
Foundation Models Enabling Multi-Scale Battery Materials Discovery: From Molecules To Devices

Vidushi Sharma

IBM Almaden Research
San Jose, CA, USA
vidushis@ibm.com

Andy Tek

IBM Almaden Research
San Jose, CA, USA
atek@us.ibm.com

Maxwell Giammona

IBM Almaden Research
San Jose, CA, USA
Maxwell.Giammona@ibm.com

Murtaza Zohair

IBM Research Almaden
San Jose, CA, USA
mzohair@ibm.com

Nathaniel Park

IBM Research Almaden
San Jose, CA, USA
npark@us.ibm.com

Tim Erdmann

IBM Research Almaden
San Jose, CA, USA
tim.erdmann@ibm.com

Linda Sundberg

IBM Almaden Research
San Jose, CA, USA
lindas@us.ibm.com

Eduardo Soares

IBM Research Brazil
Rio de Janeiro, RJ, Brazil
eduardo.soares@ibm.com

Khanh Nguyen

IBM Almaden Research
San Jose, CA, USA
khanh.vinh.nguyen@ibm.com

Young-Hye Na

IBM Almaden Research
San Jose, CA, USA
yna@us.ibm.com

Emilio Ashton Vital Brazil

IBM Research Brazil
Rio de Janeiro, RJ, Brazil
evital@br.ibm.com

Abstract

1 Recent years have seen fast emergence and adoption of chemical foundation models
2 in computational material science for property prediction and generation tasks that
3 are focused mostly on small molecules or crystals. Despite these paradigm shifts,
4 integration of newly discovered materials in real world devices continues to be a
5 challenge due to design problems. New candidate material must be optimized to
6 achieve compatibility with other components in the system and deliver the target
7 performance. Chemical foundation model benchmarks must evaluate their scope
8 in predicting macro scale outcomes that are the result of chemical interactions
9 in multi-variate design space. This study evaluates performance of chemical
10 foundation models that are pre-trained primarily with SMILES of small molecules,
11 in extrapolating learning from molecules to material design challenges across
12 multiple length scale in batteries. Ten prediction models are trained covering
13 molecular properties, formulations performance, and battery device measurement.
14 Material representations from several foundation models are compared and their
15 performance is benchmarked against conventional molecular representations such
16 as Morgan Fingerprints. The study further examines their capacity to generalize
17 to out-of-distribution cases by quantifying prediction errors for novel material
18 designs that differ substantially from the training data. Finally, interpretability of
19 the trained predictors is assessed by correlating actual outcomes and predictions
20 to the chemical moieties in the datasets, with the aim of enabling researchers to
21 interpret design rules in chemical space where model has high confidence.

22 1 Introduction

23 With evolving technologies and world economy demands, the field of material discovery has remained
24 strongly relevant. Recently, this field has acquired critical importance as new sustainable materials are
25 sought to overcome limitations of current material systems (1). Battery technologies are one strong
26 societally relevant area of research where the scope of known materials appears to be exhausted, and
27 new materials that can deliver high capacities, fast charging and longer cycle stability are continuously
28 sought to meet future demands (2; 3). Despite shifts in material research paradigms from slow, labor-
29 intensive experiments, to faster data-driven models (4; 1), it remains challenging to integrate new
30 materials in real world devices. This is due to several reasons: (i) most computational models
31 including simulations and machine learning (ML) can be used to determine intrinsic properties of
32 materials based on their chemical structure, but lack in extrapolating their outcome to meso or macro
33 scale phenomenon (5); (ii) device performance is governed by complex interactions among several
34 constituent materials, presenting vast multivariate design space difficult to screen or optimize (6); (iii)
35 limited data availability for extrinsic characteristics such as temperature and concentration dependence
36 of multi-constituent properties (7). While ML models accelerate several prediction, generative and
37 optimization problems in material science, the field continues to face challenges stemming from
38 opaque nature of the model’s decision making, impractical proposed chemical structures, scarcity of
39 quality datasets and inability to generalize out-of-distribution (OOD) (8).

40 Foundation models (*FM*s) have emerged as promising models to overcome some aforementioned
41 challenges of data scarcity and generalization. These are a class of large language models (LLMs),
42 that are pre-trained on a textual or multi-modal representations of materials in open-source databases
43 like PubChem and ZINC through self-supervised learning (9; 10). Studies have demonstrated that
44 embedding space of these transformer models segregates chemically relevant features of molecules
45 making them a suitable general-purpose tool for material science research. These base models can be
46 utilized to perform specific functions based on smaller labeled datasets with fine-tuning or transfer
47 learning (11). *FM*s are rapidly evolving, and their adoption in different application areas is on
48 the rise (12). Large portion of studies report their use in property prediction and inverse design of
49 small molecules or crystals (11). Prior studies also evaluate their scope in predicting performance
50 metrics for formulations (mixtures of more than two molecules in certain compositions) based on
51 electrolyte-performance experimental datasets curated from literature. Results demonstrate best
52 prediction accuracies from foundation models in comparison to other data-driven models (13; 14).
53 The research on representing advanced material systems such as formulations, composites and devices
54 to learning models is currently in nascent stages due to less understood chemical phenomenon and
55 lack of quality datasets. Prior studies on formulation datasets present strong evidence that foundation
56 models can extrapolate molecular features to multi-constituent properties.

57 In this work, we evaluate the capability of chemical *FM*s pre-trained with molecular representation
58 SMILES (15), to predict material properties and performance resulting from interplay of complex
59 chemical phenomenon at macroscale. We take battery electrolytes as an example where electrolyte
60 engineering has emerged as a promising approach to improve battery performance metrics such
61 as columbic efficiency (CE), cycle life and capacity. To achieve this, electrolytes are carefully
62 designed based on the individual properties of constituent molecules, their collective performance
63 as formulation and their compatibility with other battery components such as electrodes, separator
64 and current collector. Electrolyte Genome initiative in 2015 accelerated electrolyte discovery cycle
65 for new emerging battery chemistries by integrating computational workflows with experimentation
66 (16). High-throughput screening enabled selection of candidate molecules meeting threshold values
67 for HOMO-LUMO energy levels, toxicity and electrochemical stability. Once down-selection is
68 done, laborious experimentation is required to find their right combination for a functional electrolyte
69 formulation (17). Here, data availability is a primary roadblock in adoption of ML models since
70 public datasets are inconsistent and industrial datasets are propriety (18). Thus, models that can be
71 efficient with scarce datasets are desired in the domain.

72 We use *FM*s to map electrolyte formulations along with device variables to key performance
73 indicators at multiple length scale in batteries as illustrated in Figure 1. In particular,

- 74 • We target prediction of key properties that are considered in electrolyte discovery such as
75 molecular properties, formulation performance, manufacturability, surface contact char-
76 acteristics and device performance. *FM*s are used to generate input features for these

- 77 multivariate battery datasets and predictive capability is compared with standard molecular
 78 representations like Morgan Fingerprints (*MF*) (19).
- 79 • We evaluate out-of-distribution (OOD) capability of prediction models for multi-variate
 80 battery datasets.
 - 81 • Next, extrapolation capability of the models to new material designs is estimated based on
 82 the semantic similarity between train and test data. This presents a method to approximate
 83 errors in model predictions across new material landscape.
 - 84 • We investigate interpretability of *FM*-based predictors and evaluate their promise in infer-
 85 encing new material design rules.

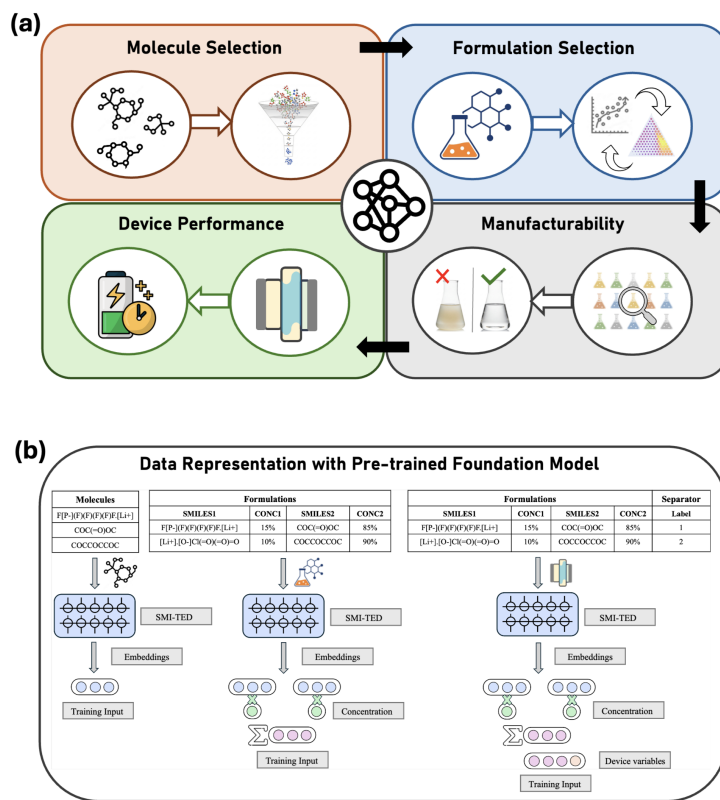


Figure 1: (a) Scheme illustrating electrolyte design problems at multiple scales. (b) Schematic summarizing the data representation for material design using pre-trained foundation models for molecules.

86 2 Electrolyte Datasets

87 Data availability is a major enabler for artificial intelligence (AI) workflows aiming for material
 88 discovery and design. To discover new material design that meets the performance goals, series
 89 of data driven predictors must be realized to allow material identification, characterization and
 90 optimization for achieving compatibility with the device. In this section, we describe battery datasets
 91 and performance indicators used across multiple length scale for electrolyte development. Most
 92 datasets are curated from literature and some are experimentally generated in the laboratory (see
 93 Supplementary Materials for details).

94 **Molecule screening:** Battery electrolytes can comprise of one or more organic solvent, and one or
 95 more salt, which facilitate Li⁺ ion transport between electrodes and electrode surface conditioning
 96 to prevent unwanted degrading side reactions. Each electrolyte component plays a crucial role in
 97 this ecosystem and is therefore selectively picked based on certain properties like HOMO-LUMO
 98 levels and redox potentials. While there is plethora of labeled dataset available in literature for

99 these properties (20; 21; 16), we use a data from a singular source to train and evaluate model’s
100 performance, i.e., D3TaLES, a database of DFT simulated properties of 40,000 organic molecules for
101 battery systems (21).

102 **Manufacturability:** Screened solvents and salts are combined in certain compositions to form
103 electrolyte formulations. These formulations must be completely miscible (or soluble) to enable ion
104 transport and manufacturing. We curate a heterogeneous dataset containing solubility information of
105 single salt-single solvent mixtures, single salt-multi solvent formulations, and multi salt- multi solvent
106 electrolytes, enabling development of a generalized model for electrolyte miscibility prediction. Refer
107 to A.1 for details on electrolyte solubility data generation. For inclusion of heterogeneous datasets,
108 we simplify approach to binary classification indicating insoluble (0) or soluble (1). The combined
109 3,300 dataset contained rich diversity of salts, solvents and electrolyte mixtures.

110 **Formulation property:** Another crucial property to consider during electrolyte design is ionic
111 conductivity (IC). The salts dissociated into ions within an electrolyte form solvation structures that
112 facilitate transport of charge between two electrodes and are responsible for battery’s charge-discharge
113 kinetics. For IC, we use 18,000 reported empirical values of electrolyte formulations at different
114 temperatures in published literature (7; 14). The dataset constitutes diverse set of solvents and salts.

115 **Surface contact characterization:** An electrolyte interfaces with multiple internal components
116 within a battery, including electrodes, separators, and current collectors. Consequently, optimizing
117 the surface interactions between the electrolyte formulation and various device constituents is crucial
118 for achieving peak performance. Traditionally, such evaluations have relied on the empirical expertise
119 of domain experts and expensive computational simulations. Nevertheless, data collected from
120 evaluation of one similar system can be used to automate future screening and assessment of
121 electrolytes. We use one such in-house generated empirical dataset of 119 electrolyte formulations
122 and their contact angle on four different separators to predict surface contact angle of electrolytes
123 (see A.2 for experimental details).

124 **Device performance:** The ultimate objective of developing a new battery electrolyte formulation is
125 to achieve superior performance metrics, such as enhanced capacity, Coulombic Efficiency (CE), and
126 cycle life. The public dissemination of such data is often limited, as its relevance is typically highly
127 specific to a particular device configuration, thereby precluding its full adherence to FAIR (Findable,
128 Accessible, Interoperable, and Reusable) data principles. To address this challenge, we leverage three
129 distinct datasets from previous publications. The first dataset, derived from a study by Kim et al. (3),
130 examines the relationship between electrolyte composition and CE across 150 datapoints. A second
131 dataset containing 125 electrolytes, originally reported by Sharma et al. (6), explores the influence of
132 electrolyte formulation on the specific capacity of a LiI conversion battery. Finally, the third dataset
133 constituting 91 datapoints focuses on capacity metric for an interhalogen conversion (Li-ICl) battery,
134 incorporating variations in cathode loading, separator type, and electrolyte compositions with fixed
135 chemicals (18).

136 3 Foundation Models for Material Representation

137 Presently, there is a plethora of pre-trained transformer models in literature that are used for specific
138 downstream scientific tasks (22; 10; 23; 24; 25). Particularly in the domain of chemistry and material
139 science, sequence prediction, molecular property prediction and chemical description generation are
140 a few tasks that are used in benchmarking *FM*. In this work, we aim to evaluate scope of *FM*s
141 pre-trained on molecular representations in addressing material design challenges across multiple
142 length scale in batteries. Comparative analyses were performed across multiple *FM* to elucidate the
143 extent to which model performance and generalization behaviors are influenced by differences in
144 pretraining modalities.

145 **SMI-TED:** SMI-TED (SMILES Transformer Encoder Decoder) is an open-source chemical *FM*
146 developed by IBM Research (10). This model has acquired a deep understanding of molecular
147 structural representations through self-supervised pre-training on a vast dataset containing string
148 representation (SMILES) of 91 million molecules, corresponding to 4 billion molecular tokens. Model
149 has been previously validated to surpass the performance of conventional data-driven alternatives in
150 downstream tasks.

151 **MolT5:** MolT5 (Molecular T5) is another open sourced chemical *FM* that is pre-trained with 100
152 million SMILES along with 33,000 natural language description of molecules (25). By correlat-
153 ing SMILES sequences to textual description of functionalities, the model has shown remarkable
154 capabilities in manipulating molecules for discovery tasks.

155 **Galactica:** Galactica is a large language model developed for general scientific tasks by Meta AI
156 (24). The model is trained on large corpus of scientific literature, natural sequences of proteins and 2
157 million chemical strings (SMILES). The inclusion of broad data makes is a reliable model for general
158 scientific tasks such as equation probing, citation prediction, reasoning, etc.

159 **GraphMVP:** GraphMVP is a graphs based pre-trained model that formulates a multi-view self-
160 supervised learning, integrating both 2D molecular graphs and rich 3D spatial arrangements of atoms
161 (26). The GraphMVP learning framework allows its encoder to integrate topological and geometric
162 information within a unified embedding space. It is worth noting that GraphMVP uses much smaller
163 graph/conformer datasets in representation learning.

164 **Morgan Fingerprints:** As a benchmark, *MF* are employed as an established molecular descriptor
165 (19). *MF* are highly effective for predicting molecular properties in ML models because they
166 efficiently capture the substructural features of a molecule (27). By representing a molecule as a
167 fixed-length binary vector, they encode the presence or absence of specific circular substructures and
168 each atom’s chemical environments. The resulting numerical representation is both computationally
169 efficient and chemically intuitive, making it an ideal input for various learning algorithms, which can
170 then identify complex patterns and relationships that are predictive of a molecule’s behavior.

171 For downstream tasks, transfer learning approach is adopted to retain chemical information from the
172 pre-trained model as molecular embeddings, and map these to the output label using a regressor model
173 such as feed forward neural networks (NN). It is noted that fine-tuning the pre-trained *FM* containing
174 several million parameters with labeled datasets can be computationally expensive. Furthermore,
175 fine-tuning current state-of-the-art *FM* is not expandable to the string representations of formulations
176 used in ref(14) as these are vastly different from the molecule representations models were pre-trained
177 on. Meanwhile, transfer learning approach is relatively robust and deliver consistently reliable results
178 (see Table S1). Therefore, embeddings from the *FMs* and *MF* are used to represent individual
179 molecules in the battery datasets. Derived molecular embeddings are aggregated into a system
180 representation based on their composition, and additional design variables in the dataset such as
181 separator, temperature and cathode loading (indicated in Figure 1b). Details of feature engineering
182 for appropriate representation of molecules, formulations and devices are described in A.3. For each
183 prediction task, feed forward neural network (NN) architectures are optimized and trained using
184 *FM*-derived and aggregated features (described in A.4). NNs were trained using five independent
185 80%-20% train-test splits, and prediction errors were quantified using the mean absolute error (MAE)
186 metric.

187 4 Results and Discussion

188 4.1 Model performance

189 We use *FMs* that recognize SMILES modality for training electrolyte design predictors due to ease of
190 chemical data representation and their demonstrated best performance in predicting molecular proper-
191 ties in several benchmark datasets (10). Prediction results for 10 battery datasets are summarized in
192 Table 1 for *FMs* and *MF*. Tabulated are the average MAE across 5 random train-test splits for all
193 models. Results show that SMI-TED and MolT5 based representations outperform *MF* in 7 out of
194 10 datasets. Meanwhile predictive capability of Galactica and GraphMVP is observed to be the lowest
195 in all 10 datasets. Particularly for molecular properties, where several prior studies have backed that
196 2048 bits of *MF* are more predictive than domain-intuitive features (27), results in Table1 indicate
197 SMI-TED outperforms *MF*. SMI-TED demonstrates notable computational efficiency despite using
198 significantly smaller feature vector size (768). This efficacy of SMI-TED embeddings testifies that
199 learnt representations encode more comprehensive set of structural features that are meaningful and
200 comprehensive.

201 In the context of more complex systems, such as formulations, we observed a systematic divergence
202 in model performance based on data size. On datasets characterized by a large volume of data, such
203 as solubility (3300 data) and IC (18,000 data), *MF* outperform all *FM* in the present evaluation,

204 categorizing miscible and immiscible electrolytes with 93.77% accuracy, and predicting log IC with
205 MAE 0.0629, surpassing previously best reported results in ref (14). This outcome is consistent with
206 the design of conventional ML methods that are optimized for large-scale data problems. *MF*'s
207 enhanced performance on these datasets suggests that the fundamental properties like IC and solubility
208 are more contingent on specific functional groups in the system that are captured precisely by *MF*.
209 This finding presents a critical consideration for the future development of foundation models.

210 SMI-TED and MolT5 demonstrated clear and consistent advantage over *MF* in low data regimes (100
211 to 200 data points), achieving superior predictive accuracy and robustness across these challenging
212 multiscale problems. Particularly MolT5, having pre-trained on largest corpus of molecular data (100
213 Million SMILES), has the lowest prediction errors for contact angle (MAE 12.944 Degrees) and LiI
214 capacity (MAE 22.408 mAh/g) datasets, and is second to *MF* for solubility (93.65% Accuracy) and
215 IC (log IC MAE 0.0722) prediction. SMI-TED demonstrates next best predictive capability among
216 *FMs*, reporting low prediction errors for all formulation datasets and outperforming all models for
217 CE dataset (6; 3). These results highlight applicability of *FM* pretrained with molecules alone to
218 multi-variate material design problems. Possible interpretation is that macroscale outcomes, such as
219 electrolyte performance, are dictated by hierarchical interactions between chemical moieties. Ion
220 aggregates and solvation substructures are examples of chemical moiety interactions responsible for
221 charge-discharge kinetics in battery electrolytes. Models such as MolT5 and SMI-TED successfully
222 predicts these macroscale outcomes due to having rich chemical vocabulary comprising of thousands
223 of unique chemical tokens or moieties as reported in ref(10). Hence, latent space of SMILES-based
224 *FM* is enriched with basic understanding of the chemical space formed by the combinations of
225 chemical moieties in molecules (10). The downstream training utilizing aggregated formulation
226 embeddings vs performance label is useful to correlate chemical moieties and compositions to the
227 label, enabling multi-scale learning (see Figure 2). This knowledge transfer is particularly useful in
228 low data regimes. Li-ICl Capacity data is a singular instance where *MF* outperforms *FMs* despite
229 low data regime, highlighting *FMs* are likely not suitable for datasets lacking chemical variability.

230 Results from MolT5 present additional interesting observations on multi-modal pre-training. Latent
231 space of MolT5 is augmented with semantic understanding of molecular string representation,
232 correlating molecule structures to specific functions (25). In Table 1, advantages of pretraining with
233 multi-modal datasets is noted in multi-variate battery datasets but not in molecular datasets. Despite
234 pre-training on largest SMILES corpus, predictive capability of MolT5 model is lower than SMI-TED
235 for molecular properties, likely due to noted functional biases and scarcity of natural language
236 datasets used during model development(25). Regardless, good predictive performance on multi-
237 variate datasets underscore the critical importance of incorporating multi-modal data representations
238 during the pretraining, enabling model to learn complex inter-dependencies and semantic nuances
239 across datasets.

240 Poor performance of Galactica in predicting material properties underline limitations of high gener-
241 ality. Despite training on large corpus of scientific knowledge and 2 Million SMILES, model
242 lacks sufficient specificity required to capture critical domain-relevant features. In lieu, GraphMVP
243 also shows poor predictive power despite high specialization in molecular geometries. The model
244 captures the 3-D topological and geometric features of molecules but lacks sufficient representational
245 capacity to resolve finer substructural moieties and their inter-dependencies. Ultimately, the choice
246 of representation is critical and must be determined by the nature of downstream task, quantity and
247 the quality of the labeled dataset.

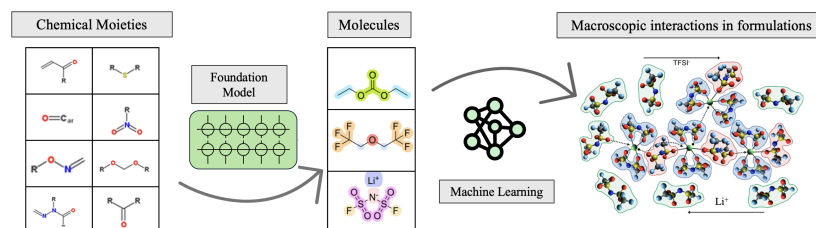


Figure 2: Multi-step training capturing complex chemical interactions at multiple length scale.

Table 1: Average mean absolute error (MAE) and prediction accuracy (%) for the battery datasets using embeddings from foundation models

Model ↓	Oxidation	Reduction	HOMO	LUMO	Solubility	IC	Contact Angle	LiI Capacity	CE	Li-ICI Capacity
MAE Units →	<i>eV</i>	<i>eV</i>	<i>eV</i>	<i>eV</i>	Accuracy %	<i>Log</i>	<i>Degrees</i>	<i>mAh/g</i>	<i>Log</i>	<i>mAh/g</i>
SMI-TED	0.2559	0.5825	0.4405	0.3663	93.11	0.0910	16.243	22.449	0.185	47.93
MolT5	0.2679	1.7375	0.4451	0.3836	93.65	0.0722	12.944	22.408	0.188	37.57
Galactica	0.2714	0.7134	0.4802	0.4283	93.05	0.1035	23.982	25.011	0.225	39.570
GraphMVP	0.3355	0.6586	0.4987	0.4432	91.17	0.0939	22.099	29.051	0.209	42.451
MF	0.2594	0.5854	0.4580	0.3746	93.77	0.0629	17.815	28.990	0.223	32.24

248 4.2 Quantifying out-of-distribution performance

249 Formulations present multi-variate design space with infinite possibilities emerging from several
 250 million known compounds, their inestimable potential combinations, and composition variations.
 251 Given this, electrolyte design discovery becomes inherently an OOD problem as novel formulations
 252 will most likely be in unseen or unfamiliar data. Thus, evaluating OOD performance is crucial
 253 for ensuring the reliability and robustness of models. One can define OOD based on divergence
 254 between train-test sets with respect to either input distribution (chemical and composition space) or
 255 output distribution (property values). Presented OOD evaluation of *FM*s for formulation and device
 256 performance datasets spans both input and output distributions.

257 First, we start with most accepted OOD evaluation based on output distribution (28). We separate test
 258 sets based on tail ends of numerical outcome distribution, for instance, lower and upper end values of
 259 ionic conductivity, capacity, contact angle, etc. Tail-end distributions used as tests in 5 electrolyte
 260 regression datasets are highlighted in A.5. This distribution estimates extrapolation capabilities of
 261 the models beyond the training data. Results of OOD predictions are presented in Table 2 along
 262 with prediction uncertainty observed across 3 predictions. Both SMI-TED and MolT5 demonstrate
 263 best OOD prediction with each having lowest MAE in 2 out of 5 datasets. Both models also had
 264 high consistency in predicted outcomes as indicated by low uncertainty. Overall extrapolation across
 265 outcome values is promising for electrolyte datasets except for Li-ICI Capacity dataset where models
 266 perform poorly as seen in previous section.

Table 2: Mean absolute error (MAE) for out-of-distribution predictions using foundation models and Morgan Fingerprints

Model ↓	CE	Contact Angle	LiI Capacity	IC	Li-ICI Capacity
MAE Units →	<i>Log</i>	<i>Degrees</i>	<i>mAh/g</i>	<i>Log</i>	<i>mAh/g</i>
SMI-TED	0.0548 ± 0.04	13.5216 ± 0.41	27.128 ± 0.70	0.1938 ± 0.01	109.21 ± 0.95
MolT5	0.0819 ± 0.00	14.0539 ± 0.98	31.2229 ± 1.61	0.1669 ± 0.01	108.2197 ± 0.93
Galactica	0.4635 ± 0.39	31.4742 ± 0.82	28.2692 ± 14.38	0.2262 ± 0.08	110.391 ± 1.50
GraphMVP	2.7758 ± 2.36	34.8031 ± 1.88	7.9974 ± 4.17	0.7429 ± 0.04	108.6611 ± 0.03
MF	0.1295 ± 0.05	19.3304 ± 1.26	29.5058 ± 2.22	0.1717 ± 0.03	114.3028 ± 31.07

267 Next, ML models frequently show poor transferability across chemical spaces and fall short in
 268 predicting properties for materials outside their training scope (29). Generalizable base models like
 269 *FM* have seen increased adoption in the community for these reasons (29). Unlike small molecules,
 270 where property can be traced to substructures and chemical motifs (10), cause-effect in formulations-
 271 like materials are more complex and intertwined in multi-variate dynamic inter-dependencies (14).
 272 Therefore, the boundaries of OOD for dynamic multi-variate chemical space is needed to be explored
 273 in a focused study. In present study, we use chemical similarity as a metric for characterizing OOD
 274 based on inputs. A chemical similarity score is employed as an approximation for how close test data
 275 is to training data in model’s latent space, and is estimated by calculating maximum of average cosine
 276 similarity (normalized) of each test datapoint with all training samples. Upon evaluating the chemical
 277 similarity between embeddings of train-test sets for tail-end OOD evaluation in Table S3, we observe
 278 there is an inverse trend between chemical similarity of OOD train-test sets and prediction MAE from
 279 the models, suggesting model prediction errors are high for chemically disparate test sets. These
 280 results confirm chemical similarity can be a reliable metric to determine distance between test and
 281 train sets in model’s latent space and characterize OOD.

282 This trend paves the way to ascertain reliability of a model when extrapolating to unexplored regions
 283 of the materials design space. By error estimation, we can systematically pinpoint regions where
 284 model lacks predictive capability, facilitating intelligent allocation of resources toward targeted
 285 experimental validation and data enrichment. We create several subsets of train-test data for battery

286 across different length scale based on their relative distance in latent space of SMI-TED, given its
 287 reliable performance in both molecules and macroscale outcomes. These subsets were carefully
 288 curated to represent a different testing scenario than the ones used in the tail-end OOD evaluation
 289 such as distinct constituent count and chemicals. Relationship between semantic similarity between
 290 the input embeddings of train-test distributions (in red) across datasets is compared with prediction
 291 MAE for the respective train-test subset (in blue) in Figure 3. Trends confirm an inverse relationship
 292 between prediction MAE and semantic proximity of test data to the training samples, yielding a
 293 linear relationship $MAE = m \cdot Similarity + c$ that estimates the approximate MAE of model
 294 predictions on new data points by quantifying their *Similarity* to the model’s training data. The
 295 slope (m) and intercept (c) for analyzed datasets are presented in Table S4. This approach enables
 296 systematic assessment of prediction uncertainty and confidence for new data, thereby supporting
 297 efficient screening in materials design and discovery.

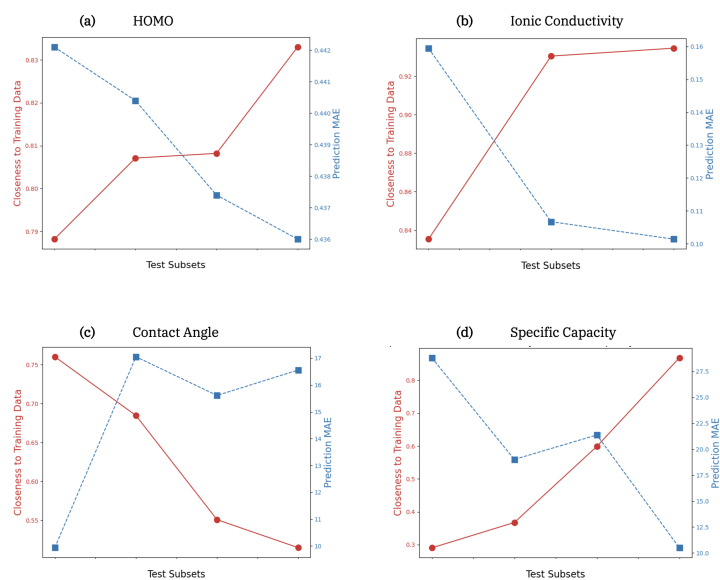


Figure 3: Relationship between prediction MAE (in blue) and chemical similarity (in red) between train and test datasets.

298 4.3 Interpretability

299 A widely embraced strategy in materials discovery involves interpreting chemical data into useful
 300 knowledge and chemical insights, uncovering conclusive design rules and trends for decision making
 301 (30; 31). The efficacy of this approach is maximized when it leverages accurate empirical data
 302 or highly reliable model-generated outputs spanning the intended design landscape. However,
 303 interpretability is frequently hindered by the intrinsic opacity of AI models, which predominantly
 304 operate as “black boxes” with internal mechanisms that remain inaccessible to researchers. This
 305 challenge is further exacerbated as training pipeline grow in complexity, for instance, input features
 306 are derived from transformer model and post processed before the training (18). Quantifying model
 307 uncertainty in new material regions can facilitate users in identifying scope of the model. However,
 308 application of these models to uncover material design rules for interpretability remains a persistent
 309 challenge.

310 We propose a method to evaluate interpretability of *FM* derived predictors by investigating correlation
 311 of performance outcomes with chemical moieties in the datasets and compare trends in train and
 312 test subsets. First, a list of several potential chemical substructures and their SMARTS (SMILES
 313 Arbitrary Target Specification) string is devised (32). Over 550 chemical substructures are defined
 314 including general and specific moieties. For instance, amine is a general functional group of material
 315 containing Nitrogen atom with lone pair of electrons, and specific derivatives for the same include
 316 aromatic amine, heterocyclic amine, tertiary amine etc. Chemical moieties in molecules are identified
 317 by matching SMARTS and presence of every moiety is indicated by a bit in a fixed length vector.
 318 This vector is taken as molecular fingerprints and aggregated for constituents in each formulation by

319 composition scaling and addition to represent concentration of each chemical moiety in a formulation.
 320 We adopt Spearman’s correlation coefficient (SCC) (33) to determine strength and direction of
 321 monotonic relationship between chemical moieties in the dataset and the outcome performance. The
 322 analysis provides meaningful insights towards the positive or negative influence of a chemical moiety
 323 in the formulation towards the outcome. Analysis is performed for data used in training and test set to
 324 correlate moieties to actual outcomes. Simultaneously, the analysis is also extended to the outcomes
 325 predicted by the models based on SMI-TED representation for the very same test set. Figure 4
 326 illustrates these correlations in three formulation datasets CE, LiI capacity and IC.

327 Comparison of correlation analysis for model prediction outcomes and actual performance within
 328 test sets is meant to demonstrate the capability of model in deriving sound chemical insights across
 329 unseen datapoints. Particularly in Figure 4, examples highlighted in green illustrate cases where the
 330 correlations in the training and test datasets were opposite, and the model correctly predicted the
 331 opposing trends. Instances highlighted in yellow represent scenarios where the model accurately
 332 identified chemical trends for the outcome, despite these trends being absent from the training data.
 333 Cases highlighted in pink show perfect alignment among all three correlations. The remaining
 334 instances in white indicate correlations that the foundation model misinterpreted. This analysis
 335 reveals the chemical insights misunderstood by the model and allows users to selectively apply these
 336 models for design interpretation and discovery within a chemical space where confidence is justified.

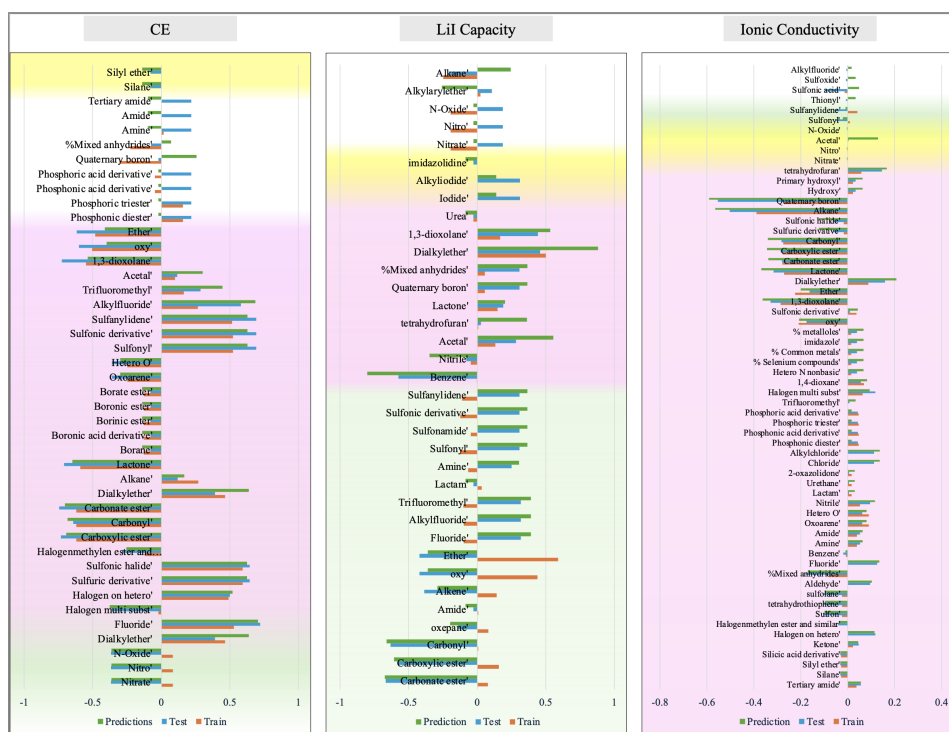


Figure 4: Correlation of chemical functional groups in formulations with performance in train (orange) - test (blue) dataset, compared with correlation to the predicted outcomes (green) in test data.

337 5 Conclusion

338 In this work, we evaluate the scope of foundation models in addressing material design challenges
 339 across multiple length scale in batteries: molecules, formulations and device. Multiple foundation
 340 models are used to derive multi-variate representations of datasets by combining molecular representations
 341 with other variables such as compositions, temperature, electrode and separator variations.
 342 Results show *FM*s pre-trained with large corpus of SMILES modality, such as SMI-TED and MolT5,
 343 can be used to extrapolate learning from moiety-level interactions to macroscopic outcomes like
 344 specific capacity, surface characteristics, and battery performance using scarce datasets. These models

345 are particularly useful in low data regimes where conventional molecular representations such as
346 Morgan Fingerprints are found to be limiting. It is also observed that pre-training on multi-modal data
347 representations has the scope to achieve superior performance in multi-variate material design space.
348 The study also presents a method to analyze model's ability to generalize out-of-distribution and
349 quantify model prediction errors across new material designs based on chemical similarity between
350 train-test sets. SMILES-based models demonstrated reliable out-of-distribution performance trends.
351 However, it is noted that out-of-distribution criterion for dynamic multi-variate chemical space
352 needs further comprehensive investigation. Lastly, we demonstrate an approach to identify chemical
353 space where model confidence is high by correlating actual outcomes and predicted outcomes to the
354 chemical moieties in the datasets. The approach allows dependable material design interpretation
355 from the model for discovery tasks.

356 References

- 357 [1] A. Merchant, S. Batzner, S. S. Schoenholz, M. Aykol, G. Cheon, and E. D. Cubuk, "Scaling deep learning
358 for materials discovery," *Nature*, vol. 624, no. 7990, pp. 80–85, 2023.
- 359 [2] J. Datta, A. Nadimpally, N. Koratkar, and D. Datta, "Generative ai for discovering porous oxide materials
360 for next-generation energy storage," *Cell Reports Physical Science*, 2025.
- 361 [3] S. C. Kim, S. T. Oyakhire, C. Athanitis, J. Wang, Z. Zhang, W. Zhang, D. T. Boyle, M. S. Kim, Z. Yu,
362 X. Gao *et al.*, "Data-driven electrolyte design for lithium metal anodes," *Proceedings of the National
363 Academy of Sciences*, vol. 120, no. 10, p. e2214357120, 2023.
- 364 [4] E. O. Pyzer-Knapp, J. W. Pitera, P. W. Staar, S. Takeda, T. Laino, D. P. Sanders, J. Sexton, J. R. Smith, and
365 A. Curioni, "Accelerating materials discovery using artificial intelligence, high performance computing
366 and robotics," *npj Computational Materials*, vol. 8, no. 1, p. 84, 2022.
- 367 [5] J. Li, K. Lim, H. Yang, Z. Ren, S. Raghavan, P.-Y. Chen, T. Buonassisi, and X. Wang, "Ai applications
368 through the whole life cycle of material discovery," *Matter*, vol. 3, no. 2, pp. 393–432, 2020.
- 369 [6] V. Sharma, M. Giammona, D. Zubarev, A. Tek, K. Nugyuen, L. Sundberg, D. Congiu, and Y.-H. La,
370 "Formulation graphs for mapping structure-composition of battery electrolytes to device performance,"
371 *Journal of Chemical Information and Modeling*, vol. 63, no. 22, pp. 6998–7010, 2023, PMID: 37948621.
372 [Online]. Available: <https://doi.org/10.1021/acs.jcim.3c01030>
- 373 [7] P. de Blasio, J. Elsborg, T. Vegge, E. Flores, and A. Bhowmik, "Calisol-23: Experimental electrolyte
374 conductivity data for various li-salts and solvent combinations," *Scientific Data*, vol. 11, no. 1, p. 750,
375 2024.
- 376 [8] A. K. Cheetham and R. Seshadri, "Artificial intelligence driving materials discovery? perspective on
377 the article: Scaling deep learning for materials discovery," *Chemistry of Materials*, vol. 36, no. 8, pp.
378 3490–3495, 2024.
- 379 [9] J. Ross, B. Belgodere, V. Chenthamarakshan, I. Padhi, Y. Mroueh, and P. Das, "Large-scale chemical
380 language representations capture molecular structure and properties," *Nature Machine Intelligence*, vol. 4,
381 no. 12, pp. 1256–1264, 2022.
- 382 [10] E. Soares, E. Vital Brazil, V. Shirasuna, D. Zubarev, R. Cerqueira, and K. Schmidt, "An open-source family
383 of large encoder-decoder foundation models for chemistry," *Communications Chemistry*, vol. 8, no. 1, p.
384 193, 2025.
- 385 [11] J. Choi, G. Nam, J. Choi, and Y. Jung, "A perspective on foundation models in chemistry," *JACS Au*, vol. 5,
386 no. 4, pp. 1499–1518, 2025.
- 387 [12] E. O. Pyzer-Knapp, M. Manica, P. Staar, L. Morin, P. Ruch, T. Laino, J. R. Smith, and A. Curioni,
388 "Foundation models for materials discovery—current state and future directions," *Npj Computational
389 Materials*, vol. 11, no. 1, p. 61, 2025.
- 390 [13] I. Priyadarsini, V. Sharma, S. Takeda, A. Kishimoto, L. Hamada, and H. Shinohara, "Improving perfor-
391 mance prediction of electrolyte formulations with transformer-based molecular representation model," in
392 *ICML'24 Workshop ML for Life and Material Science: From Theory to Industry Applications*.
- 393 [14] M. Zohair, V. Sharma, E. A. Soares, K. Nguyen, M. Giammona, L. Sundberg, A. Tek, E. A. Vital, and
394 Y.-H. La, "Chemical foundation model-guided design of high ionic conductivity electrolyte formulations,"
395 *npj Computational Materials*, vol. 11, no. 1, p. 283, 2025.

- 396 [15] D. Weininger, "Smiles, a chemical language and information system. 1. introduction to methodology and
397 encoding rules," *Journal of Chemical Information and Computer Sciences*, vol. 28, pp. 31–36, 1988.
- 398 [16] L. Cheng, R. S. Assary, X. Qu, A. Jain, S. P. Ong, N. N. Rajput, K. Persson, and L. A. Curtiss, "Accelerating
399 electrolyte discovery for energy storage with high-throughput screening," *The journal of physical chemistry
400 letters*, vol. 6, no. 2, pp. 283–291, 2015.
- 401 [17] A. Benayad, D. Diddens, A. Heuer, A. N. Krishnamoorthy, M. Maiti, F. L. Cras, M. Legallais, F. Rahmanian,
402 Y. Shin, H. Stein *et al.*, "High-throughput experimentation and computational freeway lanes for accelerated
403 battery electrolyte and interface development research," *Advanced Energy Materials*, vol. 12, no. 17, p.
404 2102678, 2022.
- 405 [18] V. Sharma, A. Tek, K. Nguyen, M. Giammona, M. Zohair, L. Sundberg, and Y.-H. La, "Improving
406 electrolyte performance for target cathode loading using an interpretable data-driven approach," *Cell
407 Reports Physical Science*, vol. 6, no. 1, 2025.
- 408 [19] D. Rogers and M. Hahn, "Extended-connectivity fingerprints," *Journal of chemical information and
409 modeling*, vol. 50, no. 5, pp. 742–754, 2010.
- 410 [20] R. Ramakrishnan, P. O. Dral, M. Rupp, and O. A. Von Lilienfeld, "Quantum chemistry structures and
411 properties of 134 kilo molecules," *Scientific data*, vol. 1, no. 1, pp. 1–7, 2014.
- 412 [21] R. Duke, V. Bhat, P. Sornberger, S. A. Odom, and C. Risko, "Towards a comprehensive data infrastructure
413 for redox-active organic molecules targeting non-aqueous redox flow batteries," *Digital Discovery*, vol. 2,
414 no. 4, pp. 1152–1162, 2023.
- 415 [22] J. Pan, "Large language model for molecular chemistry," *Nature Computational Science*, vol. 3, no. 1, pp.
416 5–5, 2023.
- 417 [23] J. Ross, B. Belgodere, S. C. Hoffman, V. Chenthamarakshan, J. Navratil, Y. Mroueh, and P. Das, "Gp-
418 molformer: A foundation model for molecular generation," *Digital Discovery*, 2025.
- 419 [24] R. Taylor, M. Kardas, G. Cucurull, T. Scialom, A. Hartshorn, E. Saravia, A. Poulton, V. Kerkez, and
420 R. Stojnic, "Galactica: A large language model for science," *arXiv preprint arXiv:2211.09085*, 2022.
- 421 [25] C. Edwards, T. Lai, K. Ros, G. Honke, K. Cho, and H. Ji, "Translation between molecules and natural
422 language," *arXiv preprint arXiv:2204.11817*, 2022.
- 423 [26] S. Liu, H. Wang, W. Liu, J. Lasenby, H. Guo, and J. Tang, "Pre-training molecular graph representation
424 with 3d geometry," *arXiv preprint arXiv:2110.07728*, 2021.
- 425 [27] H. Zhou and J. Skolnick, "Utility of the morgan fingerprint in structure-based virtual ligand screening,"
426 *The Journal of Physical Chemistry B*, vol. 128, no. 22, pp. 5363–5370, 2024.
- 427 [28] E. R. Antoniuk, S. Zaman, T. Ben-Nun, P. Li, J. Diffenderfer, B. Demirci, O. Smolenski, T. Hsu, A. M.
428 Hiszpanski, K. Chiu *et al.*, "Boom: Benchmarking out-of-distribution molecular property predictions of
429 machine learning models," *arXiv preprint arXiv:2505.01912*, 2025.
- 430 [29] M. A. Skinnider, R. G. Stacey, D. S. Wishart, and L. J. Foster, "Chemical language models enable
431 navigation in sparsely populated chemical space," *Nature Machine Intelligence*, vol. 3, no. 9, pp. 759–770,
432 2021.
- 433 [30] H. Choubisa, P. Todorović, J. M. Pina, D. H. Parmar, Z. Li, O. Voznyy, I. Tamblyn, and E. H. Sargent,
434 "Interpretable discovery of semiconductors with machine learning," *NPJ Computational Materials*, vol. 9,
435 no. 1, p. 117, 2023.
- 436 [31] J. Dean, M. Scheffler, T. A. Purcell, S. V. Barabash, R. Bhowmik, and T. Bazhurov, "Interpretable machine
437 learning for materials design," *Journal of Materials Research*, vol. 38, no. 20, pp. 4477–4496, 2023.
- 438 [32] X. Liu, S. Swaminathan, D. Zubarev, B. Ransom, N. Park, K. Schmidt, and H. Zhao, "Accfg: Accurate
439 functional group extraction and molecular structure comparison," *Journal of Chemical Information and
440 Modeling*, 2025.
- 441 [33] P. Schober, C. Boer, and L. A. Schwarte, "Correlation coefficients: appropriate use and interpretation,"
442 *Anesthesia & analgesia*, vol. 126, no. 5, pp. 1763–1768, 2018.
- 443 [34] I. Priyadarsini, S. Takeda, L. Hamada, E. V. Brazil, E. Soares, and H. Shinohara, "Self-bart: A transformer-
444 based molecular representation model using selfies," *arXiv preprint arXiv:2410.12348*, 2024.
- 445 [35] H. Zhang, T. Lai, J. Chen, A. Manthiram, J. M. Rondinelli, and W. Chen, "Learning molecular mixture
446 property using chemistry-aware graph neural network," *PRX Energy*, vol. 3, no. 2, p. 023006, 2024.

447 A Supplementary Material

448 A.1 Solubility Data Collection

449 Complete electrolyte miscibility is desired in batteries for manufacturing to ensure that the electrolyte composition is consistent batch to batch and devoid of any phase separation for uniformity in battery performance at production scale. Therefore, it is essential to identify potentially miscible formulations from the vast combinatorial design space. Heterogeneous solubility dataset is generated through experimentation:

453 **Single salt- single solvent solubility assessment:** A dataset of binary system containing single salt and a single organic solvent was collected experimentally in the laboratory. The dataset spans five most popular electrolyte salts, LiNO₃, LiFSI, LiBOB, LiFOB, and LiPF₆, and up to fifty organic solvents. The experiments were conducted in an inert glovebox (Argon, < 0.1 ppm H₂O and O₂) and all salts were dried on a hotplate at 150 °C, except for LiFSI and LiPF₆, which were used as received due to their lower thermal stability. Solvents were dried over 3Å molecular sieves for at least 24 hours prior to use. An upper salt concentration limit of 2M was set during the data collection. Salts were weighed to make 2M solution and the respective organic solvent was then added to decrease the concentration by a 0.25M interval until the solutions were visually clear without any precipitation or undissolved materials. The salt-solvent combination was considered insoluble if the solution was not clear at 0.25M concentration.

463 **Single salt- Multi solvent solubility assessment:** The dataset has measurement of the highest molar concentration of single salt dissolved in mixture of organic solvents. The data was curated during the development of electrolyte for our prior study where four salts and four solvents were shortlisted for lithium metal battery electrolyte (18). The four salts, LiCl, LiNO₃, LiTFSI and LiBOB, are individually dissolved in solvent formulations containing different compositions of ethylene carbonate, Tetraglyme, 1,3-Dimethyl-2-imidazolidinone and 1,3-Dioxolane. The solubility measurements were made as per the method described above.

469 **Multi salt-multi solvent solubility assessment:** Conventionally, functioning and high-performing electrolytes are published in literature (3; 18; 6). We also share a few "failed" non-miscible electrolytes in our previous works (18; 14). We curated 300 electrolyte formulations from these studies. Simplification of solubility metric to (0) or (1) enabled inclusion and test across widespread electrolyte dataset. The combined dataset contained rich diversity of salts, solvents and electrolyte mixtures.

474 **Post processing:** The solubility of single salt- single solvent pairs and single salt- multi solvent formulations were measured in terms of highest soluble molarity of the salt. To further add context to the solute molarity noted as metric in empirical dataset, data augmentation was done to interpolate solubility of target salt in each respective solvent system to include soluble(1) datapoints below highest soluble molarity, and insoluble(0) datapoints above recorded metric until the tested molarity. Next, the constituent moles in each formulation system were converted to molar percentage (mole%). Post data processing, there are 3300 electrolyte formulation vs solubility data that is used in the study.

481 A.2 Contact Angle Measurement Experiments

482 Electrolyte uptake by separator is an important parameter that determines ion transport and electrolyte performance. There are several separators in the commercial market based on constitution such as polymer and quartz. Within a single category like polymer separators, vast variations can be noted based in changes in polymer monomers and ratios. Electrolyte formulations are prepared inside an Ar-filled glove box (<1 ppm O₂, <1 ppm H₂O). Prior to mixing, solvents that are liquid at room temperature are dried using molecular sieves (Millipore Sigma, 3 Å) and salts are dried on a hot plate at 100 °C. Electrolytes are mixed for 24 hrs prior to contact angle measurement. Contact angle measurements were conducted using an OCA video-based contact angle goniometer (FDS Future Digital Scientific Corporation) employing the sessile drop technique. Prior to measurement, the separator was carefully placed on a flat silicon wafer substrate to ensure a uniform surface. A 2L droplet of electrolyte was then dispensed onto the separator surface and allowed to equilibrate for 800ms. Image analysis was performed on a selected video frame by manually defining the baseline and applying an ellipse-fitting algorithm to achieve optimal conformity to the droplet profile. The reported static contact angles represent the average of 3–5 independent measurements. All procedures were carried out with minimal air exposure to preserve the integrity of the electrolyte and ensure reproducibility. A dataset of 119 experiments is created using the electrolyte constituents, their respective concentrations, the experimentally measured contact angle, and a separator label. There are four different Celgard separators in the dataset, identified by unique label (1-3).

499 A.3 Feature engineering

500 The application of data-driven models in material systems rely on the correct transformation of system into
501 a numerical representation suitable for mathematical operations. Accordingly, the intricate description of a
502 battery’s formulation, which includes the identity of constituent molecules, their composition, and additional
503 configuration parameters, must be systematically converted into a relevant numerical descriptor. For this
504 purpose, pretrained *FMs* are used to acquire molecular representations which are then transformed to represent
505 multi-scale systems as described below:

506 **Molecules:** *FMs* are used to derive numerical embeddings of molecules present in the target datasets similar to
507 previous studies (10; 34).

508 **Formulations:** Three formulation datasets including solubility, CE and LiI battery capacity map electrolyte
509 formulations to the outcome. Formulation inputs constitute multiple constituents per datapoint and their
510 respective composition as mole percent (*mol%*) in the mixture. Here, constituent molecules are transformed to
511 *FM* embeddings, and are then scaled based on their *mol%* in the formulation to indicate their activity within the
512 system. The scaled embeddings are aggregated to form a formulation descriptor by addition as also summarized
513 in Figure 1. There are more than one method to aggregate formulation descriptor (18; 35; 13). Each method has
514 its own merit and preferred use. We observe that scaled addition is most convenient aggregation as the resultant
515 formulation descriptor size is invariant to the formulation constituent count. IC dataset contains temperature as
516 an additional extrinsic variable that is concatenated with the formulation descriptor for training.

517 **Surface contact characterization:** In present study, contact angle of electrolyte on several polymer-based
518 separators are measured to assess their compatibility. For best representation, a *FM* for polymer can be
519 used. However, since present study is focused on assessing molecular *FM*, separator representation has been
520 simplified by the use of labels. There are four polymer separators in the dataset labeled 0-3. These labels are
521 concatenated with formulation representation analogous to temperature in IC dataset.

522 **Device:** Li-ICl battery dataset reports specific capacity of the battery with varying compositions of 8 electrolyte
523 constituents for a range of active material loadings (30% to 60%) in cathode and varying separators (18).
524 Electrolyte formulations are aggregated as defined for formulations and additional cell variables are concatenated
525 to formulation descriptor as model inputs.

526 For each dataset, neural network (NN) architectures are individually optimized and trained using the derived
527 dataset inputs. This feature engineering for representing molecules, formulations and devices was consistent
528 across all *FMs* and *MF*.

529 A.4 Model Training

530 It is noted that fine-tuning *FMs* such as SMI-TED with string representation of formulations could result in
531 relatively higher mean squared error (MSE) than the transfer learning approach where formulation descriptor
532 aggregates pre-learned molecular embeddings scaled with the composition. MSE for both the approaches are
533 compared in Table S1 for IC dataset where finetuning achieves MSE 0.155 and transfer learning combined by
534 NN regressor achieved MSE 0.025.

Table S1: Mean squared error (MSE) for property prediction using SMI-TED

Dataset	MSE	
	Fine-tuning	Transfer learning
Reduction Potential	0.65	0.68
Oxidation Potential	0.13	0.14
Ionic Conductivity	0.155	0.025

535 **Hyperparameter Tuning:** Neural network (NN) architectures were individually optimized and trained
536 using *FM*-derived molecular embeddings or formulation descriptor. NN with 2 or 3 hidden layers, with nodes
537 500-250-100 or 500-250, and activation function relu was found optimum. Model was trained with learning rate
538 0.0001, factoring 0.5 every 200 epochs of no reduction in loss function. The model was trained for maximum of
539 2500 epochs or until 200 iterations of no improvement in validation loss. Batch size was varied based on data
540 size. For datasets < 200, batch size was kept 1, batch size was 12 for dataset <5000, and batch size of 32 was
541 used for data >5000. Regression loss was measured using mean squared error (MSE) and mean absolute error
542 (MAE) was the used metric. For binary classification of electrolyte solubility, binary cross entropy was the loss
543 function and accuracy was the metric.

Table S2: Tuning neural network hyperparameters for SMI-TED predictors

Dataset	Hidden layers	Activation Function	MAE
LCE	500-250-100	relu	0.17
LCE	500-250	relu	0.16
LCE	500-250	sigmoid	0.32
LCE	500-250-100	sigmoid	0.32
LCE	500-250-250	relu	0.16
LCE	500-500	relu	0.17
LCE	250-100	relu	0.16
IC	500-250-100	relu	0.08
IC	500-250-100	sigmoid	0.22
IC	500-250	relu	0.09
IC	500-500	relu	0.10
IC	250-250-250	relu	0.08
IC	700-700	relu	0.11
IC	500-250-100-50	relu	0.08
HOMO	500-250-100	relu	0.43
HOMO	500-250-100	sigmoid	0.44
HOMO	500-250	relu	0.44
HOMO	250-100	relu	0.44
HOMO	500-500-500	relu	0.44
HOMO	250-250-250	relu	0.44

544 A.5 Out-of-distribution (OOD) evaluation

545 Two-fold OOD evaluation is done: (1) tail end evaluation based on numerical distribution of outcome labels, and
 546 (2) chemical design evaluation based on chemical similarity between train-test sets. For tail-end evaluation, test
 547 set are created from the training data to include lower and upper end values. In certain cases such as in Figure S3
 548 and Figure S4, only one end of data was considered as the outcome label was highly biased towards the other
 549 end.

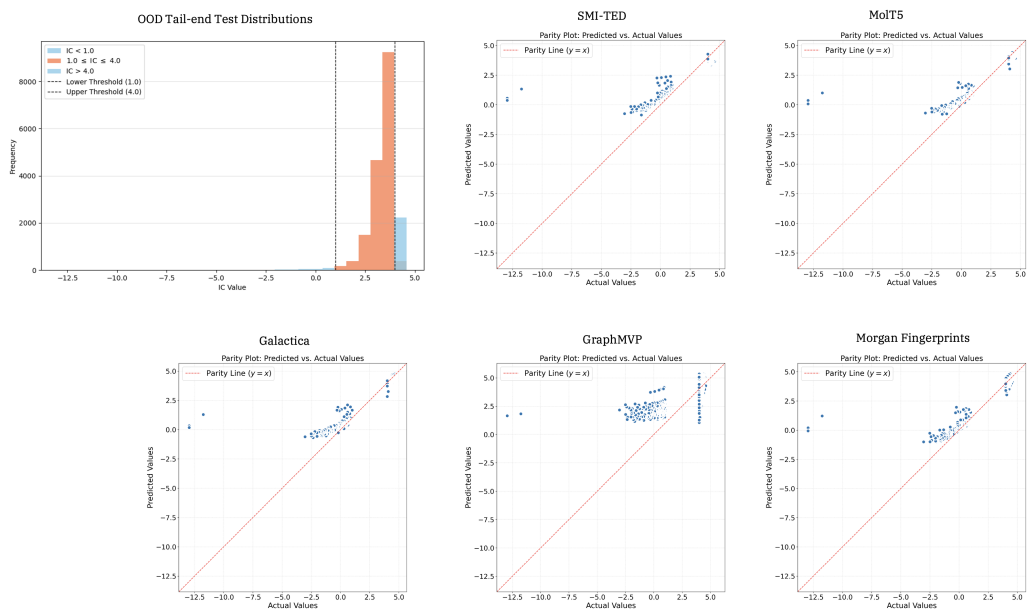


Figure S1: Tail-end OOD and parity plots for ionic conductivity test sets using benchmarking models.

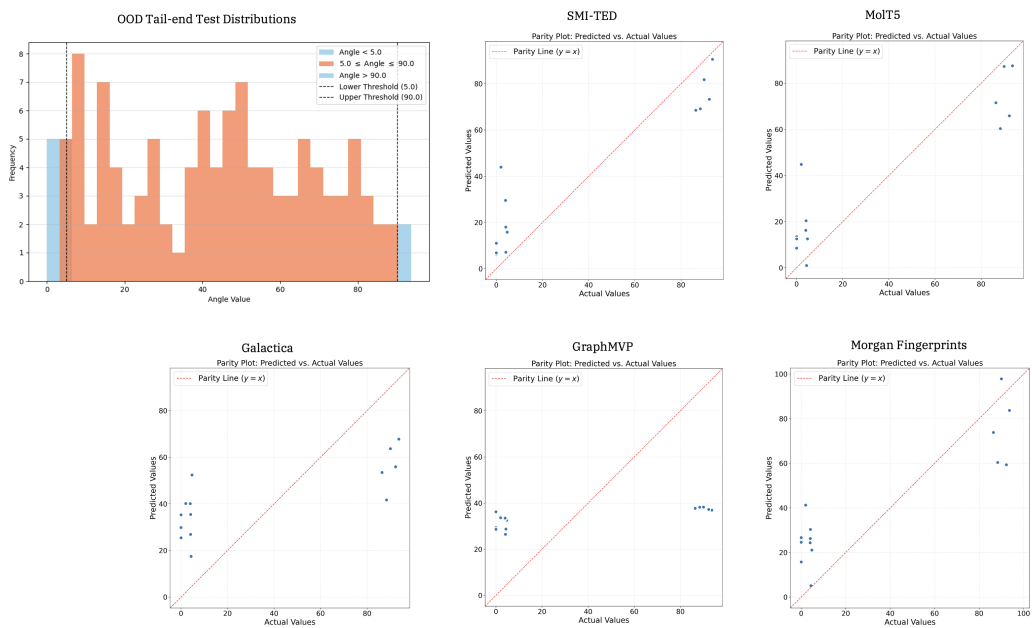


Figure S2: Tail-end OOD and parity plots for contact angle test sets using benchmarking models.

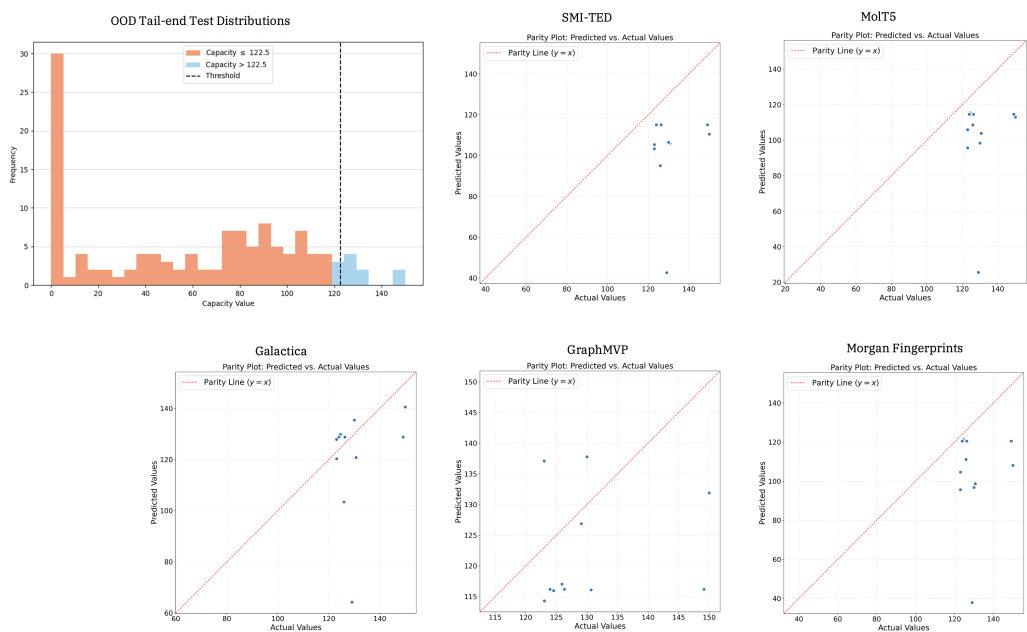


Figure S3: Tail-end OOD and parity plots for LiI capacity test sets using benchmarking models.

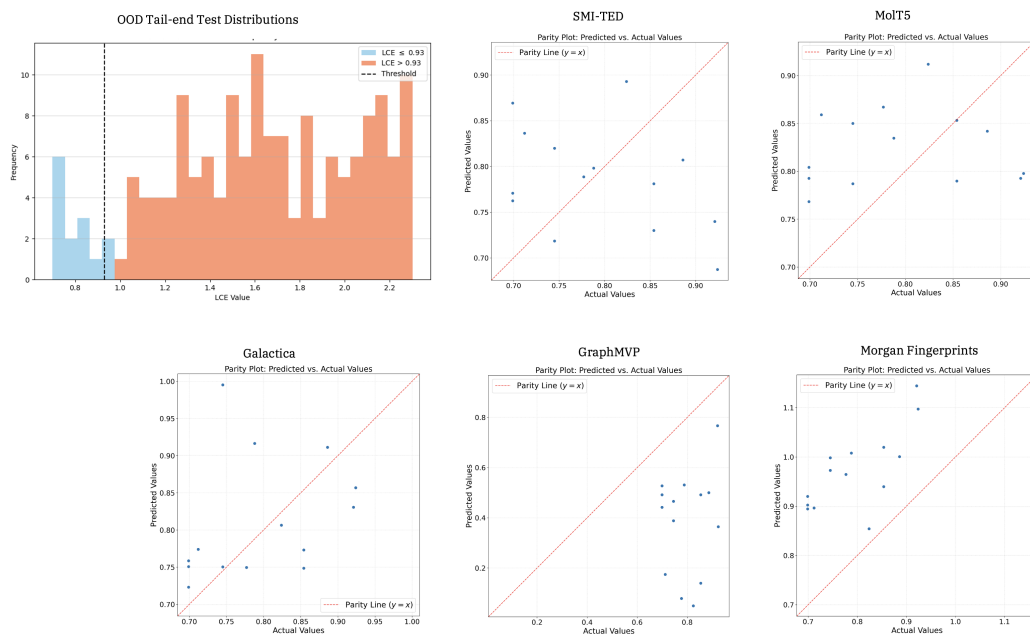


Figure S4: Tail-end OOD and parity plots for LCE test sets using benchmarking models.

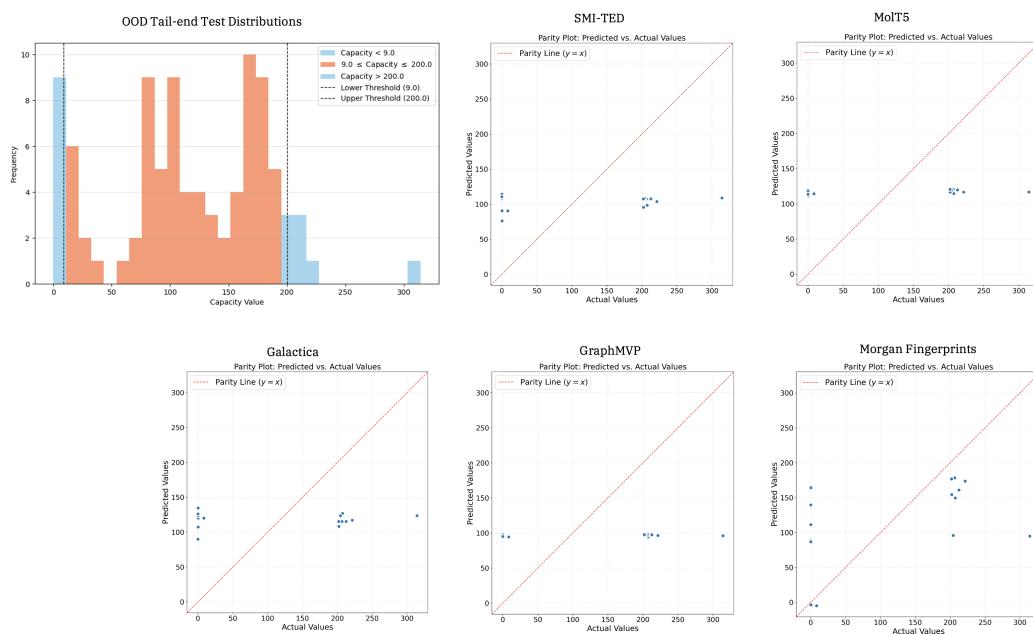


Figure S5: Tail-end OOD and parity plots for Li-ICI Capacity test sets using benchmarking models.

Table S3: Chemical similarity of out-of-distribution test datasets with training data using embeddings from foundation models and Morgan Fingerprints

Model	CE	Contact Angle	LiI Capacity	IC	Li-ICI Capacity
SMI-TED	0.3324	0.6791	0.2557	0.9244	0.6021
MolT5	0.2592	0.5472	0.1868	0.8209	0.641
Galactica	0.1925	0.6556	0.4531	0.9178	0.681
GraphMVP	0.0514	0.1099	0.0619	0.1814	0.0206
MF	0.2198	0.3281	0.1144	0.751	0.4748

Table S4: Parameters to estimate mean absolute error (MAE) in model prediction based on similarity between test-train data for SMI-TED

Datasets	Slope(m)	Intercept(c)
HOMO	-0.1602	0.5699
Ionic Conductivity	-0.5724	0.6377
Contact Angle	-19.6820	0.7601
Specific Capacity	-24.9776	33.2050