

1

## Supplementary Information

2

### 3 **Inline quality grading of commercial lithium-ion battery manufacturing via** 4 **data-efficient learning and transferable assessment**

5 Chen Liang <sup>1,2,8</sup>, Shengyu Tao <sup>2,3,8,\*</sup>, Chunqiu Xia <sup>3</sup>, Xinghao Huang <sup>2</sup>, Hang Hu <sup>4</sup>, Rui Wang <sup>5</sup>, Daoyi Dong <sup>6</sup>, Ziyang Lyu <sup>1,\*</sup>,  
6 Guangmin Zhou <sup>2,\*</sup>, Huadong Mo <sup>7,\*</sup>

7 <sup>1</sup> School of Mathematics and Statistics, University of New South Wales, Sydney, NSW 2052, Australia

8 <sup>2</sup> Tsinghua Shenzhen International Graduate School, Tsinghua University, Shenzhen, 518055, China

9 <sup>3</sup> Department of Electrical Engineering, Chalmers University of Technology, Gothenburg, 41296, Sweden

10 <sup>4</sup> Department of Mechanical Engineering, Tsinghua University, Beijing, 100084, China

11 <sup>5</sup> Faculty of Mechanical Engineering and Mechanics, Ningbo University, Ningbo, 315211, China

12 <sup>6</sup> Australian Artificial Intelligence Institute, Faculty of Engineering and Information Technology, University of Technology  
13 Sydney, Sydney, NSW 2007, Australia

14 <sup>7</sup> School of Systems and Computing, University of New South Wales, Canberra, ACT 2610, Australia

15 <sup>8</sup> These authors contributed equally to this article

16 \* Corresponding authors: shengyu.tao@chalmers.se (S.T.), ziyang.lyu@unsw.edu.au (Z.L.), huadong.mo@unsw.edu.au (H.M.),  
17 guangminzhou@sz.tsinghua.edu.cn (G.Z.).

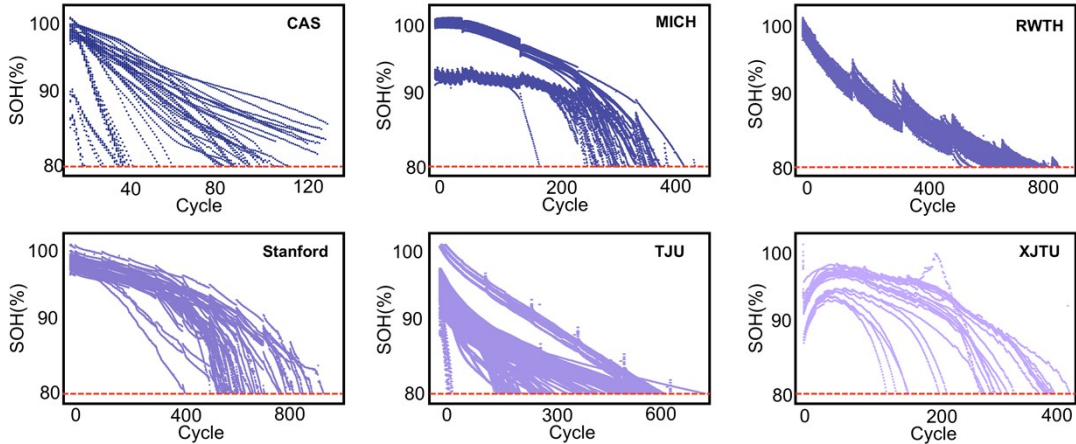
18

20	Supplementary Fig. 1. Degradation trajectories of the batch-based datasets( <i>b-datasets</i> ). .....	4
21	Supplementary Fig. 2. Degradation trajectories of the material-based datasets( <i>material_datasets</i> ). .....	5
22	Supplementary Fig. 3. Distribution of cycle life across the <i>material_dataset</i> . .....	6
23	Supplementary Fig. 4. Bimodal distributions of extracted health indicators in <i>batch_dataset</i> . .....	7
24	Supplementary Fig. 5. Bimodal distributions of extracted health indicators in <i>material_dataset</i> . .....	8
25	Supplementary Fig. 6. The diagnostic plots of the LME model for the <i>material_dataset</i> .....	9
26	Supplementary Fig. 7. Accuracy under various noise conditions.....	10
27	Supplementary Fig. 8. Accuracy comparison with baseline methods .....	11
28	Supplementary Fig. 9. Ablation experiments on label missing mechanisms .....	12
29	Supplementary Fig. 10. Ablation experiments on material features $HI_m$ .....	13
30	Supplementary Table 1. Description of the <i>batch_dataset</i> . .....	14
31	Supplementary Table 2. Description of the <i>material_dataset</i> . .....	15
32	Supplementary Table 3. Description of the base and extended sub-datasets of the <i>batch_dataset</i> . .....	16
33	Supplementary Table 4. Fifteen candidate HIs .....	17
34	Supplementary Table 5. Material grouping of the <i>batch_dataset</i> .....	19
35	Supplementary Table 6. Charging-rate grouping of the <i>batch_dataset</i> .....	19
36	Supplementary Table 7. Temperature grouping of the <i>batch_dataset</i> .....	19
37	Supplementary Table 8. Material grouping of the <i>material_dataset</i> .....	20
38	Supplementary Table 9. Charging-rate grouping of the <i>material_dataset</i> .....	20
39	Supplementary Table 10. Fixed-effects of the LME model for the <i>batch_dataset</i> . .....	21
40	Supplementary Table 11. Random-effects of LME model for the <i>batch_dataset</i> . .....	22
41	Supplementary Table 12. Fixed-effects of the LME model for the <i>material_dataset</i> .....	23
42	Supplementary Table 13. Random effects of the LME model for the <i>material_dataset</i> .....	24
43	Supplementary Table 15. Performance of DELTA for <i>material_dataset</i> .....	26
44	Supplementary Table 16. Robustness of DELTA under noisy conditions for the <i>batch_dataset</i> .....	27
45	Supplementary Table 17. Robustness of DELTA under noisy conditions for the <i>material_dataset</i> .....	29
46	Supplementary Table 18 Comparing with other methods of the <i>batch_dataset</i> .....	31
47	Supplementary Table 19 Comparing with other methods of the <i>material_dataset</i> .....	33
48	Supplementary Table 20. Computational efficiency comparison on the <i>batch_dataset</i> .....	35
49	Supplementary Table 21. Computational efficiency comparison on the <i>material_dataset</i> .....	36

50	Supplementary Table 22. Benchmark settings used for fair comparison of DELTA and baseline methods	
51	.....	37
52	Supplementary Table 23. Definitions of key datasets and classification schemes used in this study. ....	38
53	Supplementary Note 1. Relationship between battery manufacturing workflow and early-cycle quality	
54	assessment.....	39
55	Supplementary Note 2. Analysis of Data Heterogeneity .....	41
56	Supplementary Note 3. Three-class scheme based on 1sigma boundaries of a Gaussian distribution....	42
57	Supplementary Note 4. Economic analysis under end-of-line screening scenario .....	44
58	Supplementary Note 5. Feature extraction.....	48
59	Supplementary Note 6. Linear mixed-effects model .....	50
60	Supplementary Note 7. Semi-supervised Gaussian mixture model.....	51
61	Supplementary Note 8. Evaluation Metrics .....	53
62	Reference .....	54
63		
64		

65 **Supplementary Fig. 1. Degradation trajectories of the batch-based datasets(*b-datasets*).**

66 This figure shows the SOH degradation trajectories of cells from different *b\_datasets*. Each subplot  
67 corresponds to a specific dataset. The dashed red line indicates the EoL threshold at 80% SOH. Despite  
68 being tested under nominally similar operating conditions, the cells exhibit noticeable variability in  
69 degradation patterns and cycle lifetimes, reflecting the influence of batch-to-batch variations and  
70 stochastic ageing factors.



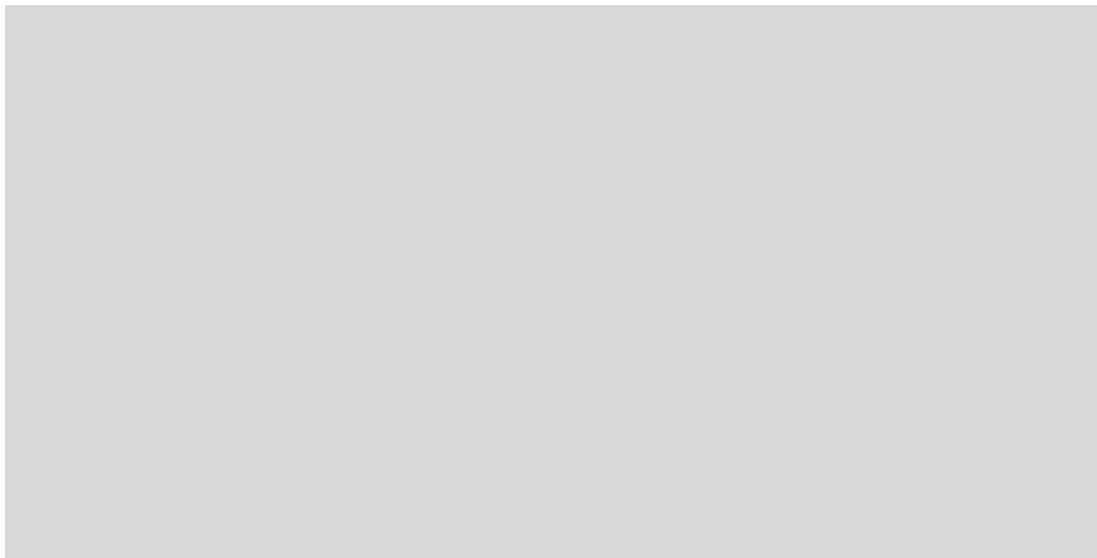
71

72

73 **Supplementary Fig. 2. Degradation trajectories of the material-based datasets(*material\_datasets*).**

74 This figure presents the SOH degradation trajectories of cells grouped according to their material.  
75 Different colors represent datasets with distinct cathode materials and operating conditions. The curves  
76 demonstrate the diversity of degradation behaviors across different materials, including variations in  
77 early-stage capacity fade and long-term ageing rates.

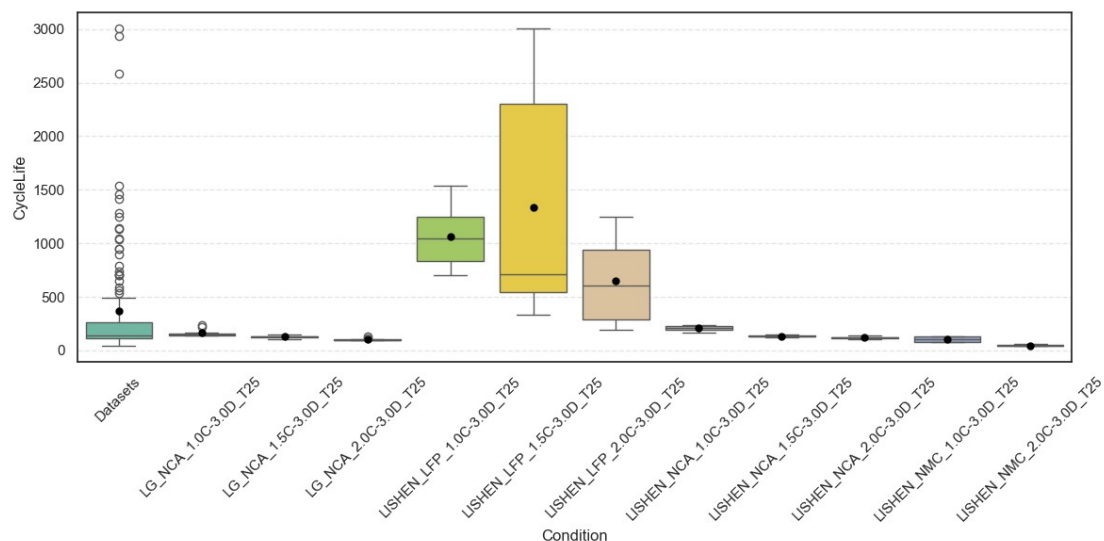
78 The legend indicates the specific dataset corresponding to each ageing trajectory. Each label follows the  
79 format *Manufacturer\_Material\_Chargerate-Dischargerate\_Temperature*, describing the testing  
80 conditions of the cells. For example, LG\_NCA\_1.0C-3.0D\_T25 denotes cells manufactured by LG with  
81 an NCA cathode chemistry, tested under a 1C charge rate and 3C discharge rate at 25 °C. Similarly,  
82 LISHEN\_LFP\_1.5C-3.0D\_T25 represents Lishen cells with an LFP cathode chemistry cycled under a  
83 1.5C charge and 3C discharge rate at 25 °C.



84

85

86 **Supplementary Fig. 3. Distribution of cycle life across the *material\_dataset***



87

88 Supplementary Fig. 3 illustrates the distribution of cycle life for different materials and testing conditions  
89 in the *material\_dataset*. Each boxplot represents the statistical distribution of EOL cycle life under a  
90 specific material-operating-condition combination, including variations in cathode chemistry (e.g., NCA,  
91 LFP and NMC) and charge/discharge rates. The aggregated dataset exhibits a wide range of cycle life  
92 values, spanning from fewer than 100 cycles to nearly 3000 cycles, indicating substantial heterogeneity  
93 across different materials and experimental conditions.

94 Significant differences can be observed among material groups. For example, cells based on LFP  
95 chemistry generally exhibit higher cycle life and larger variance, reflecting their well-known long-cycle  
96 stability but also indicating strong variability under different operating conditions. In contrast, NCA- and  
97 NMC-based cells tend to show comparatively shorter cycle life distributions, with narrower interquartile  
98 ranges in most conditions. These discrepancies highlight the strong influence of material and testing  
99 protocols on battery ageing behavior.

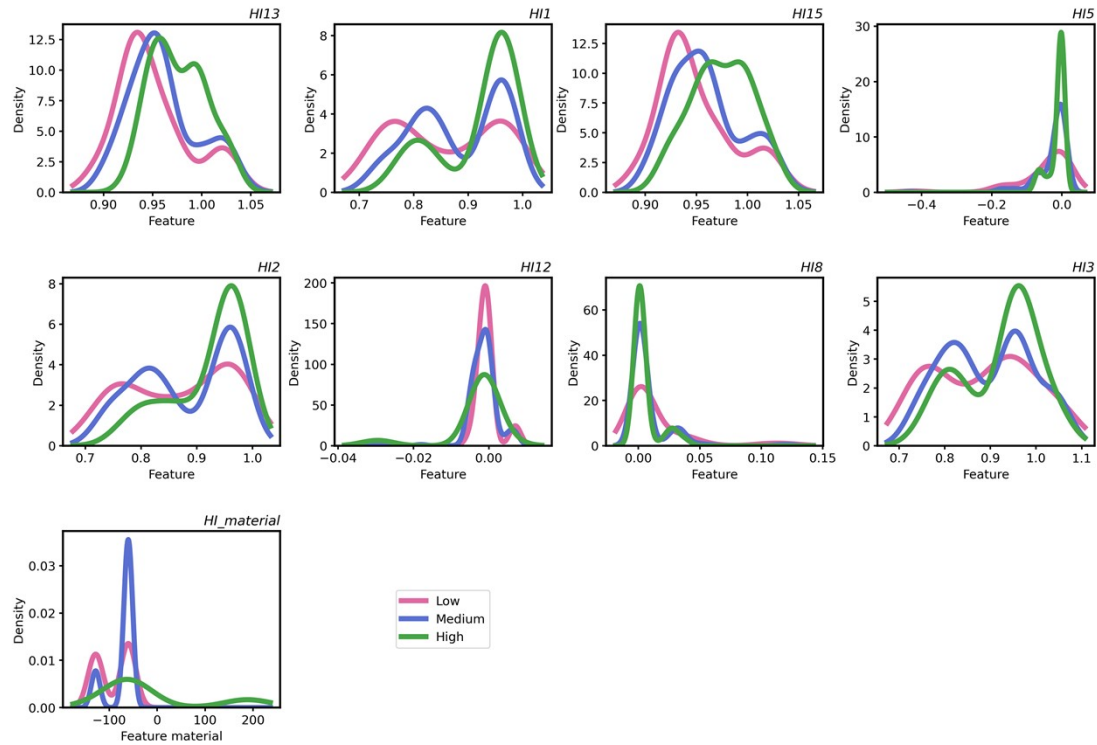
100 Moreover, even within the same material, noticeable dispersion exists across different C-rate  
101 combinations. This observation suggests that battery degradation is affected not only by material  
102 chemistry but also by operational factors such as charging rate, which can significantly alter  
103 electrochemical stress and ageing dynamics.

104 Overall, the pronounced heterogeneity observed in the *material\_dataset* further motivates the use of the  
105 Linear Mixture Model (LME) modelling strategy within the Data-Efficient Learning and Transferable  
106 Assessment (DELTA) framework. By explicitly modelling material-dependent random effects, the  
107 framework is able to capture latent variability across different materials and operating conditions, thereby  
108 improving robustness and generalization when applied to heterogeneous datasets.

109

110 **Supplementary Fig. 4. Bimodal distributions of extracted health indicators in *batch\_dataset*.**

111 This figure shows the health indicators (HIs) of *batch\_dataset* we select from 16 HIs candidates using  
112 the coefficient of association and their Kernel density estimates. The curves correspond to 3 cycle-life  
113 categories (low, normal, and high). Several HIs exhibit bimodal or multimodal distributions.

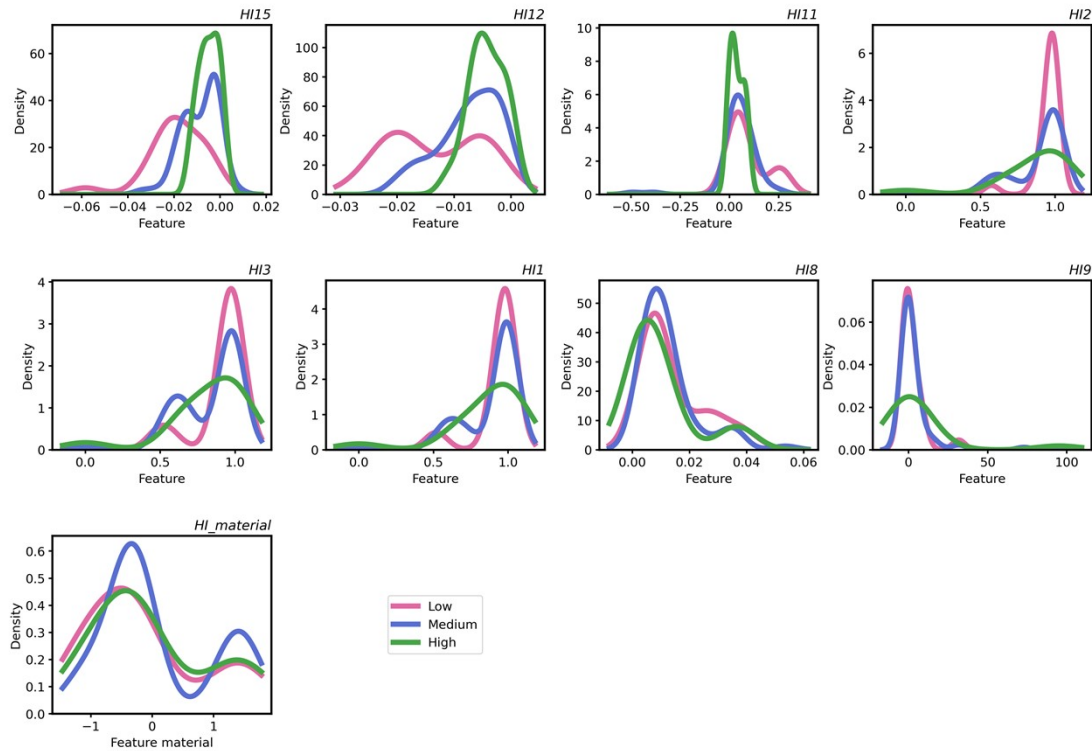


114

115

116 **Supplementary Fig. 5. Bimodal distributions of extracted health indicators in *material\_dataset*.**

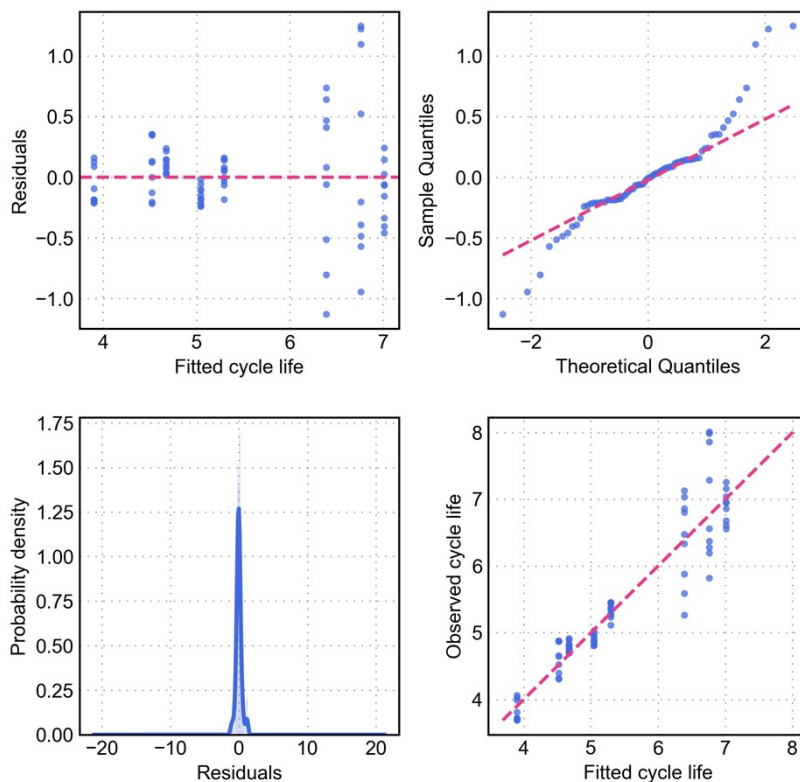
117 This figure shows the health indicators (HIs) of *material\_dataset* we select from 16 HIs candidates using  
118 coefficient of association and their Kernel density estimates. The curves correspond to 3 cycle-life  
119 categories (low, normal, and high). Several HIs exhibit bimodal or multimodal distributions. The  
120 presence of bimodal feature distributions further justifies the use of K-means clustering before Gaussian  
121 mixture model classification, which helps separate underlying feature modes and enhances the robustness  
122 of the subsequent probabilistic classification.



123

124

125 **Supplementary Fig. 6. The diagnostic plots of the LME model for the *material\_dataset***

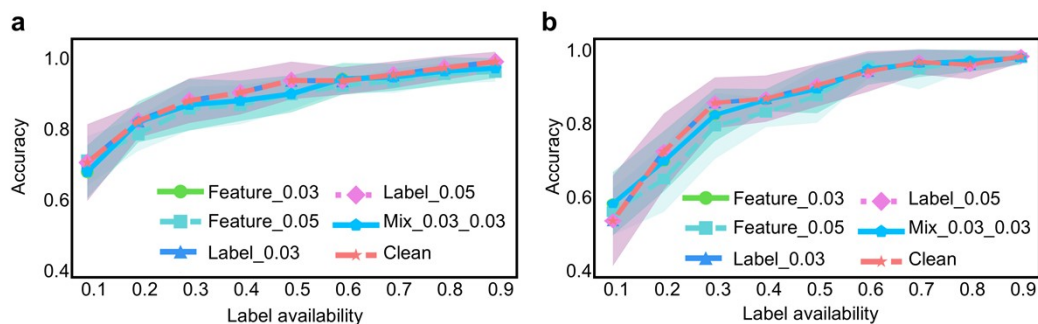


126

127 Residual diagnostic plots are presented to assess the validity of the LME model. Residuals are randomly  
128 distributed around zero without clear patterns, indicating good model fit. The Q-Q plot shows that  
129 residuals approximately follow a normal distribution, with only minor deviations at the tails. The residual  
130 density is centered around zero, suggesting limited systematic bias. In addition, the observed versus fitted  
131 values exhibit a strong linear relationship close to the identity line, confirming that the model captures  
132 the major variation in cycle life. These results support the suitability of the LME model for extracting  
133 material-related features in the DELTA framework.

134

135 **Supplementary Fig. 7. Accuracy under various noise conditions**



136

137 Supplementary Fig. 7 illustrates the classification accuracy of DELTA under different noise conditions  
138 as a function of label availability. Panel (a) corresponds to the *batch\_dataset*, while panel (b) corresponds  
139 to the *material\_dataset*. The evaluated noise settings include feature noise (Gaussian perturbations with  
140 standard deviations of 0.03 and 0.05), label noise (random label flipping with ratios of 0.03 and 0.05), a  
141 mixed noise scenario combining both feature and label perturbations, and a clean baseline without noise.

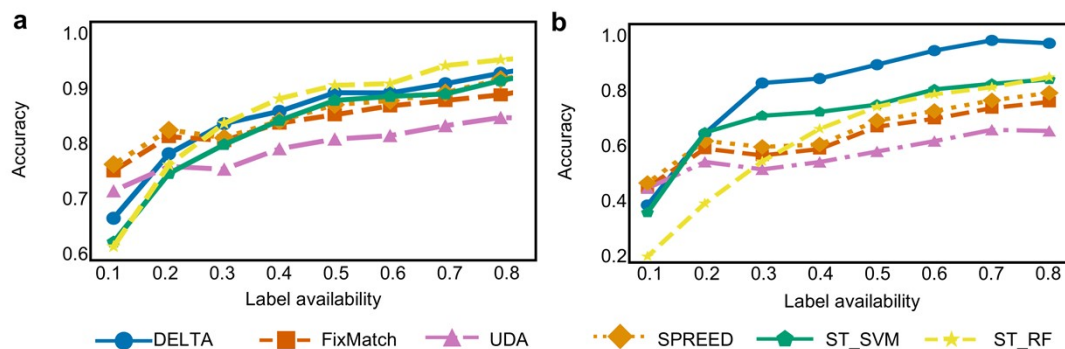
142 Overall, the classification accuracy consistently increases as label availability grows, indicating that the  
143 model effectively benefits from additional labeled information. Even when only a small proportion of  
144 labels are available (20%), the model already achieves moderate accuracy levels, demonstrating the  
145 capability of DELTA to exploit structural information in unlabeled data.

146 The results further show that feature noise has a slightly larger impact on performance than label noise,  
147 particularly at very low label availability levels. Nevertheless, the performance degradation remains  
148 relatively small across all noise configurations. As label availability increases beyond approximately 0.5,  
149 the accuracy curves under different noise settings gradually converge toward the clean baseline,  
150 indicating that the model becomes increasingly robust to noise when more labeled data are available.

151 A similar trend can be observed for the *material\_dataset* in panel (b). Although the initial accuracy values  
152 are slightly lower due to higher dataset variability, the classification performance quickly improves as  
153 label availability increases. Under higher label availability levels (70%-90%), all noise scenarios achieve  
154 accuracy values close to the clean baseline, further demonstrating the robustness of DELTA under  
155 realistic noisy conditions.

156

157 **Supplementary Fig. 8. Accuracy comparison with baseline methods**



158

159 Supplementary Fig. 8 compares the classification accuracy of DELTA with several baseline methods  
160 under different label availability ratios. Panel (a) presents the results on the *batch\_dataset*, while panel  
161 (b) shows the results on the *material\_dataset*. The baseline methods include FixMatch, UDA, SPRED,  
162 ST\_SVM, and ST\_RF.

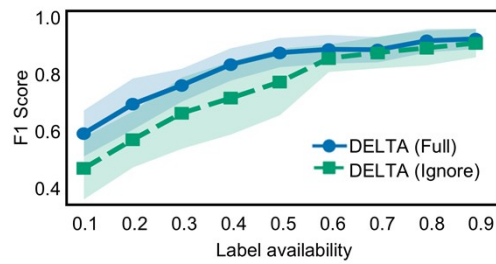
163 On the *batch\_dataset* (panel a), all methods exhibit improved performance as label availability increases,  
164 reflecting the importance of labeled data in semi-supervised classification. Among the baseline methods,  
165 ST\_RF and SPRED achieve relatively strong performance at moderate label availability levels. However,  
166 DELTA shows the fastest performance growth as label availability increases and consistently achieves  
167 higher accuracy once sufficient labeled samples are available.

168 On the *material\_dataset* (panel b), the performance gap between DELTA and the baseline methods  
169 becomes even more pronounced. While several baseline methods exhibit unstable performance or  
170 relatively slow improvement as label availability increases, DELTA achieves significantly higher  
171 accuracy across most label availability levels. In particular, when label availability exceeds  
172 approximately 0.6, DELTA approaches near-perfect classification accuracy, whereas the best baseline  
173 methods remain substantially lower.

174 These results further confirm that DELTA provides superior scalability with respect to label availability  
175 and stronger robustness across different datasets, highlighting the effectiveness of the proposed  
176 framework in handling partially labeled battery datasets.

177

178 **Supplementary Fig. 9. Ablation experiments on label missing mechanisms**



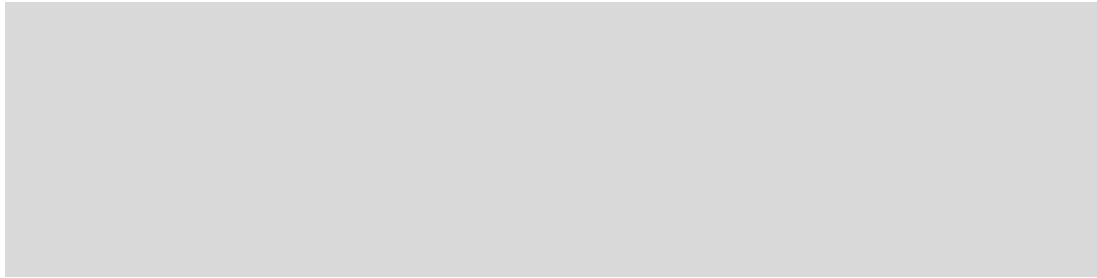
179

180 Supplementary Fig. 9 presents the accuracy comparison between the complete DELTA framework  
181 (“Full”) and the variant that ignores the label-missing mechanism (“Ignore”) under different levels of  
182 label availability. As the proportion of available labels increases from 0.1 to 0.9, the prediction accuracy  
183 of both models improves steadily, indicating that additional labeled data consistently enhances model  
184 performance.

185 Across the entire range of label availability, the full DELTA framework consistently achieves higher  
186 accuracy than the variant that ignores the label-missing mechanism. The performance gap is particularly  
187 pronounced under extremely limited labeling conditions (for example, label availability below 0.3),  
188 where the full model maintains noticeably higher accuracy and exhibits reduced performance variability.  
189 These results demonstrate that explicitly modelling the label-missing mechanism enables more effective  
190 utilization of partially labeled data, thereby improving robustness and prediction accuracy under label-  
191 scarce conditions.

192

193 **Supplementary Fig. 10. Ablation experiments on material features  $HI_m$**



194

195 Supplementary Fig. 10 shows the accuracy comparison between the DELTA framework incorporating  
196 the LME-derived material feature  $HI_m$  and a variant in which this feature is removed. As label availability  
197 increases, the prediction accuracy of both models improves gradually, reflecting the positive impact of  
198 increased supervision. However, the model incorporating  $HI_m$  consistently outperforms the variant  
199 without this feature across all levels of label availability. The improvement becomes more evident in the  
200 low-label regime, where the inclusion of  $HI_m$  leads to higher accuracy and more stable performance.  
201 This observation indicates that the random-effects modelling provided by the LME effectively captures  
202 latent inter-cell variability and material-induced heterogeneity in early-cycle data. By incorporating this  
203 uncertainty-aware feature, the DELTA framework achieves improved classification reliability and  
204 robustness.

205

206 **Supplementary Table 1. Description of the *batch\_dataset*.**

207 This table summarizes the detailed descriptions of the *batch\_dataset* used in this study. For each dataset,  
 208 the corresponding cathode material, chemical formula, rated capacity, cut-off voltage range, operating  
 209 temperature, charge/discharge rate, and the number of cells are provided. These datasets originate from  
 210 multiple public sources, including CAS<sup>1</sup>, MICH<sup>2</sup>, RWTH<sup>3</sup>, Stanford<sup>4</sup>, TJU<sup>5</sup>, and XJTU<sup>6</sup>, covering  
 211 different material systems and diverse operating conditions, most of which can be processed by  
 212 BatteryLife<sup>7</sup>.

Dataset s	Material	Chemical Formula	Qn (Ah)	Cut-off Voltage (V)	Temp eratur e (°C)	Charge/d ischarge rate (C)	N
CAS	NCA	-(LG)	2.35	2.65-4.2	25	1/3	10
		-(LISHEN)			25	1/3	10
	NCM	-(LISHEN)	2.6	2.75-4.2	0	0.5/3	10
		-(LISHEN)			25	1/3	9
		-(LISHEN)			25	2/3	10
MICH	NCM	$LiNi_{1/3}Co_{1/3}Mn_{1/3}O_2$	2.36	3.0-4.2	45	1/1	10
					25		10
					45		10
RWTH	NCM	-	1.11	3.5-3.9	25	2/2	48
Stanfor d	NCM	$LiNi_{0.5}Mn_{0.3}Co_{0.2}O_2$	0.24	3-4.4	30	1/0.75	41
					25	0.25/1	7
	NCA	$Li_{0.86}Ni_{0.86}Co_{0.11}Al_{0.03}O_2$	3.5	2.65-4.2	25	0.5/1	19
					35	0.5/1	3
					45	0.5/1	28
TJU	NCM	$Li_{0.86}Ni_{0.86}Co_{0.11}Mn_{0.07}O_2$	3.5	2.5-4.2	25	0.5/1	23
					35	0.5/1	4
					45	0.5/1	28
NCA+NC M	-	-	-	-	25	0.5/1	3
					25	0.5/2	3
XJTU	NCM	$LiNi_{0.5}Co_{0.2}Mn_{0.3}O_2$	2.0	2.5-4.2	20	2/1	8
						3/1	15

213

214 **Supplementary Table 2. Description of the *material\_dataset*.**

215 This table summarizes the detailed descriptions of the *material\_dataset* used in this study. Cells sharing  
 216 the same material are grouped together and assigned the corresponding *m\_label*. The *material\_dataset*  
 217 was constructed from the CAS dataset because it contains multiple battery chemistries, enabling material-  
 218 level grouping and cross-material comparisons. For each material group, the material type, chemical  
 219 formula, rated capacity, cut-off voltage range, operating temperature, charge/discharge rate, and number  
 220 of cells are provided.

Datasets	Material	Chemical Formula	Qn (Ah)	Cut-off Voltage (V)	Temperature (°C)	Charge/discharge rate (C)	N
CAS	NCA	-(LG)	3.35	2.65-4.2	25	1C/3C	10
						1.5C/3C	10
						2C/3C	10
						1C/3C	10
	NCA	-(LISHEN)	3.35	2.65-4.2		1.5C/3C	10
						2C/3C	10
	NCM	-(LISHEN)	2.6	2.75-4.2		1C/3C	9
						2C/3C	10
	LFP	-	1.5	2.0-4.0		1C/3C	15
						1.5C/3C	30
					2C/3C	15	

221

222 **Supplementary Table 3. Description of the base and extended sub-datasets of the *batch\_dataset*.**

223 This table summarizes the detailed descriptions of the two sub-datasets constructed from the  
 224 *batch\_dataset*, including the base dataset and the extended dataset. The *batch\_dataset* was assembled  
 225 from all eligible publicly available datasets and subsequently divided into the base and extended datasets.  
 226 For each dataset, the material system, chemical formula, rated capacity, cut-off voltage range, operating  
 227 temperature, charge/discharge rate, and the number of cells are provided.

Type	Datasets	Material	Chemical Formula	Qn (Ah)	Cut-off Voltage (V)	Temperature (°C)	Charge/discharge rate (C)	Number			
Base dataset	MICH	NCM	$LiNi_{1/3}Co_{1/3}Mn_{1/3}O_2$	2.36	3.0-4.2	25		10			
						45	1/1	10			
						25		10			
						45		10			
	RWTH	NCM	-	-	1.11	3.5-3.9	25	2/2	48		
							Stanford	NCM	$LiNi_{0.5}Mn_{0.3}Co_{0.2}O_2$	0.24	3-4.4
			25	0.25/1	7						
							25	0.5/1	19		
			NCA	$Li_{0.86}Ni_{0.86}Co_{0.11}Al_{0.03}O_2$	3.5	2.65-4.2	25	1/1	9		
							35	0.5/1	3		
							45	0.5/1	28		
		TJU					25	0.5/1	23		
			NCM	$Li_{0.86}Ni_{0.86}Co_{0.11}Mn_{0.07}O_2$	3.5	2.5-4.2	35	0.5/1	4		
							45	0.5/1	28		
		NCA+				25	0.5/1	3			
		NCM	-	-	-	25	0.5/2	3			
						25	0.5/4	3			
Extended dataset	XJTU	NCM	$LiNi_{0.5}Co_{0.2}Mn_{0.3}O_2$	2.0	2.5-4.2	20	2/1	8			
							3/1	15			
	CAS	NCA	-(LG) -(LISHEN)	-	2.35	2.65-4.2	25	1/3	10		
								2.65-4.2	25	1/3	10
									0	0.5/3	10
									25	1/3	9
		NCM	-(LISHEN) -(LISHEN)	2.6	2.75-4.2	25	1/3	9			
						25	2/3	10			

228

229 **Supplementary Table 4. Fifteen candidate HIs**

Feature	Definition	Formula	Source cycles	Physical interpretation
HI1	Early capacity ratio	$HI_1 = \frac{Q_2}{Q_{nom}}$	Cycle 2	Early-stage capacity health
HI2	Mid-cycle capacity ratio	$HI_2 = \frac{Q_{\frac{2+n}{2}}}{Q_{nom}}$	Cycle (2+n)/2	Intermediate degradation level
HI3	Late early-cycle capacity ratio	$HI_3 = \frac{Q_n}{Q_{nom}}$	Cycle n	Early degradation magnitude
HI4	Maximum $\Delta Q(V)$	$HI_4 = \frac{\max(Q_{2(V)} - Q_{n(V)})}{Q_{nom}}$	Cycle 2 vs n	Maximum capacity difference along voltage curve
HI5	Minimum $\Delta Q(V)$	$HI_5 = \frac{\min(Q_{2(V)} - Q_{n(V)})}{Q_{nom}}$	Cycle 2 vs n	Minimum capacity difference
HI6	Mean $\Delta Q(V)$	$HI_6 = \frac{\text{mean}(Q_{2(V)} - Q_{n(V)})}{Q_{nom}}$	Cycle 2 vs n	Average degradation across voltage
HI7	Variance of $\Delta Q(V)$	$HI_7 = \frac{\text{Var}(Q_{2(V)} - Q_{n(V)})}{Q_{nom}^2}$	Cycle 2 vs n	Variability of degradation
HI8	Standard deviation of $\Delta Q(V)$	$HI_8 = \frac{\text{Std}(Q_{2(V)} - Q_{n(V)})}{Q_{nom}}$	Cycle 2 vs n	Dispersion of capacity difference
HI9	Kurtosis of $\Delta Q(V)$	$HI_9 = \text{Kurtosis}(Q_{2(V)} - Q_{n(V)})$	Cycle 2 vs n	Shape changes of Q-V curve
HI10	Skewness of $\Delta Q(V)$	$HI_{10} = \text{Skew}(Q_{2(V)} - Q_{n(V)})$	Cycle 2 vs n	Asymmetry of degradation
HI11	Voltage slope of $\Delta Q(V)$	$HI_{11} = \frac{d(Q_{2(V)} - Q_{n(V)})}{dV}$	Cycle 2 vs n	Polarization evolution
HI12	Capacity fade slope	$HI_{12} = \frac{dQ}{dN}$	All cycles	Rate of capacity degradation
HI13	Capacity intercept	$Q = aN + b$	All cycles	Estimated initial capacity
HI14	Regression goodness-of-fit	$HI_{14} = R^2$	All cycles	Linearity of degradation trend
HI15	Average capacity fade per cycle	$HI_{15} = \text{mean}(\Delta Q)$	All cycles	Mean degradation per cycle

230 This table summarizes the definitions of 15 candidate HIs extracted from early-cycle battery data. The  
 231 indicators are derived from three main aspects of degradation behavior: capacity-based metrics,  
 232 differential capacity–voltage curve features, and capacity fade trend characteristics. The table lists the  
 233 definition, mathematical formulation, source cycles used for calculation, and the corresponding physical  
 234 interpretation of each indicator.

235 Specifically, HI1–HI3 describe early-stage capacity retention across different cycle stages, reflecting the  
236 initial health condition and early degradation level of the cells. HI4–HI11 are extracted from the  
237 differences between capacity-voltage (Q-V) curves at early and later cycles, capturing structural changes  
238 in electrochemical behaviour, such as capacity loss distribution along the voltage range, polarization  
239 evolution, and statistical characteristics of degradation patterns. HI12–HI15 characterize the overall  
240 capacity fade dynamics across the cycling process through regression-based metrics and statistical  
241 measures, including degradation rate, linearity of the degradation trend, and average capacity loss per  
242 cycle.

243 Together, these candidate indicators provide complementary information about battery ageing behaviour  
244 and form the basis for subsequent feature selection and health indicator construction.

246 **Supplementary Table 5. Material grouping of the *batch\_dataset***

Material	Number	Group	Num_group
NCM	206	NCA	206
NCA	56	NCM	56
NCANCM	9	NCANCM	9

247

248 **Supplementary Table 6. Charging-rate grouping of the *batch\_dataset***

Charging rate	Number	Group	Num_group
0.25	7	verylow	7
0.5	112	low	112
1	81	medium	81
2	56	high	56
3	15	veryhigh	15

249

250 **Supplementary Table 7. Temperature grouping of the *batch\_dataset***

Temperature	Number	Group	Num_group
20	23	low	23
25	145	midium	145
30	41	high	47
35	6		
45	56	Veryhigh	56

251

252 Material, temperature, and charge rate are widely recognized as the three dominant factors governing  
 253 battery ageing behavior, among which material exerts the most pronounced influence. Accordingly, in  
 254 the LME model, material is treated as a random effect to capture material-dependent variability, while  
 255 temperature and charge rate are considered fixed effects to quantify their systematic influence on  
 256 degradation behavior.

257 Because these factors are discrete experimental variables, they were grouped into representative levels  
 258 to facilitate statistical modelling. The grouping definitions for the *batch\_dataset* are summarized in  
 259 Supplementary Tables 5-7, corresponding to material type, charging rate, and temperature, respectively.

260 For the *material\_dataset*, only cells cycled under room-temperature conditions were included. Therefore,  
 261 temperature variation is not considered in this dataset. Consequently, the grouping scheme focuses on  
 262 material type and charging rate, as summarized in Supplementary Tables 8-9.

263

264 **Supplementary Table 8. Material grouping of the *material\_dataset***

Material	Number	Group	Num_group
NCM	19	NCM	19
NCA	60	NCA	60
LFP	60	LFP	60

265

266 **Supplementary Table 9. Charging-rate grouping of the *material\_dataset***

Charging rate	Number	Group	Num_group
1	44	low	44
1.5	50	medium	50
2	35	high	35

267

268 **Supplementary Table 10. Fixed-effects of the LME model for the *batch\_dataset*.**

269 This table summarizes the estimated fixed effects of the LME model for the *batch\_dataset*. In this model,  
 270 temperature and charging rate are treated as fixed effects, while the material system is modelled as a  
 271 random effect. The intercept represents the expected EOL cycle life under the reference operating  
 272 conditions, which correspond to the high temperature (highT) and high charging rate (highC) categories.  
 273 The results show that both temperature and charging rate have statistically significant effects on cycle  
 274 life ( $p < 0.001$  for all estimated coefficients). Among the temperature-related effects, both lowT and  
 275 midT show large negative coefficients relative to the reference category, indicating that these temperature  
 276 groups are associated with shorter predicted cycle life under the adopted grouping scheme. Similarly,  
 277 several charging-rate levels exhibit significant deviations from the reference condition. In particular, the  
 278 lowC and verylowC groups show substantial negative coefficients, suggesting pronounced differences  
 279 in degradation behaviour across charging-rate regimes. Overall, these results confirm that operating  
 280 conditions strongly influence battery ageing, and that temperature and charging rate introduce systematic  
 281 variations in cycle life across the dataset.

	Estimate	Std. Error	t-value	p-value
(Intercept)	1054.364551	82.13205133	12.83743112	***
lowT	-593.741737	32.1670511	-18.45807174	***
midT	-336.2625703	14.81582101	-22.69618201	***
veryhighT	-160.9742509	19.66361168	-8.186403065	***
highT	0	0	0	
lowC	-451.3203597	16.80768646	-26.85202159	***
mediumC	-307.4562084	15.5470799	-19.77581709	***
veryhighC	-152.9083333	32.73254768	-4.671446134	***
verylowC	-363.6384853	33.64721292	-10.80738801	***
highC	0	0	0	

282 Note: highT and highC are the reference categories. Standard errors are reported in the second column.

283 \*\*\*  $p < 0.001$ , \*\*  $p < 0.01$ , \*  $p < 0.05$ .

284

285 **Supplementary Table 11. Random-effects of LME model for the *batch\_dataset*.**

286 This table presents the estimated random effects associated with different material systems in the LME  
287 model for the *batch\_dataset*. The reported values correspond to the estimated best linear unbiased  
288 predictors (EBLUPs) of the random-effect terms. These coefficients represent the material-dependent  
289 deviations from the global mean cycle life after accounting for the fixed effects of temperature and  
290 charging rate. The results show that the NCANCM material system exhibits a large positive deviation,  
291 indicating that its baseline cycle life is substantially higher than the dataset average under comparable  
292 operating conditions. In contrast, both NCA and NCM show negative deviations, implying relatively  
293 shorter baseline cycle lifetimes compared with the global mean. These differences highlight the  
294 systematic influence of cathode material chemistry on battery degradation behaviour. The estimated  
295 EBLUPs are subsequently used as a static feature (*HI\_m*) to capture the material-dependent ageing  
296 contribution in the downstream classification framework.

Material	EBLUPs
NCA	-128.1777815
NCANCM	188.425596
NCM	-60.24781448

297

298

299 **Supplementary Table 12. Fixed-effects of the LME model for the *material\_dataset***

300 This table reports the fixed-effect estimates of the LME model for the *material\_dataset*. Because this  
 301 dataset only contains cells tested under room-temperature conditions, temperature is not included as a  
 302 factor in the model. Consequently, the fixed effects only account for charging-rate variations. The  
 303 intercept represents the expected cycle life under the reference charging-rate category (highC). The  
 304 coefficients for lowC and mediumC are positive and statistically significant, indicating that these  
 305 charging-rate levels are associated with longer predicted cycle life compared with the high charging-rate  
 306 condition. This result is consistent with well-established electrochemical ageing mechanisms, as higher  
 307 charging rates typically accelerate degradation processes such as polarization increase, lithium plating,  
 308 and structural stress in electrode materials.

	Estimate	Std. Error	t-value	p-value
(Intercept)	4.990191199	0.607357333	8.216236026	***
lowC	0.621381485	0.105885844	5.868409408	***
mediumC	0.370743396	0.124069312	2.988195801	**
highC	0	0	0	

309 Note: highC is the reference category. \*\*\* p < 0.001, \*\* p < 0.01, \* p < 0.05.

310

311 **Supplementary Table 13. Random effects of the LME model for the *material\_dataset***

312 This table summarizes the random-effect estimates (EBLUPs) associated with different material systems  
313 in the *material\_dataset*. These coefficients quantify the material-dependent deviations in cycle life after  
314 controlling for charging-rate effects. The results indicate that the LFP material system exhibits a positive  
315 deviation, suggesting a longer baseline cycle life relative to the dataset average. In contrast, NCM shows  
316 a negative deviation, indicating relatively shorter lifetimes under similar charging conditions. The NCA  
317 system exhibits a smaller negative deviation compared with NCM. These results further confirm that  
318 material chemistry introduces intrinsic variability in battery ageing behaviour, even when operating  
319 conditions are controlled. Incorporating these material-dependent effects as static features enables the  
320 subsequent modelling framework to better capture systematic ageing differences across material systems.

Material	EBLUPs
LFP	1.402677323
NCA	-0.314301511
NCM	-1.088375812

321

323 Supplementary Table 14. Performance of DELTA for *batch\_dataset*.

324 This table summarizes the quantitative results corresponding to **Fig. 3a**, reporting the classification  
325 accuracy and F1 score of the DELTA framework under different label availability ratios for the  
326 *batch\_dataset*.

327 As the proportion of labeled samples increases from 0.1 to 1.0, both accuracy and F1 score exhibit a clear  
328 upward trend. When only 10% of the labels are available, the classification accuracy is 0.69, indicating  
329 that the model still maintains a reasonable predictive capability even under extremely limited  
330 supervision. As the label availability increases to 0.3-0.4, the accuracy improves substantially to  
331 approximately 0.83-0.85, reflecting the increasing effectiveness of the semi-supervised learning  
332 mechanism.

333 When the label availability reaches 0.5, the accuracy rises to 0.87, demonstrating that the model can  
334 achieve relatively reliable classification performance once half of the samples are labeled. Beyond this  
335 point, the improvement becomes more gradual, with accuracy stabilizing between 0.87 and 0.92 for label  
336 availability levels between 0.6 and 1.0. The close agreement between accuracy and F1 score across all  
337 label visibility levels further indicates balanced classification performance across different classes  
338 without significant class imbalance effects.

339 Overall, these results demonstrate that DELTA exhibits strong robustness to limited label availability,  
340 maintaining stable and high performance once a moderate proportion of labeled data is available.

Label availability	Accuracy	F1 score
1.0	0.923898	0.923461
0.9	0.913931	0.913734
0.8	0.900311	0.900092
0.7	0.884889	0.885137
0.6	0.871334	0.871287
0.5	0.872260	0.872289
0.4	0.845151	0.843830
0.3	0.827076	0.826221
0.2	0.781662	0.782647
0.1	0.689156	0.687769

341

342 **Supplementary Table 15. Performance of DELTA for *material\_dataset*.**

343 reports the classification performance of DELTA under varying label availability ratios for the  
344 *material\_dataset*, corresponding to **Fig. 3b**.

345 A similar monotonic improvement trend can be observed as the label availability increases. When only  
346 10% of labels are available, the classification accuracy is 0.59, which is slightly lower than that of the  
347 *batch\_dataset* but still demonstrates meaningful predictive capability under scarce supervision. As the  
348 label availability increases to 0.2 and 0.3, the accuracy improves rapidly to 0.74 and 0.83, respectively,  
349 indicating that additional labeled samples significantly enhance the model’s ability to distinguish  
350 between cycle-life classes.

351 Once label availability reaches 0.5, the classification accuracy already exceeds 0.92 and continues to  
352 improve steadily as label visibility increases. Under near fully labeled conditions (0.9-1.0), the accuracy  
353 approaches 0.99, suggesting that the DELTA framework is capable of achieving near-perfect  
354 classification performance when sufficient labeled data are available.

355 The F1scores closely follow the same trend as the accuracy values, confirming consistent predictive  
356 performance across all classes. Compared with the *batch\_dataset*, the *material\_dataset* exhibits slightly  
357 higher classification accuracy across most label availability levels, which may be attributed to reduced  
358 variability in experimental conditions when grouped by material systems.

359 These results further confirm the scalability and robustness of the DELTA framework, demonstrating  
360 that its classification performance improves consistently with increasing label availability while  
361 remaining effective even in low-label scenarios.

Label availability	accuracy	F1 score
1.0	0.985582	0.985559
0.9	0.989153	0.989131
0.8	0.974339	0.974293
0.7	0.963624	0.962929
0.6	0.956349	0.956151
0.5	0.923942	0.924655
0.4	0.888228	0.88662
0.3	0.833069	0.832253
0.2	0.739683	0.73894
0.1	0.586376	0.57583

362

363 **Supplementary Table 16. Robustness of DELTA under noisy conditions for the *batch\_dataset***

364 Supplementary Table 16-1. Accuracy

Noise\Label availability	0.9	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
clean	0.914 ±	0.900 ±	0.885 ±	0.871 ±	0.872 ±	0.845 ±	0.827 ±	0.782 ±	0.689 ±
	0.022	0.026	0.030	0.031	0.042	0.049	0.050	0.037	0.086
feat_0.03	0.899 ±	0.893 ±	0.879 ±	0.876 ±	0.842 ±	0.827 ±	0.818 ±	0.780 ±	0.669 ±
	0.022	0.031	0.030	0.035	0.041	0.045	0.058	0.046	0.059
feat_0.05	0.892 ±	0.888 ±	0.871 ±	0.861 ±	0.842 ±	0.818 ±	0.810 ±	0.753 ±	0.693 ±
	0.023	0.023	0.027	0.013	0.036	0.043	0.049	0.037	0.055
label_0.03	0.914 ±	0.900 ±	0.885 ±	0.871 ±	0.872 ±	0.845 ±	0.827 ±	0.782 ±	0.689 ±
	0.022	0.026	0.030	0.031	0.042	0.049	0.050	0.037	0.086
label_0.05	0.914 ±	0.900 ±	0.885 ±	0.871 ±	0.872 ±	0.845 ±	0.827 ±	0.782 ±	0.689 ±
	0.022	0.026	0.030	0.031	0.042	0.049	0.050	0.037	0.086
mix_0.03_0.03	0.899 ±	0.893 ±	0.879 ±	0.876 ±	0.842 ±	0.827 ±	0.818 ±	0.780 ±	0.669 ±
	0.022	0.031	0.030	0.035	0.041	0.045	0.058	0.046	0.059

365

366 Supplementary Table 16-2. F1 Score

Noise\Label availability	0.9	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
clean	0.914 ±	0.900 ±	0.885 ±	0.871 ±	0.872 ±	0.844 ±	0.826 ±	0.783 ±	0.688 ±
	0.022	0.026	0.029	0.031	0.042	0.050	0.053	0.038	0.088
feat_0.03	0.899 ±	0.893 ±	0.879 ±	0.875 ±	0.840 ±	0.825 ±	0.816 ±	0.781 ±	0.669 ±
	0.023	0.031	0.030	0.036	0.043	0.047	0.060	0.047	0.058
feat_0.05	0.891 ±	0.887 ±	0.871 ±	0.861 ±	0.841 ±	0.818 ±	0.808 ±	0.754 ±	0.693 ±
	0.023	0.023	0.028	0.013	0.037	0.045	0.051	0.036	0.055
label_0.03	0.914 ±	0.900 ±	0.885 ±	0.871 ±	0.872 ±	0.844 ±	0.826 ±	0.783 ±	0.688 ±
	0.022	0.026	0.029	0.031	0.042	0.050	0.053	0.038	0.088
label_0.05	0.914 ±	0.900 ±	0.885 ±	0.871 ±	0.872 ±	0.844 ±	0.826 ±	0.783 ±	0.688 ±
	0.022	0.026	0.029	0.031	0.042	0.050	0.053	0.038	0.088
mix_0.03_0.03	0.899 ±	0.893 ±	0.879 ±	0.875 ±	0.840 ±	0.825 ±	0.816 ±	0.781 ±	0.669 ±
	0.023	0.031	0.030	0.036	0.043	0.047	0.060	0.047	0.058

367

368 This table summarizes the classification performance of the DELTA framework on the *batch\_dataset*  
369 under different noise settings and label availability levels. Both accuracy and F1 score are reported as  
370 mean ± standard deviation over repeated experiments. The table includes several scenarios: a clean  
371 baseline without noise, feature noise with Gaussian perturbations (standard deviation = 0.03 or 0.05),  
372 label noise where a proportion of training labels is randomly flipped (0.03 or 0.05), and a mixed setting  
373 combining feature and label noise.

374 Under the clean setting, the model achieves high performance across all label availability levels. When  
375 the label availability is 0.1, the accuracy already reaches 0.914, indicating strong predictive capability  
376 even under extremely limited supervision. As the label availability decreases toward 0.9 (i.e., fewer  
377 labeled samples), the performance gradually declines, reaching 0.689, reflecting the increasing difficulty  
378 of classification when labeled information becomes scarce.

379 When feature noise is introduced, the classification performance shows only a moderate decrease  
380 compared with the clean baseline. For instance, under feat\_0.03, the accuracy at label availability 0.1  
381 remains close to 0.90, demonstrating that the model is relatively insensitive to small perturbations in  
382 input features. Even with stronger feature noise (feat\_0.05), the performance degradation remains  
383 limited, indicating that the feature extraction and clustering steps effectively mitigate noise influence.

384 In contrast, label noise produces almost identical results to the clean setting in this experiment. This  
385 behaviour suggests that the semi-supervised structure of DELTA can effectively tolerate small levels of  
386 label corruption by leveraging the structure of the unlabeled data.

387 The mixed-noise scenario, which combines feature and label perturbations, yields results similar to those  
388 of the feature-noise-only case. This indicates that the model's performance degradation is primarily  
389 driven by feature perturbations rather than label noise. Overall, the results confirm that the DELTA  
390 framework maintains stable classification accuracy and F1 score across a range of noisy conditions,  
391 demonstrating strong robustness to both measurement noise and annotation errors.

392

393 **Supplementary Table 17. Robustness of DELTA under noisy conditions for the *material\_dataset***

394 Supplementary Table 17-1. Accuracy

Noise\Label availability	0.9	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
clean	0.989 ±	0.967 ±	0.975 ±	0.949 ±	0.913 ±	0.877 ±	0.866 ±	0.737 ±	0.552 ±
	0.017	0.038	0.033	0.054	0.057	0.062	0.068	0.101	0.120
feat_0.03	0.986 ±	0.978 ±	0.971 ±	0.956 ±	0.902 ±	0.873 ±	0.834 ±	0.710 ±	0.597 ±
	0.018	0.029	0.036	0.042	0.064	0.030	0.070	0.080	0.084
feat_0.05	0.986 ±	0.978 ±	0.956 ±	0.964 ±	0.884 ±	0.841 ±	0.804 ±	0.663 ±	0.576 ±
	0.018	0.024	0.054	0.033	0.072	0.040	0.087	0.087	0.062
label_0.03	0.989 ±	0.967 ±	0.975 ±	0.949 ±	0.913 ±	0.877 ±	0.866 ±	0.737 ±	0.552 ±
	0.017	0.038	0.033	0.054	0.057	0.062	0.068	0.101	0.120
label_0.05	0.989 ±	0.967 ±	0.975 ±	0.949 ±	0.913 ±	0.877 ±	0.866 ±	0.737 ±	0.552 ±
	0.017	0.038	0.033	0.054	0.057	0.062	0.068	0.101	0.120
mix_0.03_0.03	0.986 ±	0.978 ±	0.971 ±	0.956 ±	0.902 ±	0.873 ±	0.834 ±	0.710 ±	0.597 ±
	0.018	0.029	0.036	0.042	0.064	0.030	0.070	0.080	0.084

395

396 Supplementary Table 17-2. F1 Score

Noise\Label availability	0.9	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
clean	0.989 ±	0.967 ±	0.975 ±	0.949 ±	0.913 ±	0.876 ±	0.864 ±	0.732 ±	0.534 ±
	0.017	0.038	0.033	0.054	0.056	0.061	0.069	0.104	0.124
feat_0.03	0.986 ±	0.978 ±	0.971 ±	0.956 ±	0.903 ±	0.873 ±	0.830 ±	0.709 ±	0.595 ±
	0.018	0.029	0.037	0.042	0.062	0.030	0.072	0.082	0.087
feat_0.05	0.986 ±	0.978 ±	0.957 ±	0.964 ±	0.885 ±	0.841 ±	0.800 ±	0.659 ±	0.561 ±
	0.018	0.024	0.054	0.033	0.072	0.039	0.091	0.093	0.074
label_0.03	0.989 ±	0.967 ±	0.975 ±	0.949 ±	0.913 ±	0.876 ±	0.864 ±	0.732 ±	0.534 ±
	0.017	0.038	0.033	0.054	0.056	0.061	0.069	0.104	0.124
label_0.05	0.989 ±	0.967 ±	0.975 ±	0.949 ±	0.913 ±	0.876 ±	0.864 ±	0.732 ±	0.534 ±
	0.017	0.038	0.033	0.054	0.056	0.061	0.069	0.104	0.124
mix_0.03_0.03	0.986 ±	0.978 ±	0.971 ±	0.956 ±	0.903 ±	0.873 ±	0.830 ±	0.709 ±	0.595 ±
	0.018	0.029	0.037	0.042	0.062	0.030	0.072	0.082	0.087

397

398 This table reports the robustness evaluation of DELTA for the *material\_dataset* under the same set of  
 399 noise conditions. Similar to Table 16, the results include classification accuracy and F1 score across  
 400 different label availability ratios and noise configurations.

401 Under the clean setting, the model achieves very high classification accuracy, reaching approximately  
 402 0.99 when label availability is 0.1, and gradually decreasing as the proportion of labeled samples becomes

403 smaller. This indicates that the DELTA framework can effectively exploit both labeled and unlabeled  
404 information to maintain high classification performance.

405 When feature noise is introduced, the classification accuracy decreases slightly but remains relatively  
406 high across most label availability levels. For example, under `feat_0.03`, the accuracy at label availability  
407 0.1 remains around 0.986, and even with stronger perturbations (`feat_0.05`), the performance reduction  
408 remains moderate. These results suggest that the framework is capable of handling moderate feature  
409 perturbations without significant loss of predictive power.

410 Label noise scenarios again yield results almost identical to the clean baseline, indicating that small levels  
411 of label corruption have a limited impact on overall classification performance. This further demonstrates  
412 the robustness of the semi-supervised learning strategy used in DELTA.

413 Under the mixed-noise setting, where both feature and label perturbations are introduced simultaneously,  
414 the performance trends remain consistent with those in the feature-noise scenarios. Although accuracy  
415 and F1 score decrease slightly at higher noise levels or lower label availability, the model maintains a  
416 relatively strong classification capability.

417 Overall, the results presented in Supplementary Table 17 confirm that the DELTA framework remains  
418 robust and stable under realistic noisy conditions, even when both feature perturbations and label  
419 corruption are present. This robustness is particularly important for practical battery datasets, where  
420 measurement noise and imperfect labeling are unavoidable.

421

422 **Supplementary Table 18 Comparing with other methods of the *batch\_dataset***423 **Supplementary Table 18-1 Accuracy**

Label availability	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
DELTA	0.9003 ± 0.0265	0.8849 ± 0.0301	0.8713 ± 0.0307	0.8723 ± 0.0419	0.8452 ± 0.0488	0.8271 ± 0.0497	0.7817 ± 0.0371	0.6892 ± 0.0855
FixMatch	0.8684 ± 0.0285	0.8602 ± 0.0322	0.8521 ± 0.0291	0.8391 ± 0.0328	0.8271 ± 0.0365	0.8020 ± 0.0493	0.8074 ± 0.0428	0.7558 ± 0.0285
SPRED	0.8916 ± 0.0324	0.8720 ± 0.0347	0.8602 ± 0.0343	0.8542 ± 0.0326	0.8304 ± 0.0393	0.8059 ± 0.0451	0.8171 ± 0.0398	0.7658 ± 0.0317
ST_RF	0.9194 ± 0.0149	0.9113 ± 0.0230	0.8850 ± 0.0292	0.8823 ± 0.0255	0.8633 ± 0.0372	0.8271 ± 0.0437	0.7664 ± 0.0560	0.6459 ± 0.0634
ST_SVM	0.8886 ± 0.0158	0.8695 ± 0.0268	0.8660 ± 0.0275	0.8596 ± 0.0205	0.8306 ± 0.0310	0.7953 ± 0.0370	0.7527 ± 0.0411	0.6521 ± 0.0533
UDA	0.8354 ± 0.0392	0.8231 ± 0.0346	0.8086 ± 0.0385	0.8038 ± 0.0435	0.7896 ± 0.0453	0.7604 ± 0.0635	0.7640 ± 0.0537	0.7262 ± 0.0286

424

425 **Supplementary Table 18-2 F1 score**

Label availability	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
DELTA	0.9001 ± 0.0264	0.8851 ± 0.0293	0.8713 ± 0.0307	0.8723 ± 0.0416	0.8438 ± 0.0503	0.8262 ± 0.0526	0.7826 ± 0.0378	0.6878 ± 0.0878
FixMatch	0.8669 ± 0.0291	0.8591 ± 0.0326	0.8512 ± 0.0294	0.8379 ± 0.0334	0.8259 ± 0.0370	0.8006 ± 0.0499	0.8056 ± 0.0446	0.7548 ± 0.0280
SPRED	0.8908 ± 0.0327	0.8713 ± 0.0351	0.8597 ± 0.0346	0.8534 ± 0.0329	0.8292 ± 0.0397	0.8050 ± 0.0454	0.8164 ± 0.0400	0.7650 ± 0.0313
ST_RF	0.9187 ± 0.0149	0.9104 ± 0.0233	0.8837 ± 0.0294	0.8814 ± 0.0257	0.8624 ± 0.0374	0.8262 ± 0.0438	0.7660 ± 0.0585	0.6398 ± 0.0683
ST_SVM	0.8886 ± 0.0156	0.8691 ± 0.0275	0.8654 ± 0.0276	0.8592 ± 0.0213	0.8309 ± 0.0306	0.7960 ± 0.0375	0.7534 ± 0.0405	0.6442 ± 0.0581
UDA	0.8318 ± 0.0407	0.8203 ± 0.0358	0.8065 ± 0.0394	0.8004 ± 0.0465	0.7862 ± 0.0485	0.7557 ± 0.0673	0.7602 ± 0.0576	0.7233 ± 0.0295

426

427 Supplementary Table 18 reports the detailed accuracy and F1 score results for DELTA and several  
428 baseline semi-supervised learning methods on the *batch\_dataset* across different label-hiding ratios. The  
429 evaluated methods include FixMatch, SPRED, ST\_RF, ST\_SVM, and UDA, and the results are presented  
430 as mean ± standard deviation over repeated runs.

431 Overall, all methods show a gradual performance decline as the hide ratio increases, reflecting the  
432 increasing difficulty of the classification task when fewer labeled samples are available. Among the  
433 baseline methods, ST\_RF achieves the highest performance under very low hide ratios, with an F1 score  
434 of 0.9254 at a hide ratio of 0.1, followed closely by GELTA and SPRED. However, as the hide ratio  
435 increases, the performance of most baseline methods deteriorates more rapidly. For example, when the  
436 hide ratio reaches 0.8, the F1 score of ST\_RF decreases to approximately 0.766, while FixMatch and  
437 UDA show even lower values.

438 Compared with these baseline approaches, DELTA demonstrates better scalability with respect to label  
439 availability, as illustrated in Fig. 4a. Although some traditional semi-supervised methods achieve  
440 relatively strong performance when labeled data are abundant, their performance degrades significantly  
441 when labels become scarce. In contrast, DELTA maintains more stable classification performance across  
442 a wide range of label availability conditions, highlighting its ability to effectively exploit both labeled  
443 and unlabeled data.

444 Another notable observation from Table 18 is the close agreement between accuracy and F1 score across  
445 all methods, indicating that the classification results are relatively balanced across different classes  
446 without severe class imbalance effects. The reported standard deviations are also relatively small,  
447 suggesting that the results are stable across repeated experiments.

448 Overall, the results summarized in Supplementary Table 18 provide quantitative evidence supporting the  
449 superior scalability and stability of the DELTA framework on the batch-based dataset.

450

451 **Supplementary Table 19 Comparing with other methods of the *material\_dataset***

452 Supplementary Table 19-1 Accuracy

Label availability	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
DEALTA	0.9672 ±	0.9745 ±	0.9492 ±	0.9128 ±	0.8769 ±	0.8657 ±	0.7365 ±	0.5517 ±
	0.0384	0.0332	0.0540	0.0569	0.0617	0.0679	0.1013	0.1200
FixMatch	0.8150 ±	0.8003 ±	0.7718 ±	0.7518 ±	0.6936 ±	0.6769 ±	0.6956 ±	0.5955 ±
	0.0406	0.0657	0.0525	0.0639	0.0597	0.0589	0.0740	0.0837
SPRED	0.8379 ±	0.8197 ±	0.7915 ±	0.7675 ±	0.7056 ±	0.6987 ±	0.7153 ±	0.6064 ±
	0.0345	0.0521	0.0434	0.0602	0.0635	0.0522	0.0719	0.0778
ST_RF	0.8806 ±	0.8548 ±	0.8368 ±	0.8048 ±	0.7464 ±	0.6636 ±	0.5550 ±	0.4169 ±
	0.0598	0.0724	0.0783	0.0831	0.0664	0.0947	0.1016	0.0978
ST_SVM	0.8733 ±	0.8623 ±	0.8478 ±	0.8082 ±	0.7898 ±	0.7796 ±	0.7360 ±	0.5321 ±
	0.0399	0.0480	0.0455	0.0476	0.0560	0.0525	0.0515	0.1171
UDA	0.7399 ±	0.7436 ±	0.7136 ±	0.6873 ±	0.6596 ±	0.6404 ±	0.6608 ±	0.5929 ±
	0.0504	0.0622	0.0517	0.0624	0.0578	0.0552	0.0725	0.0816

453

454 Supplementary Table 19-2 F1 score

Label availability	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
DELTA	0.9672 ±	0.9747 ±	0.9492 ±	0.9134 ±	0.8757 ±	0.8639 ±	0.7315 ±	0.5342 ±
	0.0382	0.0326	0.0541	0.0557	0.0614	0.0686	0.1039	0.1239
FixMatch	0.8051 ±	0.7838 ±	0.7562 ±	0.7405 ±	0.6779 ±	0.6586 ±	0.6781 ±	0.5750 ±
	0.0449	0.0784	0.0681	0.0715	0.0626	0.0647	0.0786	0.1001
SPRED	0.8306 ±	0.8097 ±	0.7799 ±	0.7598 ±	0.6946 ±	0.6848 ±	0.7026 ±	0.5891 ±
	0.0387	0.0589	0.0514	0.0661	0.0660	0.0592	0.0757	0.0899
ST_RF	0.8790 ±	0.8549 ±	0.8379 ±	0.8041 ±	0.7462 ±	0.6625 ±	0.5132 ±	0.3134 ±
	0.0631	0.0736	0.0786	0.0852	0.0670	0.0955	0.1360	0.1427
ST_SVM	0.8710 ±	0.8613 ±	0.8476 ±	0.8076 ±	0.7879 ±	0.7776 ±	0.7152 ±	0.4617 ±
	0.0406	0.0482	0.0462	0.0474	0.0556	0.0546	0.0621	0.1631
UDA	0.7216 ±	0.7220 ±	0.6963 ±	0.6727 ±	0.6428 ±	0.6182 ±	0.6415 ±	0.5697 ±
	0.0530	0.0744	0.0596	0.0665	0.0593	0.0637	0.0748	0.0936

455

456 Supplementary Table 19 presents the detailed accuracy and F1 score comparisons for the  
 457 *material\_dataset*, corresponding to the results shown in Fig. 4b. Similar to Table 18, the evaluated  
 458 methods include FixMatch, SPRED, ST\_RF, ST\_SVM, and UDA.

459 In general, the performance of all methods decreases as the hide ratio increases, which corresponds to a  
 460 reduction in the proportion of labeled training samples. Among the baseline methods, GMMSSLM

461 achieves the highest performance under very low hide ratios, with an accuracy and F1 score close to 0.97  
462 when the hide ratio is 0.1–0.2. However, its performance declines rapidly as the hide ratio increases. For  
463 example, when the hide ratio reaches 0.8, its F1 score drops to approximately 0.53, indicating a  
464 significant degradation when labeled data become scarce.

465 Other baseline methods show similar trends. ST\_RF and ST\_SVM initially achieve moderate  
466 performance but experience substantial performance loss at higher hide ratios. FixMatch and UDA  
467 exhibit relatively lower performance across most label availability levels, reflecting the challenges of  
468 applying general-purpose semi-supervised learning methods to battery degradation data.

469 Compared with these baseline methods, DELTA consistently achieves higher F1 scores across nearly all  
470 label availability levels, as shown in Fig. 4b. In particular, at moderate hide ratios (e.g., 0.3–0.5), where  
471 the classification problem becomes more challenging, DELTA maintains significantly higher  
472 performance than the baseline methods. This improvement demonstrates the advantage of incorporating  
473 domain-aware feature extraction and statistical modelling in the DELTA framework.

474 Additionally, the relatively small standard deviations reported in Table 19 indicate that the performance  
475 improvements are statistically stable across repeated experiments. These results further confirm that  
476 DELTA provides more robust and reliable classification performance for partially labeled battery  
477 datasets.

478

479 **Supplementary Table 20. Computational efficiency comparison on the *batch\_dataset***

Method	Accuracy	F1 score	Training time(s)	Predict time(s)	Total time(s)
DELTA	0.8849 ± 0.0301	0.8851 ± 0.0293	0.0747 ± 0.0083	0.0001 ± 0.0002	0.0748 ± 0.0083
FixMatch	0.8602 ± 0.0322	0.8591 ± 0.0326	8.2112 ± 0.4373	0.0007 ± 0.0001	8.2119 ± 0.4374
SPRED	0.8720 ± 0.0347	0.8713 ± 0.0351	7.2506 ± 0.5109	0.0007 ± 0.0001	7.2513 ± 0.5110
ST_RF	0.9113 ± 0.0230	0.9104 ± 0.0233	0.2260 ± 0.0691	0.0000 ± 0.0000	0.2260 ± 0.0691
ST_SVM	0.8695 ± 0.0268	0.8691 ± 0.0275	0.4097 ± 0.1297	0.0000 ± 0.0000	0.4097 ± 0.1297
UDA	0.8231 ± 0.0346	0.8203 ± 0.0358	8.0285 ± 0.4345	0.0007 ± 0.0001	8.0291 ± 0.4345

480

481 Supplementary Table 20 summarizes the computational efficiency and classification performance of  
 482 different methods on the *batch\_dataset*, including accuracy, F1 score, training time per run, inference  
 483 time per run, and total runtime.

484 Among the baseline methods, ST\_RF achieves the highest predictive performance with an accuracy of  
 485  $0.9113 \pm 0.0230$  and an F1 score of  $0.9104 \pm 0.0233$ , while maintaining a relatively short training time  
 486 of 0.2260 s. DELTA also demonstrates strong performance, reaching an accuracy of  $0.8849 \pm 0.0301$   
 487 with a very short training time of 0.0747 s, making it one of the most computationally efficient baseline  
 488 methods.

489 In contrast, deep semi-supervised learning approaches such as FixMatch, SPRED, and UDA require  
 490 substantially longer training times, typically around 7-8 s per run, which is one to two orders of  
 491 magnitude higher than those of traditional machine-learning-based approaches. Despite the higher  
 492 computational cost, their classification performance remains lower than that of the best-performing  
 493 methods on this dataset.

494 Across all methods, the inference time per sample is extremely small (on the order of  $10^{-4}$ - $10^{-3}$  s),  
 495 indicating that prediction itself is computationally inexpensive once the models are trained. However,  
 496 the substantial differences in training time highlight the importance of algorithmic efficiency when  
 497 models need to be retrained repeatedly or applied to large-scale datasets.

498 Overall, these results demonstrate clear trade-offs between predictive performance and computational  
 499 cost among different semi-supervised learning strategies on the *batch\_dataset*, with simpler models  
 500 achieving significantly faster training while maintaining competitive classification accuracy.

501

502 **Supplementary Table 21. Computational efficiency comparison on the *material\_dataset***

Method	Accuracy	F1 score	Training time(s)	Predict time(s)	Total time(s)
GMMSSLM	0.8657 ± 0.0679	0.8639 ± 0.0686	0.2288 ± 0.1225	0.0001 ± 0.0002	0.2290 ± 0.1224
DELTA	0.6769 ± 0.0589	0.6586 ± 0.0647	0.9952 ± 0.0133	0.0002 ± 0.0001	0.9955 ± 0.0133
SPRED	0.6987 ± 0.0522	0.6848 ± 0.0592	0.8708 ± 0.0097	0.0002 ± 0.0001	0.8711 ± 0.0097
ST_RF	0.6636 ± 0.0947	0.6625 ± 0.0955	0.0939 ± 0.0139	0.0000 ± 0.0000	0.0939 ± 0.0139
ST_SVM	0.7796 ± 0.0525	0.7776 ± 0.0546	0.0139 ± 0.0060	0.0000 ± 0.0000	0.0139 ± 0.0060
UDA	0.6404 ± 0.0552	0.6182 ± 0.0637	0.9686 ± 0.0178	0.0003 ± 0.0001	0.9689 ± 0.0179

503

504 Supplementary Table 21 presents the computational efficiency and classification performance of  
505 different methods on the *material\_dataset*. Compared with the *batch\_dataset*, the overall computational  
506 cost on this dataset is substantially lower, with the training time of all baseline methods remaining below  
507 1s per run, reflecting the smaller dataset size and reduced computational complexity.

508 Among the evaluated methods, DELTA achieves the best predictive performance, with an accuracy of  
509  $0.8657 \pm 0.0679$  and an F1 score of  $0.8639 \pm 0.0686$ , while requiring a moderate training time of 0.2288  
510 s. ST\_SVM also shows relatively strong performance, reaching an accuracy of  $0.7796 \pm 0.0525$  with an  
511 extremely short training time of only 0.0139 s, highlighting its high computational efficiency.

512 Deep semi-supervised approaches such as FixMatch, SPRED, and UDA exhibit longer training times  
513 (approximately 0.87-1.00 s) while delivering comparatively lower predictive performance on this dataset.  
514 In contrast, ST\_RF achieves a very short training time (0.0939 s) but shows relatively lower classification  
515 accuracy.

516 Similar to the observations on the *batch\_dataset*, the inference time for all methods remains extremely  
517 small, typically on the order of  $10^{-4}$  s per sample, indicating that prediction is computationally  
518 inexpensive once the models are trained.

519 Overall, these results suggest that while computational costs are generally lower on the *material\_dataset*,  
520 there remain clear differences in the balance between training efficiency and predictive performance  
521 among different algorithms.

522

523 **Supplementary Table 22. Benchmark settings used for fair comparison of DELTA and baseline**  
 524 **methods**

Category	Parameter	Value
Data partition	Cross-validation	10-fold stratified cross-validation
Data partition	Repeated runs	3
Data partition	Random seed	42
Label availability	Hidden-label ratio	0.7
Label availability	Available-label ratio	0.3
Training	Epochs	100
Training	Evaluation interval	Every 10 epochs
Batch size	Labeled samples	16
Batch size	Unlabeled samples	32
Batch size	Test samples	64
Network architecture	Backbone	Three-layer MLP
Network architecture	Hidden dimension	64
Network architecture	Dropout	0.3
Optimization	Optimizer	Adam
Optimization	Learning rate	$1 \times 10^{-3}$
Optimization	Weight decay	$1 \times 10^{-3}$
Data augmentation	Weak augmentation noise	0.02
Data augmentation	Strong augmentation noise	0.10
FixMatch	Confidence threshold	0.95
FixMatch	Unsupervised loss weight	1.0
UDA	Temperature	0.7
UDA	Unsupervised loss weight	1.0
SPRED	Regularization weight	0.1
Self-training	Confidence threshold	0.9
Self-training	Maximum iterations	10

525 All methods were trained and evaluated using identical data partitions, label-availability settings, input  
 526 features, and evaluation metrics. Neural-network-based methods shared the same backbone architecture  
 527 and optimization settings to ensure a fair comparison. Method-specific hyperparameters were fixed  
 528 according to the unified benchmark configuration.

529

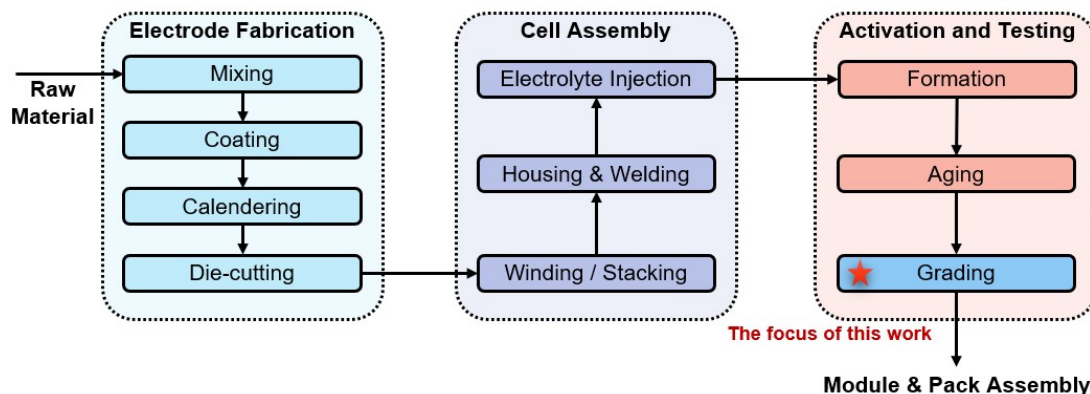
530 **Supplementary Table 23. Definitions of key datasets and classification schemes used in this study.**

Term	Definition
Base dataset	The original dataset containing cells with complete lifetime labels used for model development and evaluation.
Extended dataset	Additional battery data incorporated to increase data diversity and simulate practical limited-label scenarios.
Labeled extended dataset	Extended dataset samples whose lifetime labels are available during model training.
Unlabeled extended dataset	Extended dataset samples whose lifetime labels are intentionally treated as unavailable during model training.
Batch dataset	A dataset constructed by grouping cells according to manufacturing batch. Cells within the same batch serve as the reference population for batch-level lifetime classification.
Material dataset	A dataset constructed by grouping cells according to chemistry or material system, regardless of manufacturing batch.
Batch-based classification	Lifetime classification performed relative to the lifetime distribution of cells within the same batch, primarily used for manufacturing quality assessment and consistency evaluation.
Material-based classification	Lifetime classification performed relative to the lifetime distribution of cells sharing the same material system, primarily used for cross-batch comparison and technology benchmarking.
Label availability	The proportion of samples with accessible lifetime labels during model training.
Lifetime category (L/N/H)	Relative lifetime labels indicating Low, Normal, and High lifetime performance within a specified reference population.

531

532

533 **Supplementary Note 1. Relationship between battery manufacturing workflow and early-cycle**  
534 **quality assessment**



535

536 As illustrated in Figure, the manufacturing process of a typical lithium-ion battery consists of three major  
537 stages: Electrode Fabrication, Cell Assembly, and Activation and Testing. Electrode fabrication includes  
538 slurry mixing, coating, calendaring, and die-cutting, whereas cell assembly involves electrode  
539 winding/stacking, housing and welding, and electrolyte injection. The present study focuses on the final  
540 stage, Activation and Testing, because it represents the last manufacturing-end quality control procedure  
541 before cells are released for downstream applications and therefore provides the most direct basis for  
542 industrial quality evaluation.

543 The Activation and Testing stage mainly comprises three sequential procedures: formation, aging, and  
544 capacity grading. During formation, newly assembled cells undergo initial charge-discharge activation,  
545 promoting electrochemical conversion of active materials and the formation of stable solid electrolyte  
546 interphase (SEI) layers. Following formation, cells enter the aging stage, during which they are stored  
547 under controlled conditions for several days to facilitate electrolyte wetting and further stabilization of  
548 the SEI layer. A key quality-control objective during aging is the identification of abnormal self-  
549 discharge behaviour, which may originate from micro-short circuits, metallic contaminants, separator  
550 defects, or other latent manufacturing imperfections. Subsequently, cells undergo capacity grading,  
551 where standardized charge-discharge tests are performed to determine practical capacity and internal  
552 resistance. Based on these measurements, cells are classified into different quality grades for subsequent  
553 utilization.

554 The datasets used in this work originate from individual cells that successfully passed industrial grading  
555 procedures and had not yet been assembled into battery packs. Importantly, battery degradation is a  
556 gradual and continuous process rather than a discrete event. Consequently, the electrochemical state of a  
557 cell during its first few operational cycles remains highly similar to the state established immediately  
558 after formation, aging, and capacity grading. From this perspective, early-cycle behaviour can be  
559 regarded as a direct continuation of the manufacturing-end cell condition.

560 Although the data analyzed in this study do not contain explicit upstream manufacturing parameters such

561 as coating thickness, calendaring density, or electrolyte filling volume, the early-cycle voltage–capacity  
562 characteristics inherently reflect the cumulative effects of the preceding manufacturing and quality-  
563 control processes. Variations introduced during electrode fabrication, cell assembly, and activation  
564 procedures can influence the electrochemical state of the cell and subsequently manifest as subtle  
565 differences in early-cycle behaviour. Therefore, early-cycle signals provide observable end-of-line  
566 signatures of manufacturing-induced quality variation, even when direct process information is  
567 unavailable.

568 Current industrial quality grading primarily relies on capacity, internal resistance, and self-discharge  
569 measurements. While these indicators are effective for detecting severe defects and assessing immediate  
570 cell performance, they provide limited information regarding long-term degradation trajectories. As a  
571 result, cells exhibiting similar grading metrics may still experience substantially different ageing  
572 behaviour during subsequent operation. The central objective of this study is therefore not to reconstruct  
573 manufacturing parameters or identify specific process root causes, but rather to extract latent quality  
574 information from routinely collected early-cycle data and establish its relationship with future lifetime  
575 outcomes.

576 Accordingly, the proposed framework should be viewed as an outcome-oriented manufacturing quality  
577 assessment approach. By linking end-of-line testing signatures to subsequent degradation behaviour, the  
578 framework provides a practical means of evaluating manufacturing-related quality variation under  
579 realistic industrial conditions where direct manufacturing process data are often inaccessible.

580

## 581 Supplementary Note 2. Analysis of Data Heterogeneity

582 The DELTA framework is designed to achieve robust generalization across both dataset domains and material  
583 domains. Accordingly, 4 major material systems (LFP, NCA, NCM, and NCANCM), drawn from the six datasets,  
584 were included in the evaluation. Even when nominal material labels are identical, the underlying chemical  
585 compositions differ across datasets. For example, the NCM chemistry in the XJTU dataset is  $LiNi_{0.5}Co_{0.2}Mn_{0.3}O_2$ ,  
586 whereas the TJU dataset adopts  $Li_{0.86}Ni_{0.86}Co_{0.11}Mn_{0.07}O_2$ , and the Stanford dataset uses  $LiNi_{0.5}Mn_{0.3}Co_{0.2}O_2$ .  
587 Similarly, the NCA chemistry in the TJU dataset is  $Li_{0.86}Ni_{0.86}Co_{0.11}Al_{0.03}O_2$ , while the CAS dataset employs  
588  $LiNi_xCo_yAl_{1-x-y}O_2$ . Batteries belonging to the same material system are therefore grouped together, and material  
589 differences are treated as a primary driver of ageing behaviour. In addition, even within the same material category,  
590 cells from different datasets exhibit variations in nominal capacity and cutoff voltage. To mitigate the influence of  
591 these factors, voltage signals were uniformly sampled over the entire voltage window, and normalization was applied  
592 during HI construction. All cells follow a CCCV charging protocol and CC discharging protocol, with the main  
593 operational differences arising from charge-rate settings. As the extracted HIs are derived exclusively from charging-  
594 phase data, only charge-rate groupings are considered. Taken together, material system, temperature and charge rate  
595 are identified as the three dominant factors governing battery ageing, with material effects being the most  
596 pronounced. Consequently, in the subsequent LME model, material is treated as a random effect, while temperature  
597 and charge rate are modelled as fixed effects, enabling quantification of material-dependent ageing behaviour under  
598 controlled temperature and charge-rate conditions.

599

600 **Supplementary Note 3. Three-class scheme based on  $1\sigma$  boundaries of a Gaussian distribution**

601 In extracting EOL labels, two key questions must be addressed: (1) how many EOL cycle-life categories  
602 should be defined; and (2) which specific partitioning strategy should be adopted for a review of related  
603 literature. In early-life screening applications, the primary objective is to identify cells with potentially  
604 large performance deviations, thereby ensuring consistency within battery packs.

605 From the perspective of category design, cycle-life classification generally includes binary classification  
606 and multi-class classification. Binary classification typically distinguishes between “good” and “bad”  
607 batteries, which in practical applications often corresponds to identifying “long-life batteries” and “short-  
608 life batteries.” Multi-class classification provides a more refined categorization. In practice, batteries are  
609 commonly divided into three categories: short-life, medium-life, and long-life. The objective is to remove  
610 abnormally short-lived or excessively long-lived cells from a large population of otherwise normal  
611 batteries, thereby ensuring stable and controllable operation of battery packs.

612 Regarding specific label assignment methods, existing studies primarily adopt two strategies: proportion-  
613 based partitioning and threshold-based partitioning. Proportion-based partitioning assigns labels based  
614 on the quantiles of the cycle-life distribution across the entire dataset. For example, Lucas Murphy<sup>8</sup>  
615 defined batteries within the top 25% quantile of cycle life as long-life cells, considering them suitable  
616 for secondary use. Minzheng Hu<sup>9</sup> classified batteries into three categories, L (long-life), M (medium-  
617 life), and S (short-life), based on their quantile positions within the overall dataset.

618 Threshold-based partitioning, in contrast, classifies batteries by setting explicit cycle-life thresholds,  
619 often in combination with proportion-based methods. For instance, Xuelu Wang<sup>10</sup> first divided the dataset  
620 according to cycle-life quantiles and then defined corresponding thresholds for three-class classification.  
621 Sandro Stock<sup>11</sup> directly adopted 150 cycles as the threshold for binary classification and used 150 and  
622 300 cycles as thresholds for three-class classification. Zicheng Fei<sup>12</sup> selected 550 cycles as the binary  
623 classification threshold and used 550 and 1200 cycles for three-class classification, while Yongzhi  
624 Zhang<sup>13</sup> also used 550 cycles as the binary classification threshold.

625 It should be noted that the selection of thresholds should depend on the characteristics of the batteries  
626 under investigation; therefore, threshold-based partitioning is often combined with proportion-based  
627 methods to improve rationality and interpretability.

628 As shown in **Fig. 2c**, the distribution of battery cycle life in practical scenarios is often approximately  
629 Gaussian, with the majority of cells exhibiting intermediate lifetimes and relatively few cells at the  
630 extremes. Accordingly, a three-class scheme based on the  $1\sigma$  boundaries of a Gaussian distribution,  
631 defined by the quantiles of the dataset-specific cycle-life distribution, better reflects practical  
632 requirements.

633 Specifically, cells with cycle life falling within the interval  $\mu \pm 1\sigma$  (where  $\mu$  and  $\sigma$  denote the mean and

634 standard deviation, respectively) exhibit high concentration and limited performance variability and  
635 therefore constitute the core population for assembling highly consistent battery packs. By contrast, cells  
636 located below  $\mu - 1\sigma$  or above  $\mu + 1\sigma$ , each accounting for approximately 16% of the population, are  
637 of particular interest for early screening. Cells in the lower tail  $(0, \mu - 1\sigma)$  exhibit significantly reduced  
638 lifetime and may suffer from manufacturing-related defects, such as electrode material imperfections,  
639 non-uniform coating, or insufficient electrolyte filling, and are therefore likely to become “weak links”  
640 if incorporated into a pack. Cells in the upper tail  $(\mu + 1\sigma, +\infty)$ , although longer-lived, deviate  
641 substantially from the mainstream population and may introduce system-level issues, including capacity  
642 imbalance and increased difficulty in equalization. Consequently, both groups warrant careful review to  
643 ensure stable, consistent and reliable battery-pack performance in large-scale applications. To emulate  
644 incomplete label availability in real battery production settings, we further constructed a label-missing  
645 simulation mechanism based on feature dispersion. Specifically, the variance of the sample features was  
646 used as a proxy for sample complexity, with samples exhibiting higher dispersion assigned a higher  
647 probability of label missingness.

648

649 **Supplementary Note 4. Economic analysis under end-of-line screening scenario**

650 At the end of battery production, rapid and reliable screening of potentially risky cells is essential for  
651 ensuring the consistency, reliability, and safety of downstream battery pack integration. However, in  
652 practical manufacturing environments, end-of-life (EOL) labeled data are costly and time-consuming to  
653 obtain, and therefore typically exist only in limited quantities. In contrast, large volumes of early-cycle  
654 operational data are routinely generated during manufacturing but often remain unlabeled and  
655 underutilized.

656 To address this realistic industrial setting, as shown in Fig. S1, this study considers an end-of-line  
657 screening scenario characterized by scarce labeled data and abundant unlabeled data. The proposed  
658 DELTA framework uses the first  $k$  early cycles as input features and performs a three-class classification  
659 of battery EOL cycle life under a semi-supervised learning setting where only a fraction  $r$  of the samples  
660 are labeled. The three classes correspond to abnormally high, normal, and abnormally low lifetimes,  
661 enabling the identification of potentially risky batteries at the production stage. Batteries classified as  
662 normal can proceed directly to downstream processes, whereas batteries predicted as risky require  
663 additional inspection. Early identification of such batteries helps mitigate potential safety risks and  
664 improve pack-level consistency.

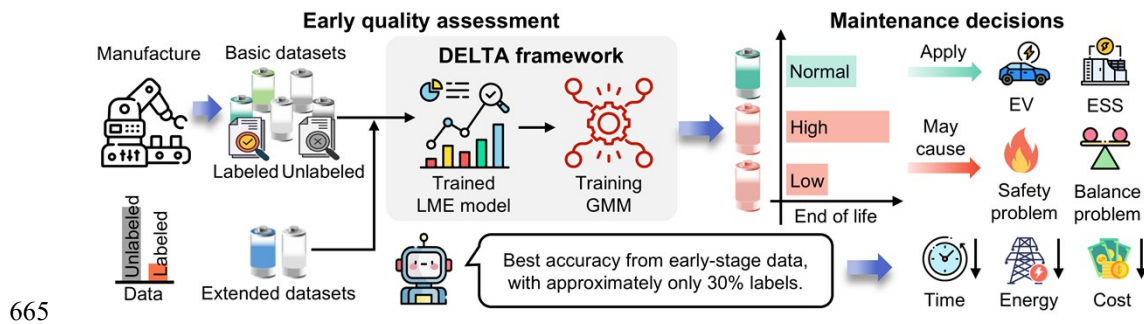


Fig. S1. Deployment of DELTA framework

667 To evaluate the economic feasibility of the proposed framework, we construct an expected economic  
668 cost model considering three main components:

- 669 1. Risk-related losses caused by undetected risky batteries
- 670 2. Additional re-inspection costs for batteries flagged as risky
- 671 3. Label acquisition costs required for model training

672 **Economic assumptions**

673 The economic analysis is conducted under a representative industrial production scale corresponding to  
674 1.3 TWh<sup>14</sup> annual battery production. Assuming a nominal energy capacity of 100 Wh per cell, the annual  
675 production volume is approximately

676 
$$N_{annual} = 1.3 \times \frac{10^{12}}{100} \approx 1.3 \times 10^{10} \text{ cells per year} \quad (1)$$

677 For batteries that are incorrectly classified as normal (false negatives), a potential field failure may occur  
 678 during downstream use. The probability that such a missed risky battery leads to a failure event is  
 679 assumed to lie within

680 
$$p_{fail} \in [10^{-5}, 2 \times 10^{-5}] \quad (2)$$

681 and the midpoint value is used in the calculation.

682 The economic loss associated with a single field failure event is assumed to lie within

683 
$$c_{event} \in [20,000, 80,000] \text{ USD} \quad (3)$$

684 Therefore, the expected loss per undetected risky battery is

685 
$$c_{risk} = p_{fail} \times c_{event} \quad (4)$$

686 For batteries predicted as risky, an additional end-of-line re-inspection procedure is required. This  
 687 includes equipment operation, labor, and additional testing costs. The unit re-inspection cost is assumed  
 688 to lie within

689 
$$c_{recheck} \in [0.08, 0.25] \text{ USD} \quad (5)$$

690 Meanwhile, labeled samples required for model training must be obtained through full lifetime testing.  
 691 The unit labeling cost is assumed to lie within

692 
$$c_{label} \in [500, 1000] \text{ USD} \quad (6)$$

693 In this study, the number of labeled samples collected per year is assumed to be

694 
$$N_{label} = 10,000 \quad (7)$$

695 In addition, the deployment of the AI-based screening system introduces an annual operational cost  
 696 including computing infrastructure, data management, and model maintenance. This annual system cost  
 697 is assumed to lie within

698 
$$C_{system} \in [1.5, 4.0] \text{ million USD} \quad (8)$$

699 Midpoint values of these ranges are used in the economic calculations.

700 **Expected economic cost model**

701 Let  $TP$ ,  $FP$ , and  $FN$  denote the numbers of correctly detected risky batteries, falsely flagged normal  
 702 batteries, and missed risky batteries in the evaluation dataset, respectively. Let  $N$  denote the total number  
 703 of evaluated samples.

704 The corresponding rates are

$$705 \quad FN \text{ rate} = \frac{FN}{N}, \text{ flag rate} = \frac{TP + FP}{N} \quad (9)$$

706 The annual expected economic cost of the proposed screening strategy is defined as

$$707 \quad S_{DELTA} = FN_{rate} \times N_{annual} \times c_{risk} + \text{flag rate} \times N_{annual} \times c_{recheck} + N_{label} \times c_{label} + C_{system} \quad (10)$$

708 This formulation captures the trade-off among risk mitigation, inspection workload, and model  
 709 deployment cost.

### 710 **Baseline strategies**

711 Two baseline strategies are introduced for comparison.

#### 712 **Baseline 1: No screening.**

713 In this strategy, all batteries are directly deployed without end-of-line screening. Although no inspection  
 714 or labeling cost is incurred, all risky batteries remain undetected. The expected annual cost therefore  
 715 arises entirely from potential risk-related losses:

$$716 \quad S_{B1} = \text{risk rate} \times N_{annual} \times c_{risk} \quad (11)$$

#### 717 **Baseline 2: Random sampling screening**

718 In this strategy, a proportion  $s$  of batteries is randomly selected for inspection. Because the sampling  
 719 process does not target risky batteries, its risk coverage capability is limited. An effective risk coverage  
 720 coefficient  $\eta_{B2} = 0.5$  is introduced to represent the fraction of risky batteries that can be detected through  
 721 random inspection. The residual risk rate is therefore

$$722 \quad (1 - \eta_{B2}s) \times \text{risk rate} \quad (12)$$

723 The total annual cost of this strategy becomes

$$724 \quad S_{B2} = (1 - \eta_{B2}s) \times \text{risk rate} \times N_{annual} \times c_{risk} + s \times N_{annual} \times c_{recheck} \quad (13)$$

725 In this study, a representative sampling ratio  $s = 0.2$  is considered.

### 726 **Economic evaluation results**

727 Under the above assumptions, the proposed DELTA framework yields an average annual economic cost  
728 of approximately 2433.67 million USD at the considered production scale of 1.3 TWh per year. In  
729 comparison, the no-screening strategy results in an expected annual cost of approximately 6500.16  
730 million USD, as all risky batteries remain undetected and may lead to downstream losses. The random  
731 sampling strategy with sampling ratio  $s = 0.2$  yields an annual cost of approximately 6279.15 million  
732 USD, providing only limited reduction in risk-related losses due to its lack of targeted detection  
733 capability. Overall, the proposed framework reduces the expected annual economic cost by  
734 approximately 4066 million USD, corresponding to a 62.6% reduction relative to the no-screening  
735 strategy.

#### 736 Discussion

737 The economic analysis further demonstrates the practical value of the proposed DELTA framework in  
738 large-scale battery manufacturing. At the production scale considered in this study (1.3 TWh annual  
739 output), the proposed method yields an expected annual economic cost of approximately 2433.67 million  
740 USD. In comparison, the no-screening baseline results in an expected cost of 6500.16 million USD per  
741 year, as all risky batteries remain undetected and may lead to downstream failures or costly corrective  
742 actions.

743 As a result, the proposed framework achieves an annual economic saving of approximately 4066.49  
744 million USD, corresponding to a 62.57% reduction in overall economic cost relative to the no-screening  
745 strategy. This large reduction primarily arises from the improved ability of the proposed method to  
746 identify potentially risky batteries early, thereby preventing costly downstream failures while  
747 maintaining a manageable inspection workload.

748 By contrast, the random sampling screening strategy with a sampling ratio of  $s = 0.2$  results in an annual  
749 cost of approximately 6279.15 million USD, which is slightly lower than the no-screening strategy but  
750 still significantly higher than the proposed method. In fact, relative to the no-screening baseline, the  
751 random sampling strategy yields an additional cost of approximately 221.02 million USD, corresponding  
752 to a 3.52% economic loss, indicating that untargeted inspection may introduce unnecessary inspection  
753 expenses without effectively mitigating risk.

754 Overall, these results demonstrate that the DELTA framework provides a significantly more cost-  
755 effective end-of-line screening strategy than both no-screening and random sampling. By combining  
756 targeted risk identification with semi-supervised learning that leverages abundant unlabeled  
757 manufacturing data, the proposed method substantially reduces risk-related losses while maintaining  
758 economically acceptable inspection and labeling costs. This highlights its strong potential for deployment  
759 in gigafactory-scale battery production systems.

760

761 **Supplementary Note 5. Feature extraction**

762 **Input signals** To enable accurate EOL cycle-life classification using early-cycle data, a set of HIs is  
763 constructed based on the evolution of the capacity degradation curve and the Q-V curve. All selected  
764 features are extracted exclusively from early-cycle data and are screened through correlation analysis  
765 with EOL cycle life. The resulting features fall into three categories: early normalized capacity features,  
766 Q-V curve-based features, and capacity degradation curve features.

767 **Early normalized capacity features** Let  $Q_{rated}$  denote the nominal capacity of a cell, and  $Q_k$  the  
768 discharge capacity at the  $k$ -th cycle. To eliminate the influence of differences in nominal capacity across  
769 cells, the normalized capacity is defined as

$$\tilde{Q}_k = \frac{Q_k}{Q_{rated}} \quad (14)$$

771 The normalized capacities at the 1st, 3rd and 5th cycles are selected as  $HI_1$ ,  $HI_2$  and  $HI_3$ , respectively.  
772 These features characterize the initial health state immediately after formation and the early capacity  
773 retention behaviour of the cell and represent the most intuitive and physically interpretable inputs for  
774 lifetime assessment.

775 **Q-V curve-based features** Let  $Q_k(V)$  denote the Q-V relationship at the  $k$ -th cycle. To capture changes  
776 in curve morphology during early cycling, a differential Q-V function between the 3rd and 5th cycles is  
777 defined as

$$Q(V) = Q_3(V) - Q_5(V) \quad (15)$$

779 This function reflects the distribution of capacity changes induced by ageing across the entire voltage  
780 window. Based on  $\Delta Q(V)$ , the following statistical descriptors are extracted as health indicators:

$$HI_4 = \frac{E[\Delta Q(V)]}{Q_{rated}}, HI_5 = Kurt(\Delta Q(V)), HI_6 = Skew(\Delta Q(V)) \quad (16)$$

782 Here, kurtosis quantifies the concentration of the distribution, while skewness captures its asymmetry.  
783 Together, these metrics reflect localized degradation heterogeneity in the Q-V curves over specific  
784 voltage regions.

785 **Capacity degradation curve features** To further characterize the trend of capacity decay with cycling,  
786 linear fitting is applied to early-cycle capacity data. For the normalized capacities from the 1st to the 3rd  
787 cycles, a linear model is fitted as

$$\tilde{Q}_k = a_1 k + b_1 \quad (17)$$

789 and the intercept term is defined as a health indicator:

790 
$$HI_7 = b_1 \tag{18}$$

791 Similarly, a linear fit is performed for the normalized capacity data from the 3rd to the 5th cycles, and  
792 the corresponding intercept is defined as

793 
$$HI_8 = b_2 \tag{19}$$

794 All health indicators are constructed solely from early-cycle data and are subjected to correlation analysis  
795 with EOL cycle life. Only features exhibiting a strong correlation with EOL are retained as model inputs.  
796 In total, 15 candidate features are initially extracted, of which the eight most informative are selected,  
797 thereby maintaining predictive performance while reducing feature redundancy.

798

799 **Supplementary Note 6. Linear mixed-effects model**

800 In practical battery operation, temperature, charge rate, and battery material are the most influential  
801 factors governing ageing behaviour, with material effects typically exhibiting the largest variability. To  
802 investigate and quantify the impact of material on battery ageing, the LME model treats temperature and  
803 charge rate as fixed effects and material as a random effect, thereby capturing and quantifying material-  
804 level random variation. The empirical best linear unbiased prediction (EBLUP) of the material random  
805 effect is subsequently extracted as the feature  $HI_m$  and used as input to the subsequent GMM-based semi-  
806 supervised classification framework.

807 Assume that for the  $i$ -th material cluster ( $i = 1, 2, \dots, g$ ), the observed cycle life of the  $j$ -th cell ( $j = 1, 2, \dots, m_i$ ) is denoted by  $Cycle_{ij}$ , with  $T_{ij}$  and  $C_{ij}$  representing the corresponding temperature and  
808 charge-rate groups, respectively. Here,  $g$  denotes the number of material clusters, and  $m_i$  the number of  
809 cells within the  $i$ -th cluster. The LME model is formulated as  
810

$$811 \quad Cycle_{ij} = \beta_0 + \beta_T T_{ij} + \beta_C C_{ij} + \alpha_i + e_{ij} \quad (20)$$

812 where  $\beta_0$  is the fixed-effect intercept,  $\beta_T$  and  $\beta_C$  are the regression coefficients for temperature and  
813 charge rate, respectively,  $\alpha_i$  denotes the random effect associated with the  $i$ -th material cluster, and  $e_{ij}$  is  
814 the residual error term. The random effects  $\{\alpha_i\}$  and residuals  $\{e_{ij}\}$  are assumed to be mutually  
815 independent with zero means and variances  $\sigma_\alpha^2$  and  $\sigma_e^2$ , respectively. No normality assumption is imposed  
816 on their distributions. The model treats different clusters as independent and assumes homogeneous  
817 within-cluster correlation. The estimated random effects  $\hat{\alpha}_i$  are extracted as the material-related health  
818 indicator  $HI_m$  and combined with the eight dynamic health indicators to form the input feature vector for  
819 classification.

820

821 **Supplementary Note 7. Semi-supervised Gaussian mixture model**

822 Battery health states are modelled using a finite Gaussian mixture model. Let the feature vector of the  
823  $j$ -th cell be

$$824 \quad y_j = (y_{j1}, y_{j2}, \dots, y_{jp})^\top \in R^p \quad (21)$$

825 which is assumed to be generated from one of  $g$  Gaussian components:

$$826 \quad Y_j \sim N(\mu_i, \Sigma_i) \text{ with probability } \pi_i \quad (22)$$

827 where  $\mu_i$  and  $\Sigma_i$  denote the mean vector and covariance matrix corresponding to the  $i$ -th health state, and

$$828 \quad \pi_i \text{ is the prior probability satisfying } \sum_{i=1}^g \pi_i = 1 .$$

829 Under the semi-supervised learning setting, the EOL labels of a subset of samples are observed, while  
830 the class memberships of the remaining samples are treated as latent variables during model estimation.  
831 For labeled samples, class information is directly incorporated into parameter estimation, whereas for  
832 unlabeled samples, class assignments are inferred via posterior probabilities.

833 **Modelling the label-missing mechanism**

834 In practical battery testing, missing health-state labels are often correlated with the degree of battery  
835 ageing. To account for this dependency, a label-missing indicator variable  $m_j$  is introduced, where  
836  $m_j = 1$  indicates that the class label of the  $j$ -th sample is missing, and  $m_j = 0$  indicates that the label is  
837 observed. The probability of label missingness is assumed to depend solely on the observed feature vector  
838  $y_j$  and is characterized through posterior classification uncertainty:

$$839 \quad Pr^{\text{[10]}}(m_j = 1 | y_j) = q(y_j; \theta, \xi) \quad (23)$$

840 where  $q(\cdot)$  is a logistic function whose input is the Shannon entropy computed from the posterior class  
841 probabilities  $\tau_{ij}$ :

$$842 \quad e_j = - \sum_{i=1}^g \tau_{ij} \log \tau_{ij} \quad (24)$$

843 By introducing the label-missing mechanism, the complete partially classified log-likelihood can be  
844 expressed as

$$845 \quad \log^{\text{[10]}} L_{pc}^{\text{full}}(\Psi; x_{pc}) = \log^{\text{[10]}} L_{pc}^{\text{ig}}(\Psi; x_{pc}) + \log^{\text{[10]}} L_{pc}^{\text{miss}}(\Psi; x_{pc}) \quad (25)$$

846 where  $\Psi$  denotes the vector of all unknown parameters,  $x_{pc}$  represents the partially classified samples,

847 and  $\log L_{pc}^{ig}$  corresponds to the likelihood obtained by ignoring the missingness mechanism. The  
 848 likelihood contribution of the missingness process is given by

$$849 \quad \log L_{pc}^{miss}(\theta, \xi; x_{pc}) = \sum_{j=1}^n [(1 - m_j) \log\{1 - q(y_j; \theta, \xi)\} + m_j \log\{q(y_j; \theta, \xi)\}] \quad (26)$$

850 Here,  $m_j$  is the missing-label indicator, and  $q(y_j; \theta, \xi)$  is the entropy-based logistic function. The detailed  
 851 procedures for computation and parameter estimation are provided in the Methods section. By explicitly  
 852 modelling the label-missing mechanism, the proposed framework is able to exploit missingness  
 853 information to improve both the accuracy and generalization of early-stage ageing classification.

854

855 **Supplementary Note 8. Evaluation Metrics**

856 **Accuracy** measures the overall proportion of correctly classified samples and is defined as

857 
$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (27)$$

858 **Recall** quantifies the model's ability to identify positive samples and is defined as

859 
$$Recall = \frac{TP}{TP + FN} \quad (28)$$

860 **Precision** represents the proportion of samples predicted as positive that are truly positive. Together with  
861 Recall, it is used to construct F-score-based composite metrics. Based on Precision and Recall, this study  
862 adopts both the F1 score and the F2 score as summary performance measures. The F1 score assigns equal  
863 importance to Precision and Recall and is defined as

864 
$$F1 = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall} \quad (29)$$

865 whereas the F2 score places greater emphasis on recall of the positive class and is defined as

866 
$$F2 = \frac{5 \cdot Precision \cdot Recall}{4 \cdot Precision + Recall} \quad (30)$$

867 To further evaluate classification consistency under class-imbalanced conditions, the Matthews  
868 correlation coefficient (MCC) is employed. MCC jointly considers TP, FP, TN and FN and is defined as

869 
$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (31)$$

870 In addition, to assess overall model performance across different decision thresholds, the area under the  
871 precision–recall curve (PR-AUC) is used. PR-AUC is obtained by integrating the precision–recall curve  
872 and provides a more informative measure of a model's ability to identify positive samples under  
873 imbalanced class distributions. By jointly considering these evaluation metrics, model performance in  
874 early-stage battery health-state classification is systematically analysed from multiple complementary  
875 perspectives.

876

877 **Reference**

- 878 1. Chen, C. *et al.* The Operation Dependence of  $C - N$  Fatigue for Lithium-Ion Batteries. *Advanced*  
879 *Energy Materials* **13**, 2300942 (2023).
- 880 2. Weng, A. *et al.* Predicting the impact of formation protocols on battery lifetime immediately after  
881 manufacturing. *Joule* **5**, 2971–2992 (2021).
- 882 3. Li, W. *et al.* One-shot battery degradation trajectory prediction with deep learning. *Journal of Power*  
883 *Sources* **506**, 230024 (2021).
- 884 4. Cui, X. *et al.* Data-driven analysis of battery formation reveals the role of electrode utilization in  
885 extending cycle life. *Joule* **8**, 3072–3087 (2024).
- 886 5. Zhu, J. *et al.* Data-driven capacity estimation of commercial lithium-ion batteries from voltage  
887 relaxation. *Nat Commun* **13**, 2261 (2022).
- 888 6. Wang, F., Zhai, Z., Zhao, Z., Di, Y. & Chen, X. Physics-informed neural network for lithium-ion  
889 battery degradation stable modeling and prognosis. *Nat Commun* **15**, 4332 (2024).
- 890 7. Tan, R. *et al.* BatteryLife: A Comprehensive Dataset and Benchmark for Battery Life Prediction.  
891 Preprint at <https://doi.org/10.48550/arXiv.2502.18807> (2025).
- 892 8. Murphy, L. & Crawford, C. Data-Driven Classification of Lithium-Ion Batteries for Second-Life  
893 Applications. *Journal of Energy Storage* <https://doi.org/10.2139/ssrn.5119726> (2025)  
894 doi:10.2139/ssrn.5119726.
- 895 9. Hu, M., Tao, S., Wang, Y. & Sun, Y. A Data-Driven Approach for Lithium-ion Battery Lifetime  
896 Classification Based on Early Cycles. in *2023 IEEE 7th Conference on Energy Internet and Energy*  
897 *System Integration (EI2)* 2208–2213 (IEEE, Hangzhou, China, 2023).  
898 doi:10.1109/EI259745.2023.10512808.
- 899 10. Wang, X., Meng, J. & Azib, T. A Comparative Study of Data-Driven Early-Stage End-of-Life  
900 Classification Approaches for Lithium-Ion Batteries. *Energies* **17**, 4485 (2024).
- 901 11. Stock, S. *et al.* Early Quality Classification and Prediction of Battery Cycle Life in Production Using  
902 Machine Learning. *Journal of Energy Storage* **50**, 104144 (2022).
- 903 12. Fei, Z., Zhang, Z., Yang, F. & Tsui, K.-L. Deep learning powered rapid lifetime classification of  
904 lithium-ion batteries. *eTransportation* **18**, 100286 (2023).
- 905 13. Zhang, Y., Zhao, M. & Xiong, R. Online data-driven battery life prediction and quick classification  
906 based on partial charging data within 10 min. *Journal of Power Sources* **594**, 234007 (2024).
- 907 14. Tao, S. *et al.* Collaborative and privacy-preserving retired battery sorting for profitable direct  
908 recycling via federated machine learning. *Nat Commun* **14**, 8032 (2023).
- 909