sup>69</sup>	Cylindrical commercial	NCA\graphite-Si	$CC - CV$ @(0.25C, 0.5C, 1C), 4.2V, \CC@1C	58	29 – 800	3.5	$M(P)$ Train\Val
KIT <sup>69</sup>	Cylindrical commercial	NCM\graphite-Si	$CC - CV$ @(0.25C, 0.5C, 1C), 4.2V, \CC@1C	55	43 – 1277	3.5	$M(P)$ Train\Val
KIT <sup>69</sup>	Cylindrical commercial	NCM+NCA\graphite	$CC - CV@0.5C, 4.2V$ , 9 \CC@(1C, 2C, 4C)	9	912 – 1031	2.5	Testing
KIT <sup>67</sup>	Coin self-made	LNO\graphite	$CC - CV@1C, 4.2V$ , 43 \CC@1C	43	82 – 505	0.004618	60% for $M(P)_f$ , 40% Testing
KIT <sup>38</sup>	Coin commercial	LCO\graphite	$CC - CV@1C, 4.25V$ , 26 \CC - CV@1C, 2.75V	26	150 – 600	0.045	$M(P)$ Train\Val
KIT <sup>39</sup>	Coin self-made	NMC622\graphite	$CC - CV@1C, 4.2V$ , 11 \CC@1C	11	228 – 501	0.00328	Testing
KIT <sup>42</sup>	Coin self-made	$Na_{0.9}[\dots]O_2$ \graphite	$CC@1C$ \CC@1C or C-rates test	44	40 – 140	0.00015	60% for $M(P)_{Na}$ and $M(B)_{Na}$ , 40% Testing
NASA <sup>50</sup>	Cylindrical commercial	NCA\graphite	$CC - CV$ @0.75C, 4.2V, \CC@(0.5C, 1C, 2C)	34	24 – 196	2.0	$M(P)$ Train\Val
HNEI <sup>21</sup>	Cylindrical commercial	LCO-NMC\graphite	$CC - CV@0.5C, 4.3V$ , 14 \CC@1.5C	14	1102 – 1133	2.8	$M(P)$ Train\Val
SNL <sup>21</sup>	Cylindrical commercial	LFP\graphite	$CC - CV$ @0.5C, 4.2V, \CC@(0.5C, 1C, 2C, 3C)	28	2621 – 19174	1.1	$M(P)$ Train\Val
SNL <sup>21</sup>	Cylindrical commercial	NCA\graphite	$CC - CV$ @0.5C, 4.2V, \CC@(0.5C, 1C, 2C)	24	463 – 7877	3.2	$M(P)$ Train\Val
SNL <sup>21</sup>	Cylindrical commercial	NMC\graphite	$CC - CV$ @0.5C, 4.2V, \CC@(0.5C, 1C, 2C, 3C)	25	388 – 11149	3.0	$M(P)$ Train\Val

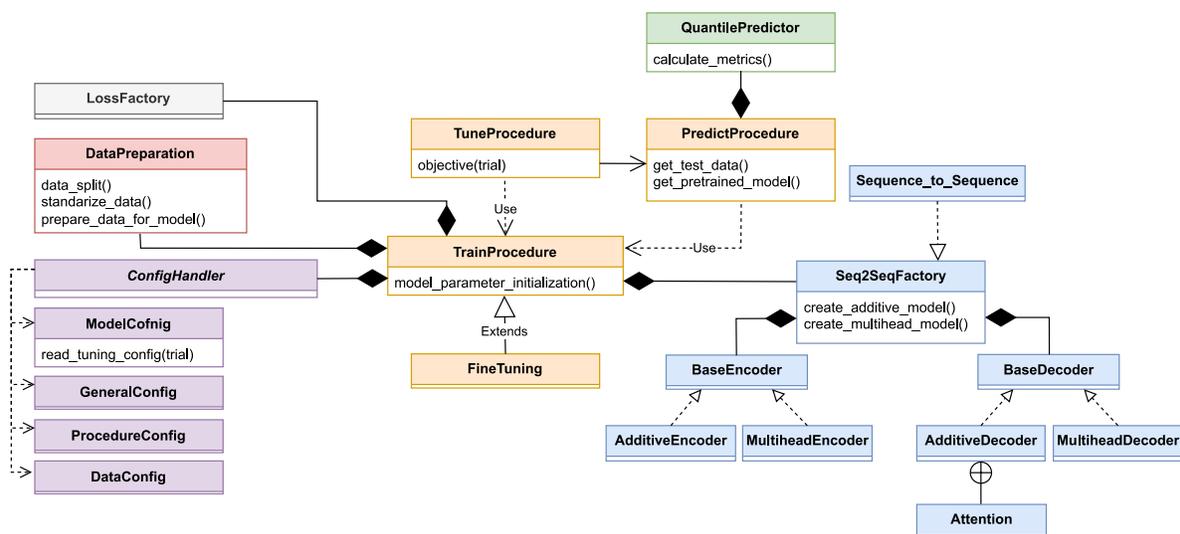


Figure 1: An UML diagram of the computational framework, designed around three principal class clusters. The first includes a `ConfigHandler` engineered to manage a comprehensive set of user-defined configurations and establishes a blueprint for handling various subconfigurations such as general settings, data properties, and model specifications. During hyperparameter optimization tasks, `ConfigHandler` interfaces with the Optuna optimization library to adaptively create and update the tuning configuration. The second key class structure includes `TrainProcedure` which serves as an architectural template for the training process. Its attributes are employed throughout the computational pipeline, starting with data preparation and extending to instantiation of specialized loss functions and Seq2Seq models via the `LossFactory` and `Seq2SeqFactory`. `FineTuning`, is a specialized subclass that inherits from `TrainProcedure` while `TuneProcedure` and `PredictProcedure`, the latter of which uses the `QuantilePredictor`, are incorporated into the pipeline depending on the desired use case and settings. The tuning operates on single trials with a `TPESampler` when multiple runs are desired. Lastly, `Seq2SeqFactory` is engineered to govern the instantiation of encoder-decoder architectures. Depending on the user-defined configurations, it can orchestrate a multihead or an additive encoder-decoder mechanism. The inclusion of custom attention mechanisms within the architecture is handled by the `AdditiveDecoder` class or the `MultiheadDecoder`, conditional upon the configuration stipulations.

150 adaptability via extensive hyperparameter optimization, Fine-Tuning for integration of a pretrained  
 151 model with selective gradient updating, and Prediction for efficient inference. Through modularity, a  
 152 logging mechanism ensures data integrity and traceability, adhering to FAIR data principles<sup>58</sup>. The  
 153 open source codebase uses the PyTorch library<sup>43</sup> for model development and the Optuna library<sup>4</sup> for  
 154 hyperparameter optimization.

## 155 The Encoder-Decoder Framework

156 The encoder (Fig.2a) initiates the Seq-to-Seq model in the ARCANA framework by processing  
 157 historical temporal segments of the past battery life cycles. Employing a LSTM network, it is designed  
 158 to capture complex, non-linear relationships and time dependencies inherent in sequence data. The  
 159 encoder processes the input tensor to accommodate sequences of different lengths, employing a padding  
 160 mechanism that enables the LSTM to efficiently process these sequences without being constrained  
 161 by their varying lengths. Within the LSTM, the temporal data is transformed into a new tensor,  
 162 constructing hidden and cell states that capture sequential information. A skip connection incorporates  
 163 the initial input into the LSTM output, thus preserving crucial temporal features and stabilizing  
 164 the learning process. Layer normalization, when applied to the LSTM output, not only accelerates  
 165 convergence but also leads to robust performance, mitigating the challenges associated with long-  
 166 sequence dependencies<sup>17</sup>. The encoder returns a rich latent representation of the historical data,  
 167 consisting of the output tensor and the updated hidden and cell states, which are then utilized by the  
 168 decoder to enable accurate forecasting in subsequent steps.

169 The decoder (Fig. 2a) takes on the task of generating future state predictions. It is initialized with  
 170 the hidden and cell states from the encoder and begins by processing the most recent historical cycle  
 171 data. The model then integrates its own previous predictions and known future conditions, such as  
 172 the expected discharge current and the cycle number. These two inputs are temporally encoded to  
 173 capture their positional relevance<sup>65</sup>, ensuring that the decoder is informed of the predefined condition  
 174 and the timing of each data point within the life cycle. The decoder employs an attention mechanism  
 175 that can dynamically adjust sequence weights, identifying critical information at each prediction step.  
 176 This approach overcomes the limitations of static-length vector representation in conventional encoder-

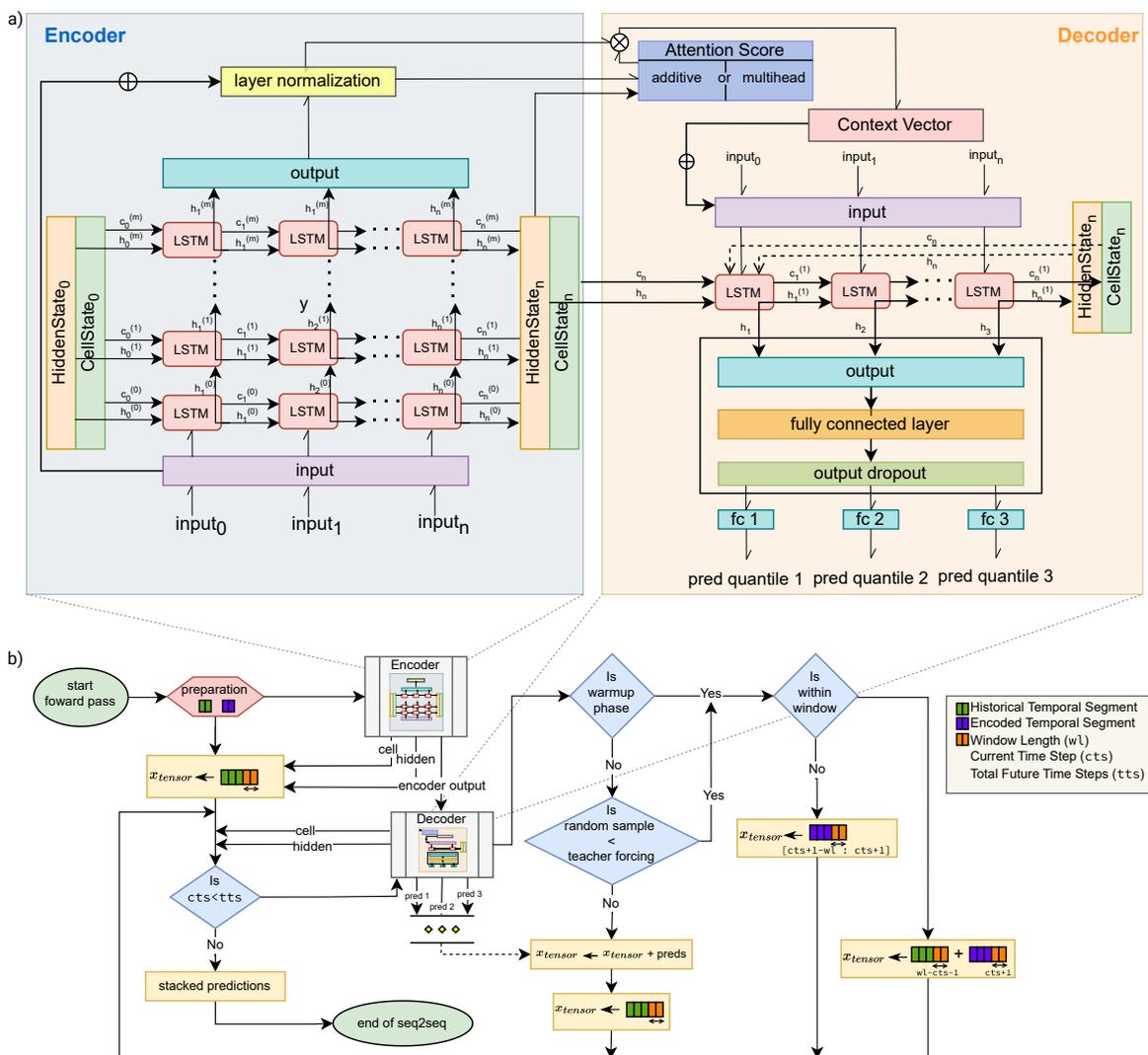


Figure 2: Architectural overview of Seq-to-Seq model. a) Presents the detailed architecture of the encoder and decoder components. The LSTM-based encoder processes historical temporal segments to capture the intricate pattern of battery life cycles. It integrates a skip-connection and layer normalization to preserve and stabilize essential key temporal features. The decoder is initialized with the encoder's final states and applies an attention mechanism to focus on relevant temporal features from the encoder output and enrich the context of its predictions. The attention-enhanced representations are combined with the initial decoder input and subsequently propagated through LSTM layers. A fully connected layer with leaky ReLU activation and a dropout layer - used solely during training and inactive during inference - for regularization, follow the LSTM outputs. The model outputs are then fed into three separate fully connected layers for predicting a specific quantile of the future distribution based on the pattern learned during training, thus providing a probabilistic characterization of the forecast. b) Illustrates the integrated Seq-to-Seq model flow, depicting the progression from encoding historical data to multi-output future forecasts. It highlights the sliding-window approach that underpins the model's capability to handle both the tail-end of historical data and the integration of self-generated forecasts with known future conditions. This process also captures the dynamic training process, which incorporates teacher forcing to enhance the predictive fidelity of the model.

177 decoder models<sup>8</sup>, allowing the decoder to focus on the most relevant parts of historical data. The  
178 attention mechanism then computes a context vector associated with the encoder’s output, which  
179 highlights the encoder sequences with the highest relevance to the current decoding task. This context  
180 vector, combined with the current input, forms a feature-rich tensor that is subsequently processed by a  
181 LSTM layer. Post-LSTM, the output layer is passed through a fully connected layer with a leaky ReLU  
182 activation function, crucial in maintaining network stability, and enhanced with a dropout layer placed  
183 to reduce overfitting risks. The culmination of this process is a decoder that generates forecasts for the  
184 0.1, 0.5, and 0.9 quantiles. These provide a probabilistic range indicative of the inherent uncertainty  
185 and offer a statistical interpretation of the potential future states of the degradation profile.

## 186 Seq-to-Seq Integration

187 In the broader Seq-to-Seq model, the encoder and decoder are orchestrated to facilitate the overall  
188 predictions, as can be seen in Fig. 2b. Here, the model processes the temporal data using a sliding  
189 window approach that enhances the ability to discern local patterns within long input sequences<sup>65</sup>.  
190 This technique allows for the integration of the last observed data or transition to the decoder’s self-  
191 generated predictions, supplemented with temporally encoded future conditions. During training, a  
192 dynamic teacher forcing strategy is employed, in which actual target outputs are used as inputs in  
193 lieu of previous predictions to promote model convergence, prediction fidelity, and generalizability  
194 in the model. This hybrid training strategy allows effective learning from the ground truth while  
195 gradually becoming equipped for self-guided predictions. At the end of the processing of this sequence,  
196 quantile-based predictions are collected into a stack of tensors, encapsulating a comprehensive forecast  
197 for subsequent decision-making processes. Thus, this forward pass provides fine-grained, probabilistic  
198 understanding of the evolving battery life-cycle stages, with the potential to inform risk assessment  
199 and optimize operational efficiency.

## 200 2.3 Experimental configuration

201 This study evaluates the ARCANA architectural model through a two-stage experimental process. Our  
202 aim is to present findings that resonate across multiple disciplines, highlighting both the complexity  
203 and versatility of our approach. The first stage involved training the model  $M$  with the coin cell  
204 dataset  $B$  from BASF. The resulting trained model is here denoted  $M(B)$ . We encoded predetermined  
205 parameters, including cycle number and discharge current, into temporal segments to capture past  
206 and future discharge conditions. The training used an additive attention mechanism in the ARCANA  
207 architecture for initial learning, with a detailed explanation in Sec. 4.1. In the second stage, the model  
208  $M$  is re-trained from scratch, with publicly available datasets as mentioned in Table. 1, and denoted as  
209  $M(P)$ . This entails various cell types, including 26 coin cells and 6 prismatic cells with Lithium-Cobalt-  
210 Oxide (LCO) cathodes, with the majority being cylindrical cells with Lithium-Iron-Phosphate (LFP),  
211 Nickel-Manganese-Cobalt (NMC), and Nickel-Cobalt-Aluminum Oxide (NCA) cathode materials. To  
212 address these cell chemistry variations, we introduced an additional predefined parameter, the nominal  
213 capacity of each cell in logarithmic format. This inclusion was critical for the model to effectively  
214 differentiate and interpret response characteristics<sup>53</sup>. The public dataset selected for  $M(P)$ , was  
215 significantly smaller, comprising 627 cell entries and accounting for only 3.35% of the total data size  
216 of the initial model  $M(B)$ . The dataset was distributed with 65% for training, 30% for validation, and  
217 5% for testing.

218 To emphasize generalizability and test model performance, we incorporated four distinct test  
219 datasets, each sourced from different locations and created by various experts. The first two test  
220 sets, denoted ( $D_{LNO}$ ) and  $D_{NMC}$ , comprise coin cell measurements made at the Institute of Physical  
221 Chemistry (IPC) of KIT, featuring the Lithium-Nickel-Oxide (LNO) and NMC materials, respectively.  
222 The third dataset consisted of cylindrical cells from Institute of Applied Materials (IAM) of KIT,  
223 containing NMC blended with NCA cathode materials ( $D_{NMC+NCA}$ ). The final dataset involved  
224 prismatic cells of the CALCE institute, with LCO materials ( $D_{LCO}$ ). The complete description of  
225 these cells is provided in the Supplementary Section 1. This approach in dataset selection and testing  
226 allowed an in-depth evaluation of  $M(P)$  for its adaptability to various cell types and experimental  
227 setups.

228 The publicly available data for  $M(P)$  presented unique challenges as they included prematurely  
229 failed cells and high experimental noise, in contrast to the high-quality data used for training  $M(B)$ .

230 These complexities required a change from an additive to a multihead attention mechanism in  $M(P)$ .  
231 We also encountered a wide range of cycles, from as few as 196 to as many as 19176. However, most of  
232 the tests we considered had fewer than 500 cycles. This variability posed a potential risk of gradient  
233 instability and inconsistent learning in the training process. To mitigate the risk of poor convergence  
234 and the possibility of overfitting, we adopted a standardization approach in which all cells were limited  
235 to a maximum of 500 cycles, ensuring better balance in the training data and reducing bias, thus  
236 increasing reliability.

237 Both  $M(B)$  and  $M(P)$  focused on predicting three parameters, which were selected for their estab-  
238 lished significance in the existing literature and their availability across the datasets. They included  
239 discharge capacity, crucial for understanding the SOH<sup>51</sup>, CE, as emphasized in studies by Burns et  
240 al.<sup>13, 14</sup> as the key to understanding the impact of electrode additives and electrode materials on  
241 battery long-term performance, and the voltage drop during the relaxation phase between charging  
242 and discharging cycles. The last parameter is less explored but, as described by e.g. Zhu et al.<sup>70</sup>, it  
243 offers valuable insights independent of the charging process. This parameter is easily calculated from  
244 cycling data, even if the studies where the data originated did not directly measure it. In this section,  
245 we evaluate our model’s performance on various scenarios, focusing on the impact of data quality on  
246 model generalization and interpretability, investigating its adaptability to different chemistries, and  
247 deriving insights from attention mechanisms and saliency analysis.

## 248 2.4 Model performance across battery types

249 The hyperparameters of  $M(P)$  were selected using Optuna’s hyperparameter tuning with 250 trials  
250 and are described in the Supp. 3. The model generalization is evaluated on two datasets; cylindrical  
251 cells of  $D_{NMC+NCA}$  and prismatic cells of  $D_{LCO}$ , neither of which were seen by the model during  
252 training. Here, the objective was to determine how effectively the model generalizes across different  
253 battery configurations despite the presence of noisy data.

254 As shown in Fig. 3, the model handles multidimensional predictions for both  $D_{NMC+NCA}$  and  
255  $D_{LCO}$  well. For  $D_{NMC+NCA}$ , it accurately forecasts up to 500 cycles based on 24 input cycles (see  
256 Panel I, Fig. 3) even though the extracted data exhibits occasional jumps, despite the discharge  
257 current remaining constant throughout. Given that these unexpected jumps are not annotated in the  
258 original dataset, we have chosen to acknowledge their presence, but not alter them for the sake of  
259 data integrity. Aggregated attention weights in early cycles indicate their importance for long-term  
260 forecasting. Emblematic is  $D_{LCO}$ , that starts from a 23 cycle profile (Panel II, Fig. 3); the model  
261 demonstrates robustness even in the presence of more complex noise patterns. Here, the attention  
262 weights are distributed not only in the initial cycles but also in later cycles, proving the necessity of  
263 incorporating an attention mechanism.

264 Illustrating the model’s generalization capabilities, a detailed analysis of  $Q_{dis}$  in Fig. 4 is presented.  
265 In both  $D_{NMC+NCA}$  and  $D_{LCO}$ , there is good agreement between the model’s predictions and actual  
266 values (Panel I & II, Fig. 4a), as complemented by the density graphs in Fig. 4b. For  $D_{NMC+NCA}$ , the  
267 predicted and actual densities closely overlap. For  $D_{LCO}$ , the predicted density is highly similar, with  
268 a slightly skewed distribution towards lower  $Q_{dis}$ . The better density distributions for  $D_{NMC+NCA}$   
269 are likely attributable to the larger proportion of cylindrical cells in the training data, which accounts  
270 for 94.9% of the total.

271 A detailed evaluation of the uncertainty of the model  $M(P)$  is provided in Fig. 4c-e for both  
272 datasets. Panel I & II of Fig. 4c evaluate the calibration by comparing the observed quantile proportions  
273 to the expected proportions under the assumption of a normal distribution. This continuous curve  
274 indicates the model’s general performance across the entire probability distribution. The miscalibration  
275 area, quantified by the degree of deviation from the ideal diagonal line, represents the aggregate of  
276 discrepancies<sup>28</sup>. For  $D_{NMC+NCA}$ , the predicted distribution of  $Q_{dis}$  is well calibrated around the  
277 median but diverges at the tail, with calibration points showing underconfidence at higher quantiles.  
278 For  $D_{LCO}$  the individual calibration points suggest a slight overconfidence in the 10th-50th percentile  
279 and underconfidence in the ranges 50th-90th and 10th-90th percentile. The miscalibration areas for  
280  $D_{LCO}$  is 0.16, which is slightly higher than  $D_{NMC+NCA}$ , likely due to noisier data. The overall  
281 calibration performance across both datasets is comparable. Fig. 4 e) shows a histogram of prediction  
282 interval quantiles, revealing the spread between the 10th and 90th percentiles and evaluating the  
283 concentration of its predictive distribution as indicated by sharpness. The lower values suggest higher  
284 confidence in the prediction<sup>25</sup>. For  $D_{NMC+NCA}$ , a bimodal distribution highlights variable prediction

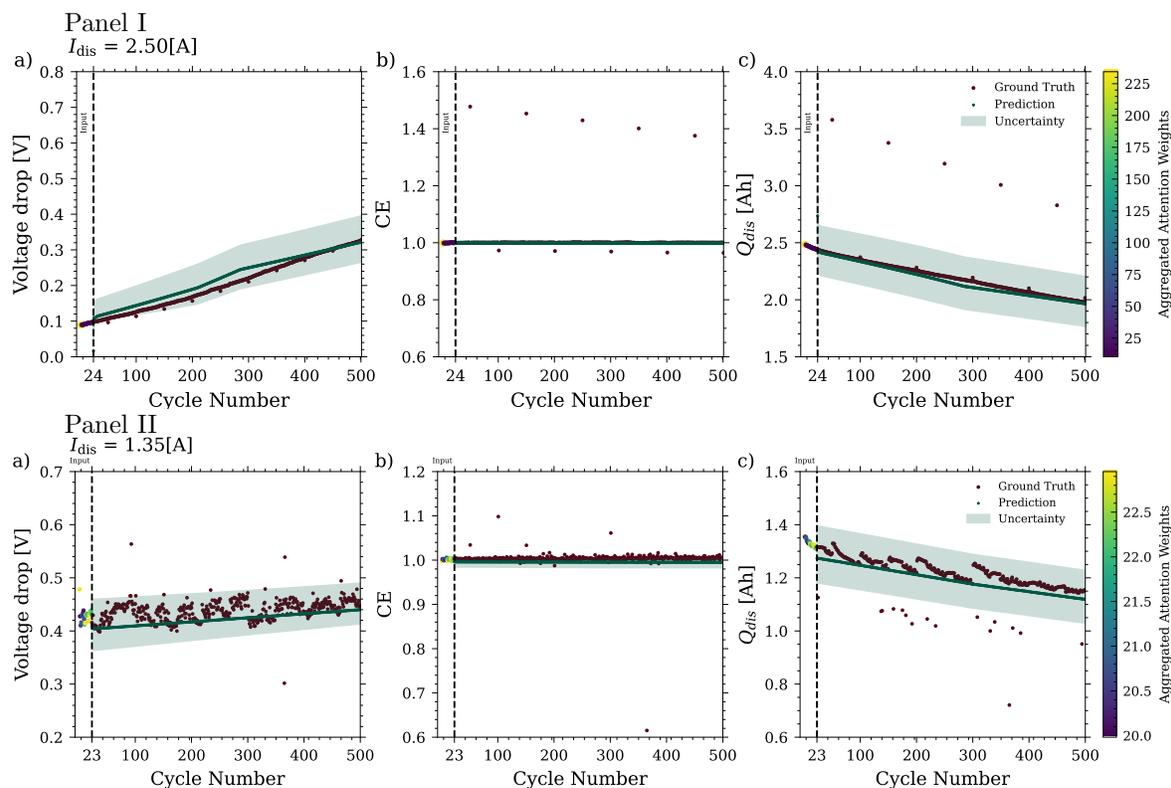


Figure 3: ARCANAs predictive performance on two datasets, namely cylindrical  $D_{NMC+NCA}$  in Panel I and prismatic  $D_{LCO}$  in Panel II, when predicting battery behaviour over 500 cycles for three predictors of Voltage drop [V], CE and  $Q_{dis}$  [Ah]. The model emphasizes its robust noise filtering and adaptive attention mechanism across different data characteristics.

285 certainty across cycles, suggesting potential fluctuations in battery behavior.  $D_{LCO}$  shows two clusters  
 286 of distributions, mostly around a central quantile with a sharpness of 0.19, indicative of consistent  
 287 uncertainty. Fig. 4d further supports these findings by illustrating the model’s median prediction  
 288 uncertainty and the variability of these predictions by interquantile range (IQR). Here,  $D_{NMC+NCA}$   
 289 in Panel I shows varying IQR, suggesting changes in model confidence over lifespan. In contrast,  
 290  $D_{LCO}$  maintains a more uniform IQR, indicating steady prediction uncertainty and aligning with the  
 291 model’s attention on later cycles to contend with the increased complexity and noise. These metrics  
 292 complement the information provided in Fig. 4c-e, serving as a benchmark for the model’s reliability  
 293 and its capacity to generalize within a precise estimate range.

294 The multitasking capabilities of  $M(P)$  are further highlighted by its performance in predicting  
 295 the second parameter, voltage drop (SI). The model exhibits strong prediction accuracy with both  
 296 datasets.  $D_{NMC+NCA}$  shows a smaller range of predictions over increasing cycles, and  $D_{LCO}$  shows  
 297 a stable range with decreasing median intervals, while the calibration accuracy and the reliability of  
 298 the predictions remain high across both datasets. The performance on the third predictor, CE (SI),  
 299 shows consistency and low prediction uncertainty, although the high measurement noise present in this  
 300 dimension poses a challenge and makes convergence more demanding<sup>26</sup>. The evaluation metrics for  
 301  $M(P)$  (Supp. Table. 1) demonstrate its predictive strengths for both  $D_{NMC+NCA}$  and  $D_{LCO}$ . For  
 302 the  $D_{LCO}$  dataset, the voltage drop is predicted with root mean square error (RMSE) of 0.0335 and  
 303 mean absolute percentage error (MAPE) of 6.6052. However,  $D_{NMC+NCA}$  outperforms in CE with  
 304 significantly lower error rates of 0.0256 and 0.2489 for the RMSE and MAPE, respectively. However,  
 305 both datasets present higher error rates in the predicted discharge capacity. To counteract the impact  
 306 of systematic noise, Median Absolute Error (medAE) is used along with MAE for a more robust error  
 307 analysis. These metrics highlight  $M(P)$ ’s versatile predictive capabilities, in handling diverse dataset  
 308 requirements for multiple features and long-term predictions<sup>15, 33</sup>.

309 We further examine  $M(P)$ ’s performance on unseen coin cell datasets,  $D_{LNO}$  and  $D_{NMC}$ . The  
 310 model predicts the voltage drop and CE well, but shows limitations and high uncertainty when pre-  
 311 dicting the discharge capacity with an RMSE of 0.5827. This may stem from the low representation  
 312 of coin cells in the training data; just 4.1% of the total. To alleviate this problem, we fine-tuned the

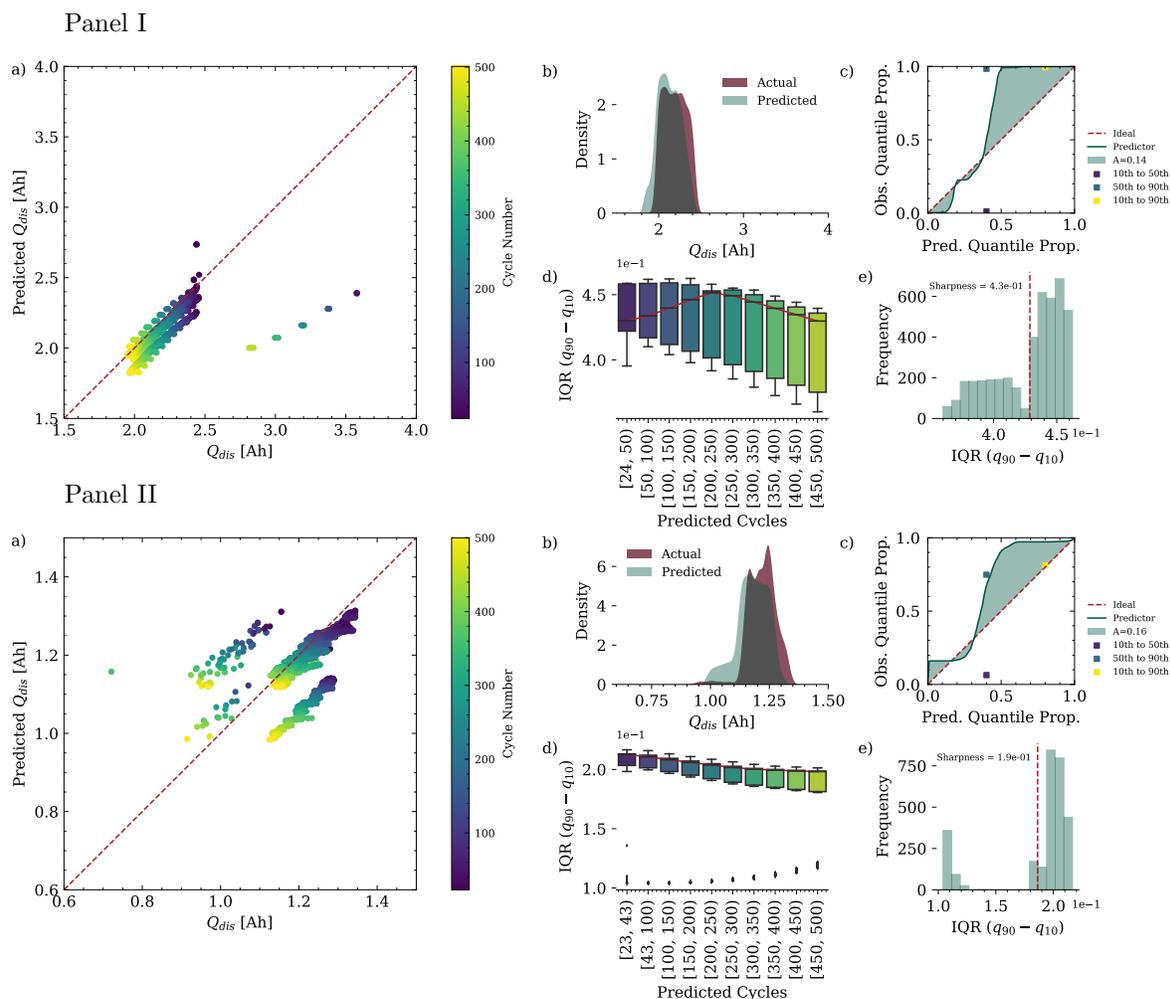
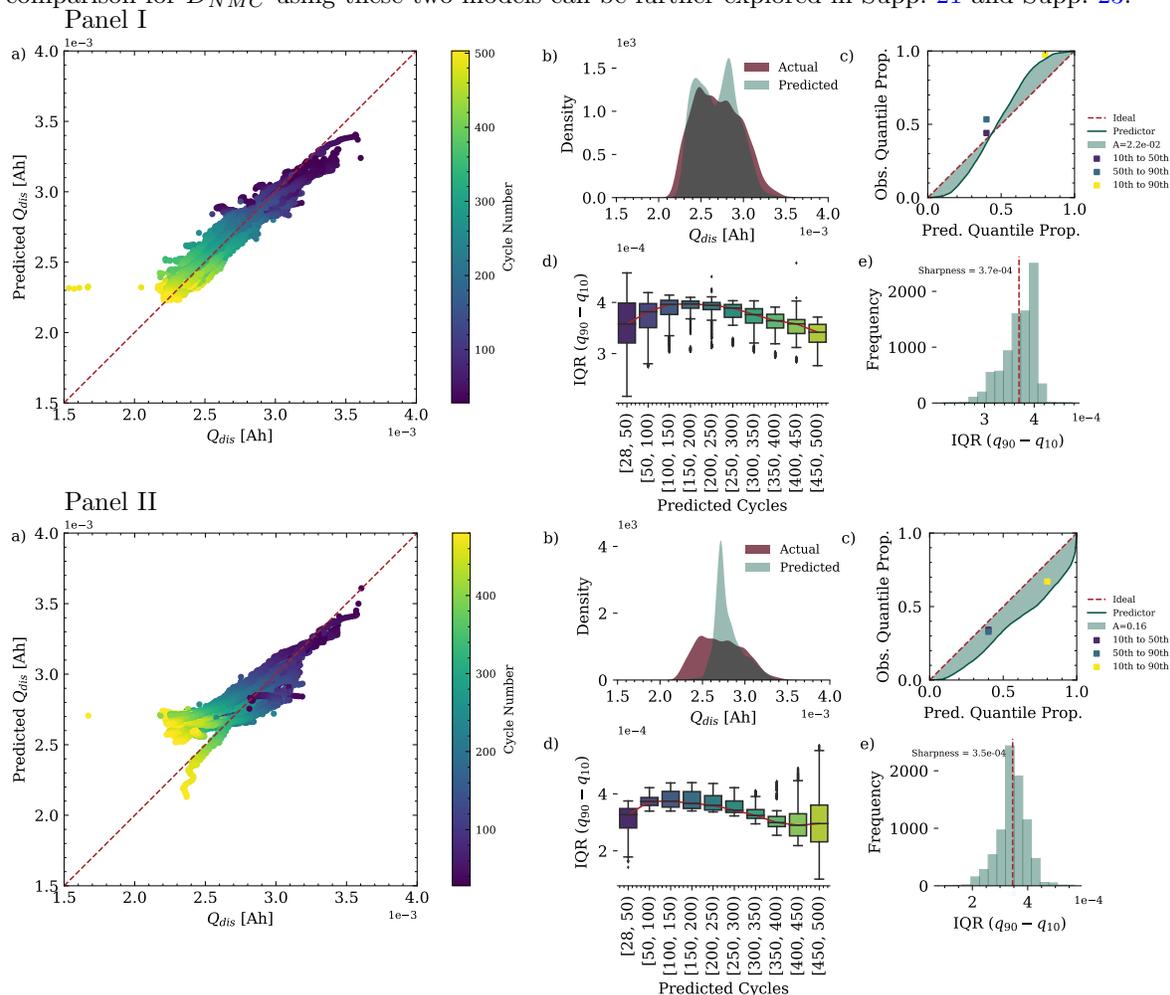


Figure 4: Comparative analysis of model predictions and its uncertainty and calibration for  $Q_{dis}$  for two datasets;  $D_{NMC+NCA}$  (Panel I) and  $D_{LCO}$  (Panel II), where a) depicts the relationship between predicted and actual values of  $Q_{dis}$ , with the diagonal dashed line indicating perfect prediction accuracy, b) illustrates the density distributions of predicted versus actual  $Q_{dis}$ . The calibration plot in c) assumes a normal distribution, where the mean and standard deviation are estimated from the 10th, 50th, and 90th percentiles of predictions. It depicts the cumulative proportion of actual  $Q_{dis}$  values that fall at or below the predicted quantile values, rather than within symmetric intervals around the predictions. The ideal diagonal line represents perfect calibration with the shaded area indicating the degree of miscalibration, denoted  $A$ . The approximately diagonal trend of the calibration line up to the 0.5 quantile shows that data with residuals below the median are well described by the predictive distribution. The jump from 0.5 to 1 indicates that the predictive distribution extends further to positive values than the observed distribution of residuals; almost all test data are already covered by the predicted 0.6 quantile for both datasets. However, the overall miscalibration areas for both datasets are quite similar, indicating that despite different patterns of over- and underconfidence at specific quantiles, the general calibration performance across both datasets is comparable. Box plots at d) show the prediction intervals over multiple cycles, demonstrating the median and variability of the model prediction uncertainty over the battery's lifespan. e) provides histograms that depict the quantile-based prediction interval width between the 10th and 90th percentiles, as a measure of sharpness. The red dashed line indicates the sharpness as the mean interval width and shows the concentration of the predictive distributions that indicate narrower distribution and consequently higher confidence in predicting  $Q_{dis}$  for  $D_{NMC+NCA}$  in Panel I.

313 decoder weights of  $M(P)$  using the data of 17 coin cells from  $D_{LNO}$ , resulting in an updated model,  
 314  $M(P)_f$ . This fine-tuning process is detailed in Supp. 14, and led to a substantial improvement in  
 315 predicting  $Q_{dis}$ , dropping the RMSE to 0.0002, indicating a significantly enhanced precision.  $M(P)_f$ 's  
 316 performance will be compared with  $M(B)$ , trained with the BASF dataset  $B$ , in the following section.

## 317 2.5 Model Performance on Coin Cell Data for Generalization Insights

318 While comparing the predictive performance of models  $M(B)$  and  $M(P)_f$  on subsets of unseen  $D_{LNO}$   
 319 and  $D_{NMC}$  dataset (Supp. 15),  $M(P)_f$  demonstrates reliable predictive alignment for voltage drop,  
 320 CE, and  $Q_{dis}$ . In contrast,  $M(B)$ , shows a divergent pattern in voltage drop predictions, which may be  
 321 due to its training on data with inherently long relaxation time profiles compared to those in  $D_{LNO}$ ,  
 322 where measurements are taken shortly after state changes. However, it maintains consistency in CE  
 323 predictions and adjusts  $Q_{dis}$  predictions in response to changes in the test protocol. The performance



**Figure 5:** Performance analysis of  $M(P)_f$  (Panel I) and  $M(B)$  (Panel II) on  $D_{LNO}$  for  $Q_{dis}$  prediction. Plots a) illustrate the relationship between models' predictions and the actual  $Q_{dis}$  with the diagonal line representing perfect prediction accuracy, plots b) compare the density distribution of actual and predicted  $Q_{dis}$ , plots c) present calibration curves that reflect the degree of alignment between predicted probabilities and observed frequencies under a normal distribution assumption. The discrete points on the calibration curve show the observed proportions of actual values that fall within three specific intervals based on the quantiles: between the 10th and 50th, 50th and 90th, and 10th and 90th percentiles. Model  $M(P)_f$  shows a high level of calibration for predicting  $Q_{dis}$  of  $D_{LNO}$  samples with a minimal miscalibrated area of 0.022. The points for 10th and 50th and 50th and 90th percentiles lie close to the diagonal line, indicating a nearly perfect calibration for these intervals.  $M(B)$  exhibits a slight overconfidence by deviating from the ideal line, with a miscalibration area of 0.16. The three calibration markers for  $M(B)$ , are all positioned just below the diagonal line, showing uniform overconfidence across these quantile ranges, yet they remain close to this line, indicating a generally well-calibrated model. Plots d) show the prediction intervals across lifespan cycles, highlighting models' uncertainty over time and plots e) detail the distribution of prediction intervals' quantiles between the 10th and 90th percentiles, which convey the models' prediction uncertainty; a distribution skewed towards the lower quantiles suggests a higher confidence in predictions at these quantiles. The sharpness, as a measure of mean interval width, is approximately similar for both models at  $3.7 \times 10^{-4}$  and  $3.5 \times 10^{-4}$  for  $M(P)_f$  and  $M(B)$ , respectively. Together, these plots demonstrate the  $M(P)_f$ 's precision in capturing discharge capacity behavior and  $M(B)$ 's robust generalization.

Table 2: Summary of evaluation metrics for  $D_{LNO}$

Metrics	$M(P)_f$			$M(B)$		
	Voltage drop	CE	$Q_{dis}$	Voltage drop	CE	$Q_{dis}$
<b>RMSE</b>	0.0703	0.0331	0.0002	0.1247	0.0588	0.0003
<b>MAPE</b>	9.2285	1.1922	20.7946	34.8638	4.4560	8.8914
<b>MAE</b>	0.0353	0.0076	0.0001	0.0867	0.0335	0.0002
<b>medAE</b>	0.0181	0.0021	0.0001	0.0513	0.0104	0.0001

325 In our analysis of  $D_{LNO}$  for  $Q_{dis}$ , Fig. 5 demonstrates that  $M(P)_f$  achieves high predictive fi-  
326 delity. This is evident from the dense alignment of the predictions with the actual values in the scatter

327 plot (Fig. 5a), and the significant overlap in distributions seen in the density plot (Fig. 5b). The  
 328 model's precision is further highlighted by concentrated prediction intervals and a calibration curve  
 329 that closely traces the diagonal (Fig. 5c-e). It achieves a high proportion of data points within the pre-  
 330 dictive bounds, indicative of accuracy, without excessively wide intervals that could decrease the utility  
 331 of the predictions. Panel II for  $M(B)$  also demonstrates a close tracking of the actual values, with  
 332 a marginally broader prediction interval and higher miscalibrated area of 0.16 compared to  $M(P)_f$ 's  
 333 of 0.022 (Panel I). Despite this variance,  $M(B)$  maintains a reasonable estimate range. Qualitatively  
 334 (Table 2),  $M(P)_f$  achieves better accuracy in predicting  $Q_{dis}$  with a lower RMSE and MAE.  $M(B)$   
 335 shows higher RMSE, especially in voltage drop, and lower MAPE in  $Q_{dis}$  (8.8914) suggesting effective  
 336 capture of proportional changes in the data, despite a larger absolute error. Detailed analyses of addi-  
 337 tional predictive dimensions for  $D_{LNO}$  for both models and the complete dataset  $D_{NMC}$  are available  
 338 in the supplementary materials. Despite the  $D_{LNO}$  data originating from another institute, the gen-  
 339 eralization of  $M(B)$  highlights the potential of well-trained DL models to overcome the variability of  
 340 data sources.

## 341 2.6 Achieving chemical agnosticism

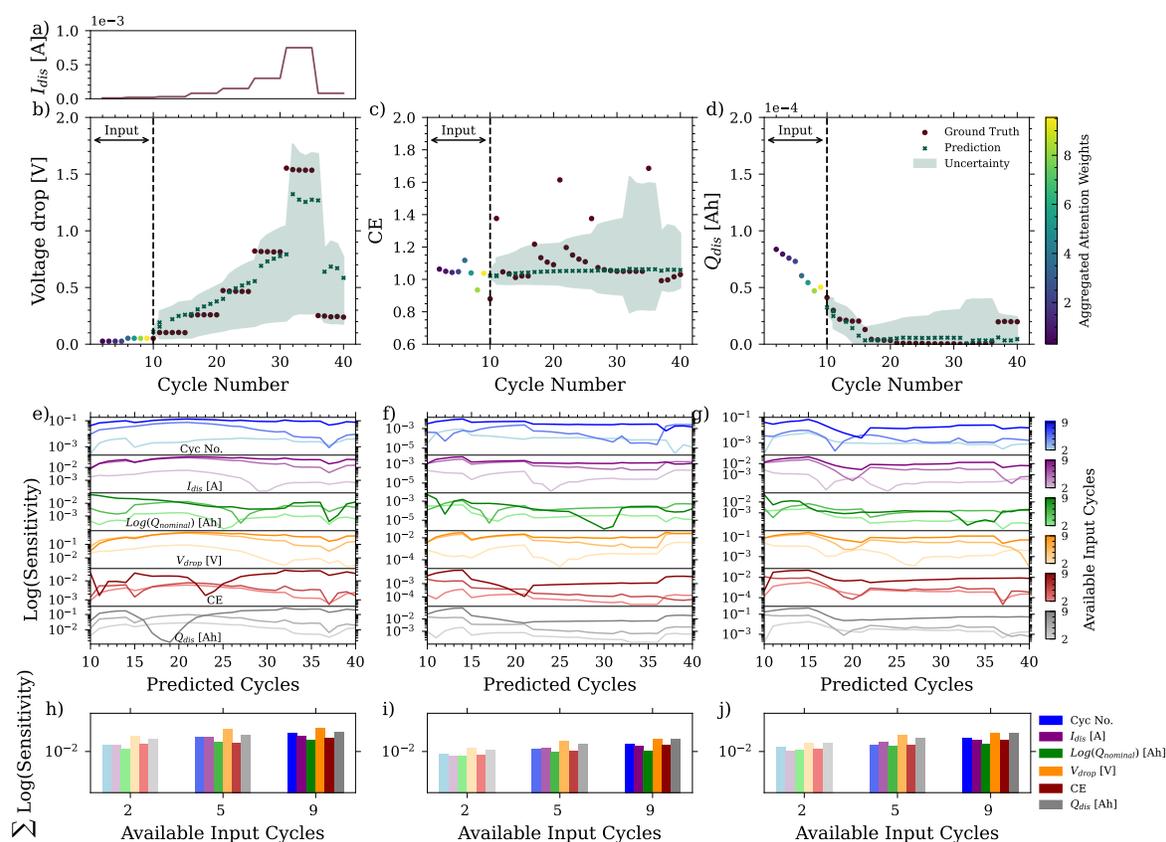


Figure 6: Analysis of  $M(P)_{Na}$ 's predictive accuracy and input sensitivity on Na-ion data. Plot a) presents the C-rate profile for cycling one battery, while plots b), c), and d) compare the model's prediction to actual data, showing consistency and adaptability. Sensitivity to input parameters across predicted cycles is analyzed in plots e), f), and g) on a logarithmic scale. The color intensity in these plots denotes the specific cycles from which the input parameter originates. Plots h), i), and j) show the sum of the logarithmic contribution of each input parameter towards predicting future cycles with a selective representation of three past cycle data. These visualizations confirm the model's attentive adjustment to the latest available input data and its capacity for generalization, despite the high experimental noise and limited battery performance.

342 ARCANa was so far demonstrated to generalize well across battery formats, electrolyte formu-  
 343 lations, cathode chemistries and cycling procedures for LIBs. The ultimate generalization would be  
 344 achieved if the model could also be deployed to Na-ion batteries. Since the underlying degradation  
 345 mechanism of Na-ion batteries is very different, we performed fine-tuning to test whether  $M(B)$  and  
 346  $M(P)$  are capable of achieving "true chemistry agnosticism" <sup>30, 16</sup>. These fine-tuned models are den-  
 347 oted  $M(B)_{Na}$  and  $M(P)_{Na}$ , and are trained on Na-ion cycling data with CC-CV and pulse discharge

348 settings. Details on the fine-tuning parameters are available in SI.

349 In Figures 6 and 7, we evaluate the fine-tuned  $M(B)_{Na}$  and  $M(P)_{Na}$  models on an unseen C-rate  
 350 test protocol (Fig. 6a and Fig. 7a). Both models demonstrate flexibility in adjusting to changes in  
 351 C-rates, with voltage drop, CE and  $Q_{dis}$  depicted in Fig. 6b-d and Fig. 7b-d. The model  $M(B)_{Na}$   
 352 shows narrower prediction intervals, indicative of lower uncertainty and greater predictive robustness.  
 353 This trend is consistent across all predictive dimensions and the model is probably benefiting from the  
 354 larger initial dataset on which it was trained, since it provided a richer learning environment for the  
 355 model to become more 'protocol-agnostic'. Its precision is especially notable in predicting the voltage  
 356 drop and CE estimations, closely following the ground-truth despite the substantial experimental noise.  
 357 The aggregated attention mechanism in  $M(B)_{Na}$  (Fig. 7d) also appears more fine-tuned, with greater  
 358 weights on the latest cycle data, which is consistent with its precise predictions. While  $M(P)_{Na}$  is  
 359 adaptable, it shows a marginally wider uncertainty (Fig. 6b-d).

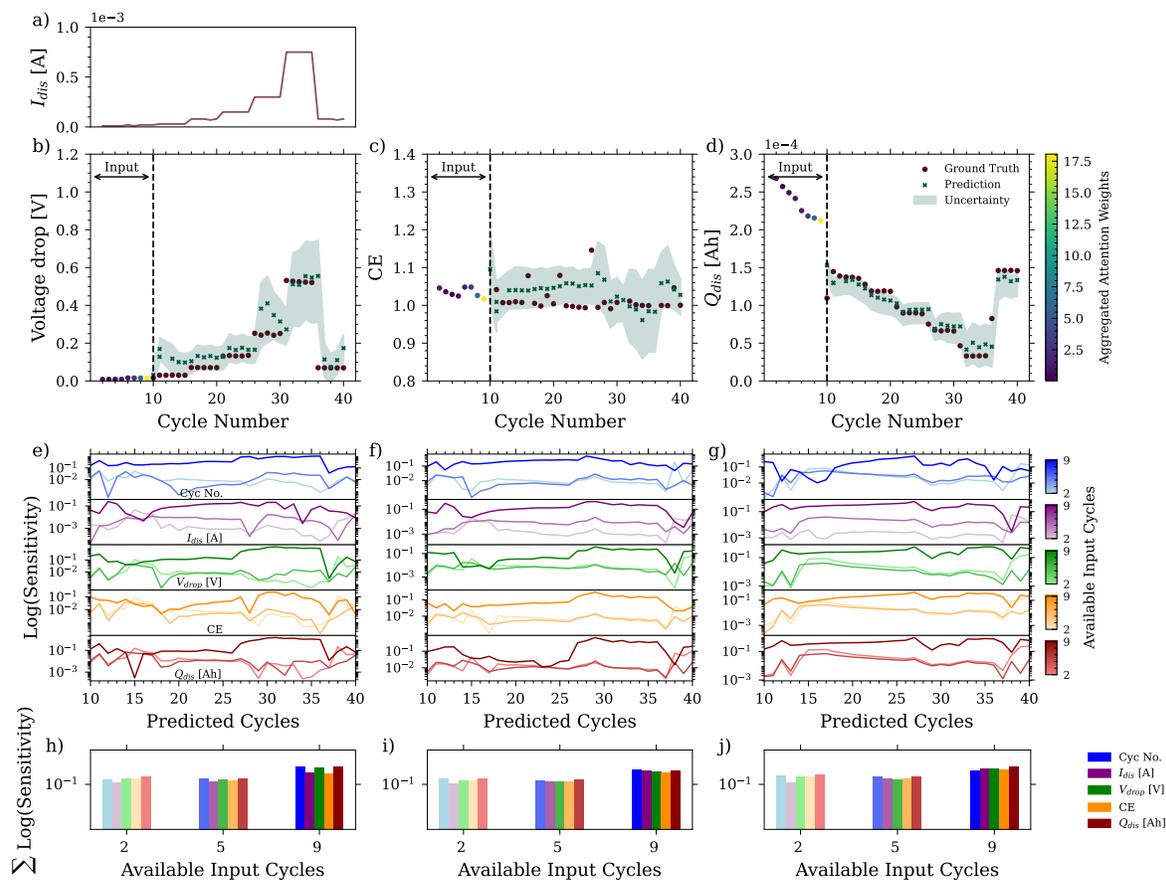


Figure 7: Evaluation of  $M(B)_{Na}$ 's predictive performance and input sensitivity on our own Na-ion data. Plot a) shows the discharge current profile, while plots b), c), and d) depict the predictions for voltage drop, CE and  $Q_{dis}$  against the ground truth. The colorbar here shows the aggregated attention weights across the input data. Plots e-g) provide a detailed logarithmic sensitivity analysis per predictive cycle for each input parameter, and plots h-j) aggregate these sensitivities, highlighting the model's focus on different input cycles, especially the most recent ones, reflecting  $M(B)_{Na}$ 's protocol adaptability and robust response to experimental noise.

360 Sensitivity analysis, as shown in Figures 7e-g and 6e-g evaluates the input parameter influence  
 361 on future predictions for  $M(B)_{Na}$  and  $M(P)_{Na}$ . Both models demonstrate increased sensitivity to  
 362 the most recent input data, i.e. cycles 7 to 9 in this provided example, aligned with their attention  
 363 distributions, with cycle 9 receiving the highest attention. This increased emphasis on the last input  
 364 cycles corresponds to the rapid degradation patterns in this sodium coin cell. As the model receives  
 365 each successive cycle, the most recent data, here cycle 9, becomes important in shaping its predictions,  
 366 allowing the model to more accurately predict ongoing trends.

367 In Fig. 7,  $M(B)_{Na}$  shows a greater overall sensitivity across input cycles, particularly for the di-  
 368 mensions of voltage drop and  $Q_{dis}$ . This is further illustrated in sensitivity profiles and cumulative  
 369 plots (Fig. 7h-j) highlighting a refined input-response relationship and a lower uncertainty interval in

370 the primary prediction (Fig. 7a-c). Such a distinct sensitivity indicates  $M(B)_{Na}$ 's ability to precisely  
 371 identify and respond to subtle variations. Despite the high experimental noise and limited battery per-  
 372 formance, the saliency and attention trends of both models remain remarkably similar. This suggests  
 373 that both mechanisms are intrinsic to the model's architecture, enabling them to perform consistently  
 374 at diverse scenarios.

375 To further substantiate our initial findings, the plots in Fig. 8, show both models'  $Q_{dis}$  predictions  
 376 aligning well with the ground-truth.  $M(P)_{Na}$  exhibit a tighter clustering around the actual values,  
 377 while  $M(B)_{Na}$  exhibits a broader spread. The prediction intervals and the distribution of quantiles  
 378 across the 10th and 90th percentile for both models confirm their consistency and calibrated confi-  
 379 dence. These evaluations provide insights into the model's robustness. The performance of  $M(B)_{Na}$ 's  
 380 especially underscores the advantage of extensive and diverse pretraining datasets in enhancing model  
 381 generalization across different battery chemistries.

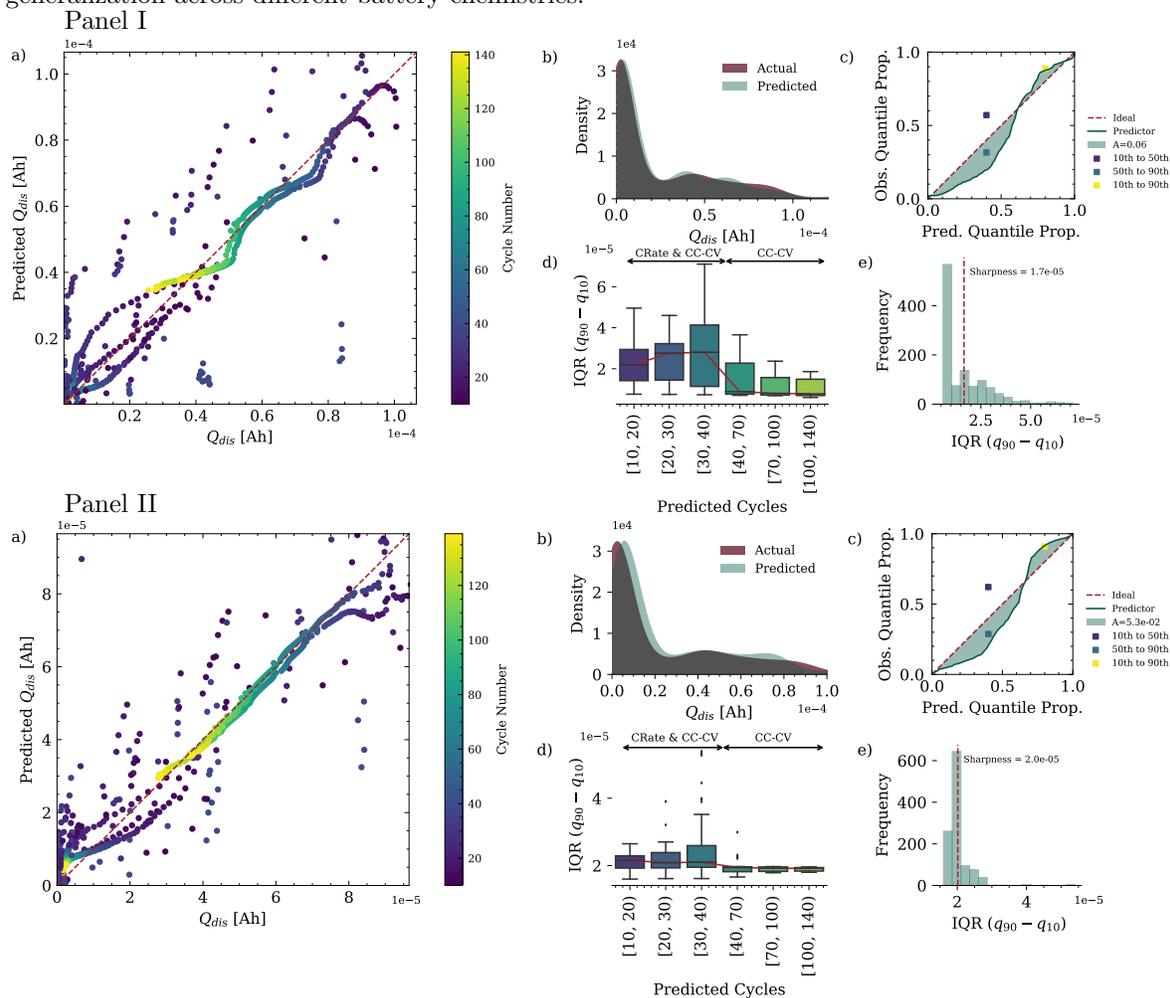


Figure 8: comparative  $Q_{dis}$  prediction analysis for Na-ion batteries using  $M(P)_{Na}$  (Panel I) and  $M(B)_{Na}$  (Panel II). The scatter plots a) illustrate the models' alignment with actual measurements. Density plots b) compare the distributions of predicted and actual values, demonstrating the models' accuracy in estimating  $Q_{dis}$ . Calibration plots in c) depict how well the predicted probabilities match the observed outcomes against the benchmark line, with the discrete points representing the observed proportions of actual values that fall within three quantile intervals. Both models demonstrate a pattern of marginal overconfidence below the 70th percentile and a slight underconfidence above this percentile, as evidenced by the calibration points's positions beneath and above the diagonal line, respectively.  $M(P)_{Na}$  shows a larger area of divergence,  $A = 0.06$ , while  $M(B)_{Na}$  presents a closer fit with a miscalibration of 0.053, highlighting both models' well-calibrated prediction capabilities across different chemistries. Boxplots d) visualize the spread and consistency of prediction intervals across predicted cycles. Histograms in e) represent the distribution of the quantile intervals of the models' prediction, highlighting uncertainty; these distributions indicate where, within the prediction range, the models' confidence is concentrated, with sharpness values of  $1.7 \times 10^{-5}$  for  $M(P)_{Na}$  and  $2.0 \times 10^{-5}$  for  $M(B)_{Na}$ , demonstrating a precise estimation of uncertainty.

### 3 Discussion

We demonstrated the chemistry-, format- and cycling procedure-agnostic ARCANA framework, and its ability to reliably monitor battery life and SOH by utilizing multitask learning with an attention mechanism. ARCANA excelled across three predictive settings, demonstrating that augmenting the model with diverse knowledge streams enhances its generalization across virtually all variations possible in batteries such as anode, cathode, electrolyte and shuttle ion chemistry and format. The ARCANA model integrates uncertainty quantification and attention mechanisms for each and every cycle to elucidate the model’s focus for each prediction and is essential for uncovering complex patterns associated across multiple factors. Further evaluation involves saliency and sensitivity assessments, allowing us to understand the impact of perturbation of input parameter on output predictions. By examining whether saliency and attention are directly correlated or orthogonal to each other, we gain a comprehensive understanding of input-output relationships, increasing the model’s explainability and reliability in extrapolation. Incorporating raw data and failed experiments, as suggested in prior studies<sup>45, 15</sup> is a deliberate strategy to teach our models to recognize variations across similar cell types and manufactures. This inclusion not only enables uncertainties to be quantified more accurately but also deepens reliability insights, reduces bias, and offers a more meaningful understanding of the data. A conceptually straightforward extension to this work would be to incorporate additional features, such as the rate of change of voltage with respect to capacity ( $dQ/dV$ )<sup>35, 12</sup>, and leverage different characterization methods like spectroscopy, to enhance the predictive power of the models. This will not only enhance multi-feature predictions, but also deepen the understanding of degradation processes<sup>15, 33, 51</sup>.

We observed that  $M(P)$ , trained on public data, offers broader generalization across various battery types and protocols, albeit with increased uncertainty.  $M(B)$ , trained on a more extensive dataset, demonstrates a lower uncertainty. This further motivates the importance of data sharing and management. Our findings also reveal that fine-tuning the models with few labels significantly improves their generalization to new chemistries, especially for  $M(B)$ . The methodology outlined in this paper presents an opportunity for other researchers to create their own high-performance models. By retraining or fine-tuning with different datasets, researchers can tailor these predictive models to their specific experimental setups and desired outcomes. This flexibility allows for the exploration of different perspectives and approaches, facilitating the development of more accurate and specialized models. One could envision a model-sharing and transfer-learning community similar to those found today in the fields of computer vision and language modeling. Furthermore, the performance metrics explored here raise the tantalizing prospect of further improving model quality via a federated learning approach. This could enable researchers from diverse backgrounds and institutions to pool their data and expertise, leading to more powerful models.

The modular design of the ARCANA pipeline enables real-time monitoring of battery degradation profiles, promoting timely and cost-effective interventions. This proactive approach prevents prolonged suboptimal testing conditions, improving the R&D process, and contributes to more informed material selection and protocol optimization. By automating data collection, processing, and analysis, researchers can streamline their experimental workflows and reduce human error. Furthermore, ML models can continuously learn from new data, adapt to evolving experimental conditions, and provide real-time insights. This integration of ML and laboratory workflows has the potential to transform battery research, enabling researchers to make data-driven decisions, uncover novel insights more rapidly and accelerating the pace of discovery.

Overall, we demonstrated that incorporating multitask learning with an attention mechanism creates a framework that can achieve chemistry agnosticism as envisioned by Battery 2030+<sup>5</sup> and the interesting fact that a DL architecture trained on a smaller, noisier, but more diverse dataset yields better generalization at the cost of higher uncertainty. We hope that the pipeline will emerge as an indispensable and transformative tool to bridge the gap between lab-scale research and commercial viability, and will become essential for development of applications and insightful predictive models in the energy storage field.

## 4 Methods

### 4.1 Model compartments Dynamics

In the following section, some of the key components of the ARCANA framework are explained to underscore their contribution to the overall efficacy and reliability of the model. This includes an exploration of attention mechanisms, a teacher forcing scheduler, methods to quantify predictive uncertainty, a strategic early stopping protocol, a training procedure, and evaluation metrics.

#### Attention mechanism

Within the proposed ARCANA framework, two distinct attention mechanisms are implemented. The first, termed additive attention, is also known as Bahdanau attention<sup>8</sup>. This mechanism aligns the hidden state of the decoder  $h_t$  at each time step  $t$  with the hidden states of the encoder ( $h_s$ ), thus producing a context vector that encapsulates the weighted relevance of each historical temporal segment from the past cycles. This vector provides a dynamically focused representation of the input sequence pertinent to the current decoding step. This mechanism is functional through a parameterized attention model. The model calculates an attention score  $e_{ts}$  (Eq. 1) for each encoder state  $h_s$  given by:

$$e_{ts} = v^T \tanh(W_1 h_t + w_2 h_s) \quad (1)$$

where  $W_1$  and  $W_2$  are the weight matrices that transform the respective hidden states into a common feature space and  $v$  is a weight vector that projects the activated sum into a scalar score. Attention weights  $\alpha_{ts}$  are then determined by normalizing these scores using the softmax function (Eq. 2).

$$\alpha_{ts} = \frac{\exp(e_{ts})}{\sum_{k=1}^{T_e} \exp(e_{tk})} \quad (2)$$

here,  $T_e$  is the total number of time steps in the encoder sequence.

The context vector  $c_t$  results from aggregating the encoder hidden states, each weighted by its respective attention weight, as can be seen in Eq. 3, and can improve the model's capacity for handling Seq-to-Seq predictions<sup>41</sup>.

$$c_t = \sum_{s=1}^{T_e} \alpha_{ts} h_s \quad (3)$$

Another attention mechanism that can be employed within the ARCANA architecture is multihead attention. This mechanism expands the model's capacity to focus on different positions of the input sequence simultaneously<sup>49</sup>, which is crucial for capturing a wider range of dependencies inherent in battery lifetime data. This attention mechanism operates by projecting the decoder's hidden states and the encoder outputs, representing the past cycle's information, into multiple subspaces. This is formulated as: (Eq. 4)

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^0 \quad (4)$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (5)$$

where each head ( $\text{head}_i$ ) captures different aspects of the input data and is computed as shown in Eq.5. The operation applied in each head is defined by the attention of the scaled dot product and is presented in Eq. 6.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \quad (6)$$

Here,  $Q$ ,  $K$ , and  $V$  are the query, key, and value matrices, respectively.  $Q$  is generated from the hidden states of the decoder, while  $K$  and  $V$  are derived from the encoder outputs. This arrangement enables the decoder to integrate the current state information with historical data provided by the encoder. The parameter matrices  $W_i^Q$ ,  $W_i^K$ , and  $W_i^V$  for each head  $i$ , along with the output weight matrix  $W^0$ , are optimized during the training process. These matrices are instrumental in transforming the input data into different representational subspaces to capture various aspects and dependencies within the data. The parameter  $d_k$ , representing the dimension of the key vectors, scales the dot product within the attention mechanism. In Eq. 6, the softmax function is applied to these scaled attention scores, which originate from the interactions between the query and key matrices. This process results in the production of a context vector, which integrates information from different representational subspaces and allows the model to consider multiple aspects of historical data<sup>65, 60</sup>.

## Teacher forcing

Teacher forcing optimizes the learning of temporal dependencies. By integrating the real data from previous time steps, the technique promotes rapid stabilization and convergence of the model. In the present study, the implementation of the teacher forcing strategy is applied through a calculated division of training epochs. This division is reflective of the model's incremental improvement in processing sequences with varying lengths over time by prioritizing shorter sequences at the early stages of training to ensure intensive guidance. This preferential focus ensures that the model does not prematurely plateau when learning to predict longer-term dependencies.

To quantitatively define this approach, the training period consisting of  $E$  epochs is divided into  $D$  equal segments  $s$ . Within the  $i$ -th segment, the teacher forcing ratio is adjusted through a decay parameter  $\lambda$ , which represents how quickly the training procedure switches from using real data as decoder inputs to using model predictions from the previous cycle, as depicted in Fig.2b. The allocation of epochs per division  $d_i$  is calculated as can be seen in Eq. 7

$$d_i = \text{round}\left(\frac{s \cdot e^{-\lambda i}}{\sum_{j=0}^{D-1} s \cdot e^{-\lambda j}} \cdot E\right) \quad (7)$$

Following this, the teacher forcing ratio for the  $t$ -th epoch in the  $i$ -th segment is linearly reduced from a starting ratio  $R_{start}$  to an ending ratio  $R_{end}$ , using the following equation, Eq. 8.

$$A = \left(\frac{R_{start} - R_{end}}{d_i + \epsilon}\right) \quad (8)$$
$$R_{t_i} = R_{start} - A \cdot t$$

Here,  $R_{t_i}$  indicates the teacher forcing ratio at epoch  $t$  for the  $i$ -th segment. The expression  $A$  represents the decrease per epoch in that segment. To ensure numerical stability and avoid division by zero, a small constant  $\epsilon$ , set to  $10^{-8}$ , is included in the calculation as indicated in Eq. 8. The teacher forcing ratio, as a probabilistic measure, represents the likelihood that the model will utilize the actual observation from the training data at a given prediction step. This approach modulates the ratio to facilitate a smooth transition from guided to self-generated sequence prediction. The adjusted ratios are indicative of the model's learning trajectory, enhancing its independent predictive accuracy across different sequence lengths.

## Uncertainty quantification

The pinball loss, in this study, provides a robust metric for predicting a range of potential outcomes, rather than a single point estimation. This is an effective measure for forecasting scenarios where the impacts of overprediction and underprediction are asymmetric<sup>57</sup>. It is defined for a set of quantiles  $Q = \{q_1, q_2, q_3\}$  where  $q_1 < q_2 < q_3$  and in this study, we select  $Q = \{0.1, 0.5, 0.9\}$  corresponding to the 10-th, 50-th, and 90-th percentiles, respectively. For a given predicted value  $\hat{y}$  and the actual target value  $y$ , the pinball loss for a single quantile  $q$  is calculated as:

$$L_q(\hat{y}, y) = \begin{cases} (1 - q) \cdot (\hat{y} - y) & \text{if } y < \hat{y} \\ q \cdot (y - \hat{y}) & \text{if } y \geq \hat{y} \end{cases} \quad (9)$$

In the implementation of this loss function, a mask is provided and applied to each quantile's loss to selectively evaluate certain predictions, allowing for the exclusion of outliers. The total pinball loss for multiple quantiles is then the sum of the individual losses for each quantile, averaged over all predictions, as shown in Eq. 10, reflecting the model's performance across the specified range of quantiles.

$$L(Q, \hat{Y}, Y) = \frac{1}{N} \sum_{i=1}^N \sum_{q \in Q} L_q(\hat{y}_{qi}, y_i) \quad (10)$$

Here,  $N$  is the number of observations,  $\hat{Y}$  is a stack of vectors, with each vector containing the predictions for all observations at one of the specified quantiles, and  $Y$  is the vector of the true target values. Each element  $\hat{y}_{qi}$  in  $\hat{Y}$  denotes the predicted value for the  $i$ -th observation at quantile  $q$ . This configuration not only facilitates efficient computation of the loss function across multiple quantiles and observations, but also captures the central tendency and variability of the predictions, making it a comprehensive loss function for probabilistic forecasting<sup>57, 37</sup>.

## 511 Early stopping

512 To optimize training, a rigorous early stopping approach is incorporated. This method was originally  
513 proposed by Prechelt et al. <sup>44</sup> and combines criteria to prevent overfitting while ensuring substantial training  
514 progress, especially in the presence of noisy data. Here, a dual-criteria strategy is implemented. The first  
515 criterion assesses the ratio between generalization loss (GL) and training progress, which is shown in Eq.11,  
516 where  $E_{val}$  represents the validation error at the current epoch,  $E_{min\ val}$  is the lowest validation error obtained  
517 up to the current epoch, and  $E_{train\ strip}$  denotes the training errors within a recent sequence of epochs. This  
518 sequence, or strip, is a designated period in which progress quotient (PQ) is measured. If the generalization-  
519 loss-to-progress-quotient-ratio (GL/PQ) surpasses a predefined value, it may indicate that further training will  
520 not be beneficial for the model’s generalizability.

$$512 \quad GL = 100 \cdot \left( \frac{E_{val}}{E_{min\ val}} - 1 \right) \quad (11)$$
$$513 \quad PQ = 1000 \cdot \left( \frac{Mean(E_{train\ strip})}{Min(E_{train\ strip})} - 1 \right)$$

521 The second criterion implements a conventional check and is applied to monitor the trend in validation  
522 error. An increased trend over the epoch sequence suggests that overfitting could be occurring. Training  
523 is discontinued when both the ratio criterion and the error-trend criterion indicate that further training is  
524 unlikely to yield significant gains. In general, this strategy offers a control mechanism that aligns the duration  
525 of training with the achievement of a well-generalized model capable of accurate predictions.

## 526 Training Procedure

527 Expanding on Seq-to-Seq integration, the training phase begins by initializing the data loaders for batch  
528 processing and configuring the parameters of the Seq-to-Seq model, the loss criteria, the optimizer, and a  
529 dynamic learning rate scheduler <sup>26</sup>. Hyperparameter optimization, through a series of trials using Optuna’s <sup>4</sup>  
530 Tree-structured Parzen Estimator (TPE) Sampler, employs a probabilistic model to specify the most promising  
531 parameter configuration, navigating the search space while balancing exploration and exploitation within a  
532 complex and high-dimensional domain <sup>10</sup>. Training unfolds over several epochs, with each iteration starting  
533 with a reset of the model’s hidden states and zeroing gradients to ensure clean computation for the forward  
534 pass. The pinball loss function is selected for its effectiveness in probabilistic forecasting, eliminating the need  
535 for a presumptive data distribution model <sup>37</sup> unlike traditional metrics <sup>57</sup>, which are more sensitive to noise and  
536 anomalies. These asymmetric and non-parametric criteria assess forecast accuracy by penalizing deviations  
537 from three targeted quantiles, namely 0.1, 0.5, and 0.9, enhancing robustness to outliers and the efficacy for  
538 LSTM-based networks <sup>57</sup>. At the same time, a masking technique <sup>33</sup> is implemented to filter out padding-  
539 induced distortions from the loss calculation, ensuring the integrity of the learning signal. Backpropagation  
540 follows loss computation, incorporating gradient clipping to prevent divergence and gradient explosion in  
541 recurrent network architectures. Additionally, learning rate adjustments encourage robust convergence. The  
542 validation phase alternates with training, where performance is assessed and early stopping criteria are applied  
543 to mitigate overfitting. Optuna enhances optimization by pruning the less promising trials. Once the training  
544 is completed, the model parameters are saved and a comprehensive report is generated detailing the training  
545 results. The training procedure steps described are schematically depicted in Supp.1

## 546 Evaluation metrics

547 For this study, the following metrics are implemented, including both average errors and variability of indi-  
548 vidual predictions, to evaluate the performance of the model. These metrics are RMSE (Eq. 12) which provides  
549 a measure of the magnitude of prediction errors, MAPE (Eq. 13) which measures the average magnitude of  
550 errors as a percentage, medAE (Eq. 14) to capture the median error, reducing the influence of outliers, and  
551 mean absolute error (MAE) (Eq. 15) which represents the mean absolute differences.

$$547 \quad RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (12)$$

$$548 \quad MAPE = \frac{100\%}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \quad (13)$$

$$549 \quad medAE = \text{median}(|y_i - \hat{y}_i| : i = 1, 2, \dots, n) \quad (14)$$

$$550 \quad MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (15)$$

## Data and code availability

Open source data supporting the findings of this study are available online, with access details provided in Table. 1. The ARCANA framework can be installed using `pip install arcana` or cloned from <https://github.com/basf/ARCANA>. In addition, public pre-trained model weights can be accessed at <https://doi.org/10.5281/zenodo.10293072>.

## Acknowledgements

This work contributes to the research performed at CELEST (Center for Electrochemical Energy Storage Ulm-Karlsruhe) and was partly funded by the German Research Foundation (DFG) under Project ID 390874152 (POLiS Cluster of Excellence). This project also received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 957189 (BATTERY 2030+).

## Author contributions statement

K.M. and B.B. provided the comprehensive BASF dataset, L.M. and L.N. conducted all cycling data for Li-ion and Na-ion batteries at KIT/ IPC respectively. Data assembly, data cleaning, model idea including the design architecture, implementation, training, evaluation, and package creation is conducted by F.R.. R.L., D.L. supervised the model development. K.M., B.B., H.S., R.L, and D.L. supervised this research. All authors reviewed the manuscript.

## Conflicts of interest

The authors have no competing or financial interests to declare.

## References

- [1] Toyota Research Institute. Experimental data platform. 2021. *Project Closed-loop optimization of extreme fast charging for batteries using machine learning*. URL: <https://data.matr.io/1/projects/5d80e633f405260001c0b60a>.
- [2] Toyota Research Institute. Experimental data platform. 2021. *Project Data-driven prediction of battery cycle life before capacity degradation*. URL: <https://data.matr.io/1/projects/5c48dd2bc625d700019f3204>.
- [3] Haruna Adamu et al. “Artificial intelligence-navigated development of high-performance electrochemical energy storage systems through feature engineering of multiple descriptor families of materials”. In: *Energy Advances* 2.5 (2023), pp. 615–645. DOI: [10.1039/D3YA00104K](https://doi.org/10.1039/D3YA00104K).
- [4] Takuya Akiba et al. “Optuna: A next-generation hyperparameter optimization framework”. In: *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. Ed. by Association for Computing Machinery. Royal Society of Chemistry, 2019, pp. 2623–2631. DOI: [10.1145/3292500.3330701](https://doi.org/10.1145/3292500.3330701).
- [5] Julia Amici et al. “A roadmap for transforming research to invent the batteries of the future designed within the European large scale research initiative BATTERY 2030+”. In: *Advanced energy materials* 12.17 (2022), p. 2102785. DOI: [10.1002/aem.202102785](https://doi.org/10.1002/aem.202102785).
- [6] Peter M Attia et al. ““Knees” in lithium-ion battery aging trajectories”. In: *Journal of The Electrochemical Society* 169.6 (2022), p. 060517. DOI: [10.1149/1945-7111/ac6d13](https://doi.org/10.1149/1945-7111/ac6d13).

- 591 [7] Peter M Attia et al. “Closed-loop optimization of fast-charging protocols for batteries with  
592 machine learning”. In: *Nature* 578.7795 (2020), pp. 397–402. DOI: [10.1038/s41586-020-1994-  
593 5](https://doi.org/10.1038/s41586-020-1994-5).
- 594 [8] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. “Neural Machine Translation by  
595 Jointly Learning to Align and Translate”. In: (2016). DOI: [10.48550/arXiv.1409.0473](https://doi.org/10.48550/arXiv.1409.0473). arXiv:  
596 [1409.0473](https://arxiv.org/abs/1409.0473).
- 597 [9] Thorsten Baumhöfer et al. “Production caused variation in capacity aging trend and correlation  
598 to initial cell performance”. In: *Journal of Power Sources* 247 (2014), pp. 332–338. DOI: [10.  
599 1016/j.jpowsour.2013.08.108](https://doi.org/10.1016/j.jpowsour.2013.08.108).
- 600 [10] James Bergstra et al. “Algorithms for Hyper-Parameter Optimization”. In: *Advances in Neu-  
601 ral Information Processing Systems*. Ed. by J. Shawe-Taylor et al. Vol. 24. Curran Associates,  
602 Inc., 2011. URL: [https://proceedings.neurips.cc/paper\\_files/paper/2011/file/  
603 86e8f7ab32cfd12577bc2619bc635690-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf).
- 604 [11] Arghya Bhowmik et al. “Implications of the BATTERY 2030+ AI-Assisted Toolkit on Future  
605 Low-TRL Battery Discoveries and Chemistries”. In: *Advanced Energy Materials* 12.17 (2022),  
606 p. 2102698. DOI: [10.1002/aenm.202102698](https://doi.org/10.1002/aenm.202102698).
- 607 [12] Ira Bloom et al. “Differential voltage analyses of high-power, lithium-ion cells: 1. Technique  
608 and application”. In: *Journal of Power Sources* 139.1-2 (2005), pp. 295–303. DOI: [10.1016/j.  
609 jpowsour.2004.07.021](https://doi.org/10.1016/j.jpowsour.2004.07.021).
- 610 [13] JC Burns et al. “Evaluation of effects of additives in wound Li-ion cells through high precision  
611 coulometry”. In: *Journal of The Electrochemical Society* 158.3 (2011), A255. DOI: [10.1149/1.  
612 3531997](https://doi.org/10.1149/1.3531997).
- 613 [14] JC Burns et al. “Predicting and extending the lifetime of Li-ion batteries”. In: *Journal of The  
614 Electrochemical Society* 160.9 (2013), A1451. DOI: [10.1149/2.060309jes](https://doi.org/10.1149/2.060309jes).
- 615 [15] Yunhong Che et al. “Health prognostics for lithium-ion batteries: mechanisms, methods, and  
616 prospects”. In: *Energy & Environmental Science* 16.2 (2023), pp. 338–371. DOI: [10.1039/  
617 D2EE03019E](https://doi.org/10.1039/D2EE03019E).
- 618 [16] Guzhong Chen et al. “Generalizing property prediction of ionic liquids from limited labeled data:  
619 a one-stop framework empowered by transfer learning”. In: *Digital Discovery* 2.3 (2023), pp. 591–  
620 601. DOI: [10.1039/D3DD00040K](https://doi.org/10.1039/D3DD00040K).
- 621 [17] Tim Cooijmans et al. “Recurrent Batch Normalization”. In: (2017). arXiv: [1603.09025 \[cs.LG\]](https://arxiv.org/abs/1603.09025).
- 622 [18] JR Dahn, JC Burns, and DA Stevens. “Importance of coulombic efficiency measurements in  
623 R&D efforts to obtain long-lived Li-ion batteries”. In: *The Electrochemical Society Interface* 25.3  
624 (2016), p. 75. DOI: [10.1149/2.F07163if](https://doi.org/10.1149/2.F07163if).
- 625 [19] Valerio De Angelis, Yuliya Preger, and Babu R Chalamala. *Battery Lifecycle Framework: A  
626 Flexible Repository and Visualization Tool for Battery Data from Materials Development to Field  
627 Implementation*. Mar. 2021. DOI: [10.1149/osf.io/h7c24](https://doi.org/10.1149/osf.io/h7c24). URL: [osf.io/preprints/ecsarxiv/  
628 h7c24](https://osf.io/preprints/ecsarxiv/h7c24).
- 629 [20] Zhongwei Deng et al. “Battery health estimation with degradation pattern recognition and trans-  
630 fer learning”. In: *Journal of Power Sources* 525 (2022), p. 231027. DOI: [10.1016/j.jpowsour.  
631 2022.231027](https://doi.org/10.1016/j.jpowsour.2022.231027).
- 632 [21] Sandia National Laboratories Grid Energy Storage Department. *Battery Archive. Homepage of  
633 Battery Archive*. [Online; accessed 31-July-2023]. 2021. URL: [https://www.batteryarchive.  
634 org](https://www.batteryarchive.org).
- 635 [22] Gonçalo Dos Reis et al. “Lithium-ion battery data and where to find it”. In: *Energy and AI* 5  
636 (2021), p. 100081. DOI: [10.1016/j.egyai.2021.100081](https://doi.org/10.1016/j.egyai.2021.100081).
- 637 [23] Claudia Draxl and Matthias Scheffler. “NOMAD: The FAIR concept for big data-driven materials  
638 science”. In: *Mrs Bulletin* 43.9 (2018), pp. 676–682. DOI: [10.1557/mrs.2018.208](https://doi.org/10.1557/mrs.2018.208).
- 639 [24] Maximilian Fichtner et al. “Rechargeable batteries of the future - the state of the art from a  
640 BATTERY 2030+ perspective”. In: *Advanced Energy Materials* 12.17 (2022), p. 2102904.

- 641 [25] Tilmann Gneiting, Fadoua Balabdaoui, and Adrian E Raftery. “Probabilistic forecasts, calibration and sharpness”. In: *Journal of the Royal Statistical Society Series B: Statistical Methodology* 69.2 (2007), pp. 243–268. DOI: [10.1111/j.1467-9868.2007.00587.x](https://doi.org/10.1111/j.1467-9868.2007.00587.x).
- 642
- 643
- 644 [26] Yoav Goldberg. “A primer on neural network models for natural language processing”. In: *Journal of Artificial Intelligence Research* 57 (2016), pp. 345–420. DOI: [10.1613/jair.4992](https://doi.org/10.1613/jair.4992).
- 645
- 646 [27] Qingrui Gong, Ping Wang, and Ze Cheng. “An encoder-decoder model based on deep learning for state of health estimation of lithium-ion battery”. In: *Journal of Energy Storage* 46 (2022), p. 103804. DOI: [10.1016/j.est.2021.103804](https://doi.org/10.1016/j.est.2021.103804).
- 647
- 648
- 649 [28] Chuan Guo et al. “On Calibration of Modern Neural Networks”. In: *Proceedings of the 34th International Conference on Machine Learning*. Ed. by Doina Precup and Yee Whye Teh. Vol. 70. Proceedings of Machine Learning Research. PMLR, Aug. 2017, pp. 1321–1330. URL: <https://proceedings.mlr.press/v70/guo17a.html>.
- 650
- 651
- 652
- 653 [29] CALCE battery research group homepage. URL: <https://calce.umd.edu/data>.
- 654 [30] Swarn Jha et al. “Learning-assisted Materials Development and Device Management in Batteries and Supercapacitors: Performance Comparison and Challenges”. In: *Journal of Materials Chemistry A* 11.8 (2023), pp. 3904–3936. DOI: [10.1039/D2TA07148G](https://doi.org/10.1039/D2TA07148G).
- 655
- 656
- 657 [31] MM Kabir and Dervis Emre Demirocak. “Degradation mechanisms in Li-ion batteries: a state-of-the-art review”. In: *International Journal of Energy Research* 41.14 (2017), pp. 1963–1986. DOI: [10.1002/er.3762](https://doi.org/10.1002/er.3762).
- 658
- 659
- 660 [32] Weihan Li et al. “Digital twin for battery systems: Cloud battery management system with online state-of-charge and state-of-health estimation”. In: *Journal of energy storage* 30 (2020), p. 101557. DOI: [10.1016/j.est.2020.101557](https://doi.org/10.1016/j.est.2020.101557).
- 661
- 662
- 663 [33] Weihan Li et al. “Forecasting battery capacity and power degradation with multi-task learning”. In: *Energy Storage Materials* 53 (2022), pp. 453–466. DOI: [10.1016/j.ensm.2022.09.013](https://doi.org/10.1016/j.ensm.2022.09.013).
- 664
- 665 [34] Weihan Li et al. “One-shot battery degradation trajectory prediction with deep learning”. In: *Journal of Power Sources* 506 (2021), p. 230024. DOI: [10.1016/j.jpowsour.2021.230024](https://doi.org/10.1016/j.jpowsour.2021.230024).
- 666
- 667 [35] Xiaoyu Li, Zhenpo Wang, and Jinying Yan. “Prognostic health condition for lithium battery using the partial incremental capacity and Gaussian process regression”. In: *Journal of power sources* 421 (2019), pp. 56–67. DOI: [10.1016/j.jpowsour.2019.03.008](https://doi.org/10.1016/j.jpowsour.2019.03.008).
- 668
- 669
- 670 [36] Chen Ling. “A review of the recent progress in battery informatics”. In: *npj Computational Materials* 8.1 (2022), p. 33. DOI: [10.1038/s41524-022-00713-x](https://doi.org/10.1038/s41524-022-00713-x).
- 671
- 672 [37] Bidong Liu et al. “Probabilistic load forecasting via quantile regression averaging on sister forecasts”. In: *IEEE Transactions on Smart Grid* 8.2 (2015), pp. 730–737. DOI: [10.1109/TSG.2015.2437877](https://doi.org/10.1109/TSG.2015.2437877).
- 673
- 674
- 675 [38] Leon Merker. *2023 Commercial Coincell 45mAh*. Version 1.0. Nov. 2023. DOI: [10.5281/zenodo.10102627](https://doi.org/10.5281/zenodo.10102627). URL: <https://doi.org/10.5281/zenodo.10102627>.
- 676
- 677 [39] Leon Merker. *InZePro InForm 300 Cycles CCCV after EOL*. Version 1.0. Nov. 2023. DOI: [10.5281/zenodo.10102508](https://doi.org/10.5281/zenodo.10102508). URL: <https://doi.org/10.5281/zenodo.10102508>.
- 678
- 679 [40] Man-Fai Ng et al. “Predicting the state of charge and health of batteries using data-driven machine learning”. In: *Nature Machine Intelligence* 2.3 (2020), pp. 161–170. DOI: [10.1038/s42256-020-0156-7](https://doi.org/10.1038/s42256-020-0156-7).
- 680
- 681
- 682 [41] Zhaoyang Niu, Guoqiang Zhong, and Hui Yu. “A review on the attention mechanism of deep learning”. In: *Neurocomputing* 452 (2021), pp. 48–62. DOI: [10.1016/j.neucom.2021.03.091](https://doi.org/10.1016/j.neucom.2021.03.091).
- 683
- 684 [42] Leah Nuss et al. *Formation and cycling data for Na-ion batteries from high-throughput synthesis, coating, and assembly*. Version v1. May 2023. DOI: [10.5281/zenodo.7981011](https://doi.org/10.5281/zenodo.7981011). URL: <https://doi.org/10.5281/zenodo.7981011>.
- 685
- 686
- 687 [43] Adam Paszke et al. “PyTorch: An Imperative Style, High-Performance Deep Learning Library”. In: *Advances in Neural Information Processing Systems* 32. Curran Associates, Inc., 2019, pp. 8024–8035. URL: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- 688
- 689
- 690

- [44] Lutz Prechelt. “Early stopping-but when?” In: *Neural Networks: Tricks of the trade*. Springer, 2002, pp. 55–69. DOI: [10.1007/978-3-642-35289-8\\_5](https://doi.org/10.1007/978-3-642-35289-8_5).
- [45] Paul Raccuglia et al. “Machine-learning-assisted materials discovery using failed experiments”. In: *Nature* 533.7601 (2016), pp. 73–76. DOI: [10.1038/nature17439](https://doi.org/10.1038/nature17439).
- [46] Fuzhan Rahmanian et al. “Conductivity experiments for electrolyte formulations and their automated analysis”. In: *Scientific Data* 10.1 (2023), p. 43. DOI: [10.1038/s41597-023-01936-3](https://doi.org/10.1038/s41597-023-01936-3).
- [47] Laura Hannemose Rieger et al. “Uncertainty-aware and explainable machine learning for early prediction of battery degradation trajectory”. In: *Digital Discovery* 2.1 (2023). DOI: [10.1039/D2DD00067A](https://doi.org/10.1039/D2DD00067A).
- [48] Darius Roman et al. “Machine learning pipeline for battery state-of-health estimation”. In: *Nature Machine Intelligence* 3.5 (2021), pp. 447–456. DOI: [10.1038/s42256-021-00312-3](https://doi.org/10.1038/s42256-021-00312-3).
- [49] Jerret Ross et al. “Large-scale chemical language representations capture molecular structure and properties”. In: *Nature Machine Intelligence* 4.12 (2022), pp. 1256–1264. DOI: [10.1038/s42256-022-00580-7](https://doi.org/10.1038/s42256-022-00580-7).
- [50] Bhaskar Saha and Kai Goebel. *NASA Prognostics Data Repository*. 2007. URL: <https://www.nasa.gov/content/prognostics-center-of-excellence-data-set-repository>.
- [51] Kristen A Severson et al. “Data-driven prediction of battery cycle life before capacity degradation”. In: *Nature Energy* 4.5 (2019), pp. 383–391. DOI: [10.1038/s41560-019-0356-8](https://doi.org/10.1038/s41560-019-0356-8).
- [52] AJ Smith et al. “Precision measurements of the coulombic efficiency of lithium-ion batteries and of electrode materials for lithium-ion batteries”. In: *Journal of The Electrochemical Society* 157.2 (2009), A196. DOI: [10.1149/1.3268129](https://doi.org/10.1149/1.3268129).
- [53] Anna Smith et al. “Potential and limitations of research battery cell types for electrochemical data acquisition”. In: *Batteries & Supercaps* 6.6 (2023), e202300080. DOI: [10.1002/batt.202300080](https://doi.org/10.1002/batt.202300080).
- [54] Helge Sören Stein. “Nonlinear potentiodynamic battery charging protocols for fun, education, and application”. In: *ChemRxiv* (2023). DOI: [10.26434/chemrxiv-2023-vj5n0](https://doi.org/10.26434/chemrxiv-2023-vj5n0).
- [55] Calum Strange and Goncalo Dos Reis. “Prediction of future capacity and internal resistance of Li-ion cells from one cycle of input data”. In: *Energy and AI* 5 (2021), p. 100097. DOI: [10.1016/j.egyai.2021.100097](https://doi.org/10.1016/j.egyai.2021.100097).
- [56] Zheming Tong et al. “Early prediction of remaining useful life for Lithium-ion batteries based on a hybrid machine learning method”. In: *Journal of Cleaner Production* 317 (2021), p. 128265. DOI: [10.1016/j.jclepro.2021.128265](https://doi.org/10.1016/j.jclepro.2021.128265).
- [57] Yi Wang et al. “Probabilistic individual load forecasting using pinball loss guided LSTM”. In: *Applied Energy* 235 (2019), pp. 10–20. DOI: [10.1016/j.apenergy.2018.10.078](https://doi.org/10.1016/j.apenergy.2018.10.078).
- [58] Mark D Wilkinson et al. “The FAIR Guiding Principles for scientific data management and stewardship”. In: *Scientific Data* 3.1 (2016), pp. 1–9. DOI: [10.1038/sdata.2016.18](https://doi.org/10.1038/sdata.2016.18).
- [59] Billy Wu et al. “Battery digital twins: Perspectives on the fusion of models, data and artificial intelligence for smart battery management systems”. In: *Energy and AI* 1 (2020), p. 100016. DOI: [10.1016/j.egyai.2020.100016](https://doi.org/10.1016/j.egyai.2020.100016).
- [60] Changwen Xu, Yuyang Wang, and Amir Barati Farimani. “TransPolymer: a Transformer-based language model for polymer property predictions”. In: *npj Computational Materials* 9.1 (2023), p. 64. DOI: [10.1038/s41524-023-01016-5](https://doi.org/10.1038/s41524-023-01016-5).
- [61] Yuzhi Xu, Jiankai Ge, and Cheng-Wei Ju. “Machine learning in energy chemistry: introduction, challenges and perspectives”. In: *Energy Advances* 2.7 (2023), pp. 896–921. DOI: [10.1039/D3YA00057E](https://doi.org/10.1039/D3YA00057E).
- [62] Fangfang Yang et al. “A coulombic efficiency-based model for prognostics and health estimation of lithium-ion batteries”. In: *Energy* 171 (2019), pp. 1173–1182. DOI: [10.1016/j.energy.2019.01.083](https://doi.org/10.1016/j.energy.2019.01.083).
- [63] Yixin Yang. “A machine-learning prediction method of lithium-ion battery life based on charge process for different applications”. In: *Applied Energy* 292 (2021), p. 116897. DOI: [10.1016/j.apenergy.2021.116897](https://doi.org/10.1016/j.apenergy.2021.116897).

- 742 [64] Zhenpeng Yao et al. “Machine learning for a sustainable energy future”. In: *Nature Reviews*  
743 *Materials* 8.3 (2023), pp. 202–215. DOI: [10.1038/s41578-022-00490-5](https://doi.org/10.1038/s41578-022-00490-5).
- 744 [65] Jaekyun Yoo et al. “An artificial neural network using multi-head intermolecular attention for  
745 predicting chemical reactivity of organic materials”. In: *Journal of Materials Chemistry A* 11.24  
746 (2023), pp. 12784–12792. DOI: [10.1039/D2TA07660H](https://doi.org/10.1039/D2TA07660H).
- 747 [66] Yunwei Zhang et al. “Identifying degradation patterns of lithium ion batteries from impedance  
748 spectroscopy using machine learning”. In: *Nature communications* 11.1 (2020), p. 1706. DOI:  
749 [10.1038/s41467-020-15235-7](https://doi.org/10.1038/s41467-020-15235-7).
- 750 [67] Zhang et al. *Cycling Data of 64 Cells manufactured by AutoBASS*. Version 2. Nov. 2022. DOI:  
751 [10.5281/zenodo.7299473](https://doi.org/10.5281/zenodo.7299473). URL: <https://doi.org/10.5281/zenodo.7299473>.
- 752 [68] Chunxiang Zhu et al. “Prognosis of Lithium-Ion Batteries’ Remaining Useful Life Based on a  
753 Sequence-to-Sequence Model with Variational Mode Decomposition”. In: *Energies* 16.2 (2023),  
754 p. 803. DOI: [10.3390/en16020803](https://doi.org/10.3390/en16020803).
- 755 [69] Jiangong Zhu et al. *Data-driven capacity estimation of commercial lithium-ion batteries from*  
756 *voltage relaxation*. Apr. 2022. DOI: [10.5281/zenodo.6405084](https://doi.org/10.5281/zenodo.6405084). URL: [https://doi.org/10.](https://doi.org/10.5281/zenodo.6405084)  
757 [5281/zenodo.6405084](https://doi.org/10.5281/zenodo.6405084).
- 758 [70] Jiangong Zhu et al. “Data-driven capacity estimation of commercial lithium-ion batteries from  
759 voltage relaxation”. In: *Nature communications* 13.1 (2022), p. 2261. DOI: [10.1038/s41467-](https://doi.org/10.1038/s41467-022-29837-w)  
760 [022-29837-w](https://doi.org/10.1038/s41467-022-29837-w).