

1 **Supplementary Information**

2 **Machine learning-guided scalable manufacturing of high-**  
3 **efficiency perovskite solar modules**

4 *Yuchen Bu<sup>1</sup>, Junjie Gu<sup>1</sup>, Yu Shen<sup>1</sup>, Bairu Li<sup>1</sup>, Cheng Hu<sup>2,\*</sup>, Jianning Ding<sup>1,\*</sup> and Long*  
5 *Luo<sup>1,\*</sup>*

6 <sup>1</sup>Y. Bu, J. Gu, Y. Shen, B. Li, J. Ding, L. Luo

7 College of Physical Science and Technology, Institute of Technology for Carbon

8 Neutralization, Yangzhou University, Jiangsu 225009, China

9 <sup>2</sup>Dr. C. Hu

10 Business school, Institute of Technology for Carbon Neutralization, Yangzhou

11 University, Jiangsu 225009, China

12

13 \*Corresponding author:

14 [longluo@yzu.edu.cn](mailto:longluo@yzu.edu.cn); [dingjn@yzu.edu.cn](mailto:dingjn@yzu.edu.cn); [hucheng@yzu.edu.cn](mailto:hucheng@yzu.edu.cn)

15

## 16 **Dataset construction**

17 We compiled a dataset of 339 data points extracted from 113 peer-reviewed  
18 articles published between 2016 and 2025, with reported power conversion efficiencies  
19 (PCEs) ranging from 4.2% to 26.93%. The data were systematically collected using the  
20 keyword "perovskite solar modules," covering both formal (*n-i-p*) and inverted (*p-i-n*)  
21 device architectures, as well as six fabrication techniques: spin-coating, blade-coating,  
22 slot-die coating, evaporation, inkjet printing, and vacuum-assisted solution processing  
23 (VASP). The dataset comprises 255 large-area devices (active area  $\geq 10 \text{ cm}^2$ ) and 77  
24 related small-area devices (active area  $< 10 \text{ cm}^2$ ), with each entry characterized by 23  
25 input features.

26 Following initial data compilation, a rigorous multi-step preprocessing protocol  
27 was applied, resulting in a refined dataset of 332 valid entries. The preprocessing  
28 procedure included the following steps: (1) Removal of statistical outliers that deviated  
29 significantly from the majority of the data distribution. (2) Exclusion of samples with  
30 excessive missing feature values that could not be reliably imputed. (3) Recalculation  
31 and standardization of perovskite chemical compositions based on precursor ion  
32 stoichiometries. (4) Imputation of missing continuous variables, such as perovskite  
33 component ratios, using mean values to enhance data completeness. (5) Conversion of  
34 categorical charge transport layer materials into continuous variables by assigning  
35 representative conduction and valence band energy values. (6) Systematic classification  
36 of diverse additives and passivators into four major functional categories according to  
37 their chemical characteristics and roles in device operation.

38

## 39 **Model selection**

40 **Light Gradient Boosting Machine (LightGBM) algorithm:** LightGBM is a high-  
41 efficiency open-source framework based on gradient boosting decision trees. It is  
42 designed for large-scale data and high-speed training, and while maintaining high  
43 accuracy, it significantly improves training efficiency and reduces memory usage.  
44 Compared with other gradient boosting libraries, LightGBM's most significant feature  
45 is that it adopts a leaf-growing decision tree growth strategy. It selects the leaf node

46 with the largest gain for growth at each split, thus reducing the loss function more  
47 quickly.

48 **Support vector regression (SVR) algorithm:** SVR is a regression method based on  
49 support vector machines. The core goal is to find a function that allows for an error  
50 band of  $\epsilon$  precision when fitting data, and strives to make the function as flat as possible.  
51 It mainly uses the samples outside the error band to determine the decision boundary  
52 through optimization process. This method can effectively model the complex  
53 nonlinear relationship between features and targets by introducing kernel functions.

54 **Random forest (RF) algorithm:** Random forest is an ensemble learning algorithm that  
55 constructs multiple decision trees and aggregates their prediction results. It adopts a  
56 dual random strategy during training: Bootstrap sampling and feature random selection,  
57 effectively improving model diversity and generalization ability, and effectively  
58 suppressing overfitting.

59 **Multilayer perceptron (MLP) algorithm:** MLP is a feedforward neural network that  
60 consists of multiple layers. It consists of an input layer, at least one hidden layer, and  
61 an output layer, connected by adjustable weights between layers. The model calculates  
62 the output through forward propagation and updates the weight parameters using the  
63 backpropagation algorithm. This structure enables MLP to automatically learn complex  
64 nonlinear transformations between input and output, and has strong function  
65 approximation capabilities.

66 **Adaptive boost (AdaBoost) algorithm:** AdaBoost is an iterative ensemble learning  
67 algorithm that trains a series of weak learners sequentially, adjusting sample weights  
68 based on the previous round's prediction results each time, so that subsequent learners  
69 focus more on the previously incorrectly predicted samples. Finally, the prediction  
70 results of all weak learners are integrated through weighted voting.

71 **Extreme gradient boosting (XGBoost) algorithm:** XGBoost is an optimization  
72 implementation based on the gradient boosting framework, which constructs a  
73 sequence of decision trees in sequence, with each tree dedicated to correcting the  
74 prediction error of the preceding model. Its innovation lies in using second-order Taylor  
75 expansion to approximate the loss function and introducing regularization terms to

76 control model complexity. This algorithm has the advantages of high computational  
77 efficiency and strong ability to prevent overfitting. Furthermore, its sparsity-aware split  
78 finding enables robust performance even with small sample sizes by automatically  
79 learning handling strategies for missing or sparse data. Additionally, XGBoost  
80 effectively addresses imbalanced data distributions through customizable scale  
81 weights, ensuring high predictive accuracy and generalization in highly skewed  
82 datasets.

83

#### 84 **Model evaluation**

85 **Coefficient of determination ( $R^2$ ):**  $R^2$  represents the degree to which the quantitative  
86 regression model explains the variance changes of the target variable. Its numerical  
87 value represents the proportion of the total variance of the dependent variable that can  
88 be explained by the independent variables in the model. The value range of  $R^2$  is  $[0,1]$ .  
89 The closer the value is to 1, the better the model's fit to the data, indicating that the  
90 model can explain more data fluctuations. On the contrary, if  $R^2$  is low, it means that  
91 the model has failed to effectively capture the inherent variability of the data.

92 **Root mean square error (RMSE):** RMSE is the square root of the mean sum of  
93 squared differences between predicted and true values, with the same unit as the  
94 original data. Due to squaring during the calculation process, it has a stronger penalty  
95 for larger prediction errors and is therefore more sensitive to outliers.

96 **Mean absolute percentage error (MAPE):** MAPE evaluates the relative prediction  
97 error of the model by calculating the average ratio of the absolute error of each data  
98 point to the corresponding observation value. It is expressed in percentage form,  
99 making it a dimensionless indicator that facilitates comparison between prediction tasks  
100 at different scales or units. The lower the MAPE value, the smaller the average relative  
101 error of the model.

#### 102 **Code abstract**

##### 103 **Model training**

```
104 import shap  
105 import pandas as pd
```

```

106 import numpy as np
107 from sklearn.metrics import r2_score, mean_squared_error,
108 mean_absolute_percentage_error, accuracy_score, classification_report,
109 confusion_matrix
110 import matplotlib.pyplot as plt
111 import seaborn as sns
112 from matplotlib import rcParams
113 import warnings
114 from sklearn.utils import resample
115 from sklearn.model_selection import RandomizedSearchCV, KFold
116 from sklearn.isotonic import IsotonicRegression
117 import xgboost as xgb
118 from tqdm import tqdm
119 warnings.filterwarnings('ignore')
120 rcParams['font.family'] = 'Times New Roman'
121 rcParams['mathtext.fontset'] = 'stix'
122 df = pd.read_excel("data set.xlsx", sheet_name="Sheet1")
123 columns_to_drop = [df.columns[0], df.columns[-1]]
124 for col in df.columns:
125     if 'FF1' in str(col): columns_to_drop.append(col)
126 df = df.drop(columns=columns_to_drop)
127 df.columns = [col.strip() for col in df.columns]
128 df['Cation_sum'] = df[['Cs', 'FA', 'MA']].sum(axis=1)
129 df['Cs'] /= df['Cation_sum']; df['FA'] /= df['Cation_sum']; df['MA'] /= df['Cation_sum']
130 df = df.drop(columns=['Cation_sum'])
131 def adjust_halogens(row):
132     br, i, cl = row['Br'], row['I'], row['Cl']
133     if br > 0 and cl > 0:
134         if br >= cl: cl = 0
135         else: br = 0
136     total = br + i + cl
137     if total > 0:
138         scale = 3 / total
139         br *= scale; i *= scale; cl *= scale
140     return pd.Series([br, i, cl], index=['Br', 'I', 'Cl'])
141 df[['Br', 'I', 'Cl']] = df.apply(adjust_halogens, axis=1)
142 additive_1_values = ['none', 'Organic ammonium salt', 'others', 'Inorganic salts']
143 additive_2_values = ['none', 'others', 'Inorganic salts']
144 passivator_values = ['none', 'Organic ammonium salts', 'others', 'polymer', 'Inorganic
145 molecules']
146 for col, vals in {'additive_1': additive_1_values, 'additive_2': additive_2_values,
147 'Surface passivator': passivator_values}.items():
148     df[col] = df[col].fillna('none')
149     for v in vals:

```

```

150         df[f'{col}_{v}'] = df[col].apply(lambda x: 1 if str(x).strip() == v else 0)
151 process_dict = {'Spin coating': 0, 'Blade Coating': 1, 'Slot-Die Coating': 2,
152 'Evaporation': 3, 'Inkjet printing': 4, 'VASP': 5}
153 df['Preparation process'] = df['Preparation process'].map(process_dict).fillna(-1)
154 numeric_cols = ['Voc', 'Jsc', 'FF', 'Area', 'Substrate function', 'perovskite Eg',
155                 'ETL_1 CB', 'ETL_1 VB', 'ETL_2 CB', 'ETL_2 VB',
156                 'HTL_1 CB', 'HTL_1 VB', 'HTL_2 CB', 'HTL_2 VB', 'electrode
157 function']
158 for col in numeric_cols:
159     if col in df.columns:
160         df[col] = pd.to_numeric(df[col], errors='coerce')
161         df[col] = df[col].fillna(df[col].median())
162 df = df.drop(columns=['additive_1', 'additive_2', 'Surface passivator'])
163 X = df.drop(columns=['Efficiency', 'Voc', 'Jsc', 'FF'])
164 y_pce = df['Efficiency']
165 y_voc_raw = df['Voc']
166 VOC_LOW = float(y_voc_raw.quantile(0.25)) # ≈ 1.020V
167 VOC_HIGH = float(y_voc_raw.quantile(0.75)) # ≈ 1.140V
168 def voc_to_class(v):
169     if v < VOC_LOW: return 0
170     elif v <= VOC_HIGH: return 1
171     else: return 2
172 y_voc = y_voc_raw.apply(voc_to_class)
173 y_jsc = df['Jsc']
174 y_ff = df['FF']
175 train_idx = resample(range(len(df)),
176 n_samples=int(0.8*len(df)), random_state=68)
177 test_idx = resample(range(len(df)),
178 n_samples=int(0.2*len(df)), random_state=76)
179 X_train, X_test = X.iloc[train_idx], X.iloc[test_idx]
180 distribution = {
181     0: {0: 32, 1: 44, 2: 97, 3: 7, 4: 17},
182     1: {0: 4, 1: 9, 2: 23, 3: 4, 4: 5},
183     2: {0: 8, 1: 4, 2: 22, 3: 2, 4: 1},
184     3: {0: 4, 1: 2, 2: 7, 3: 1, 4: 0},
185     4: {0: 5, 1: 3, 2: 2, 3: 1, 4: 11},
186     5: {0: 7, 1: 5, 2: 6, 3: 2, 4: 2},
187 }
188 _all_counts = [v for p in distribution.values() for v in p.values() if v > 0]
189 _median = float(np.median(_all_counts))
190 MAX_W = 3.0
191 def get_area_bin(area):
192     if area < 1: return 0
193     elif area < 15: return 1
194     elif area < 60: return 2
195     elif area <= 100: return 3
196     else: return 4
197 def compute_weights(X_sub):

```

```

194     w = []
195     for i in range(len(X_sub)):
196         proc = int(X_sub.iloc[i]["Preparation process"])
197         abin = get_area_bin(float(X_sub.iloc[i]["Area"]))
198         cnt = distribution.get(proc, {}).get(abin, 0)
199         w.append(min(_median / cnt, MAX_W) if cnt > 0 else 1.0)
200     return np.array(w)
201 def add_highval_weights(w_base, y_series, boost=2.5):
202     q75 = float(y_series.quantile(0.75))
203     high_mask = (y_series.values > q75)
204     w = w_base.copy()
205     w[high_mask] *= boost
206     return w
207 w_base = compute_weights(X_train)
208 w_pce = add_highval_weights(w_base, y_pce.iloc[train_idx], boost=2.0)
209 w_jsc = add_highval_weights(w_base, y_jsc.iloc[train_idx], boost=3.0)
210 w_ff = add_highval_weights(w_base, y_ff.iloc[train_idx], boost=2.5)
211 w_voc = None
212 param_dist = {
213     'n_estimators': [100, 200, 300, 400, 500],
214     'max_depth': [3, 4, 5, 6, 7],
215     'learning_rate': [0.01, 0.05, 0.08, 0.1, 0.15],
216     'subsample': [0.6, 0.7, 0.8, 0.9],
217     'colsample_bytree': [0.6, 0.7, 0.8, 0.9],
218     'reg_alpha': [0.0, 0.1, 0.5, 1.0],
219     'reg_lambda': [0.5, 1.0, 2.0, 3.0],
220     'min_child_weight': [1, 3, 5, 7],
221     'gamma': [0.0, 0.1, 0.2, 0.5],
222 }
223 param_dist_reg = {
224     'n_estimators': [200, 300, 400, 500, 600],
225     'max_depth': [3, 4, 5],
226     'learning_rate': [0.01, 0.03, 0.05, 0.08],
227     'subsample': [0.5, 0.6, 0.7, 0.8],
228     'colsample_bytree': [0.5, 0.6, 0.7, 0.8],
229     'reg_alpha': [0.5, 1.0, 2.0, 5.0],
230     'reg_lambda': [2.0, 3.0, 5.0, 8.0],
231     'min_child_weight': [5, 7, 10, 15],
232     'gamma': [0.1, 0.3, 0.5, 1.0],
233 }
234 def tune_model(estimator, X_tr, y_tr, sw=None, n_iter=80, pd=None):
235     cv = KFold(n_splits=5, shuffle=True, random_state=42)
236     fit_kw = {'sample_weight': sw} if sw is not None else {}
237     _pd = pd if pd is not None else param_dist

```

```

238     search = RandomizedSearchCV(
239         estimator, _pd, n_iter=n_iter, scoring='r2',
240         cv=cv, random_state=42, n_jobs=-1, refit=True, verbose=0
241     )
242     search.fit(X_tr, y_tr, **fit_kw)
243     return search.best_estimator_, search.best_params_
244 def tune_classifier(X_tr, y_tr, sw=None, n_iter=80):
245     param_dist_cls = {**param_dist,
246                       'use_label_encoder': [False],
247                       'eval_metric': ['mlogloss']}
248     cv = KFold(n_splits=5, shuffle=True, random_state=42)
249     fit_kw = {'sample_weight': sw} if sw is not None else {}
250     search = RandomizedSearchCV(
251         xgb.XGBClassifier(random_state=42, n_jobs=-1, verbosity=0),
252         param_dist_cls, n_iter=n_iter, scoring='accuracy',
253         cv=cv, random_state=42, n_jobs=-1, refit=True, verbose=0
254     )
255     search.fit(X_tr, y_tr, **fit_kw)
256
257     return search.best_estimator_, search.best_params_
258 def fit_isotonic(model, X_tr, y_tr):
259     y_raw = model.predict(X_tr)
260     iso = IsotonicRegression(out_of_bounds='clip')
261     iso.fit(y_raw, y_tr)
262     return iso
263 def predict_with_iso(model, iso, X):
264     y_raw = model.predict(X)
265     return iso.predict(y_raw)
266 iso_pce = fit_isotonic(pce_model, X_train, y_pce.iloc[train_idx])
267 iso_jsc = fit_isotonic(jsc_model, X_train, y_jsc.iloc[train_idx])
268 iso_ff = fit_isotonic(ff_model, X_train, y_ff.iloc[train_idx])
269 def get_metrics(y_true, y_pred):
270     r2 = r2_score(y_true, y_pred)
271     rmse = np.sqrt(mean_squared_error(y_true, y_pred))
272     mape = mean_absolute_percentage_error(y_true, y_pred) * 100
273     return r2, rmse, mape
274 y_pce_tr_pred = predict_with_iso(pce_model, iso_pce, X_train)
275 y_pce_te_pred = predict_with_iso(pce_model, iso_pce, X_test)
276 y_voc_tr_pred = voc_model.predict(X_train); y_voc_te_pred =
277 voc_model.predict(X_test)
278 y_voc_tr_proba = voc_model.predict_proba(X_train); y_voc_te_proba =
279 voc_model.predict_proba(X_test)
280 y_jsc_tr_pred = predict_with_iso(jsc_model, iso_jsc, X_train)
281 y_jsc_te_pred = predict_with_iso(jsc_model, iso_jsc, X_test)

```

```

282 y_ff_tr_pred = predict_with_iso(ff_model, iso_ff, X_train)
283 y_ff_te_pred = predict_with_iso(ff_model, iso_ff, X_test)
284 pce_tr_r2, pce_tr_rmse, pce_tr_mape = get_metrics(y_pce.iloc[train_idx],
285 y_pce_tr_pred)
286 pce_te_r2, pce_te_rmse, pce_te_mape = get_metrics(y_pce.iloc[test_idx],
287 y_pce_te_pred)
288 voc_tr_acc = accuracy_score(y_voc.iloc[train_idx], y_voc_tr_pred) * 100
289 voc_te_acc = accuracy_score(y_voc.iloc[test_idx], y_voc_te_pred) * 100
290 jsc_tr_r2, jsc_tr_rmse, jsc_tr_mape = get_metrics(y_jsc.iloc[train_idx], y_jsc_tr_pred)
291 jsc_te_r2, jsc_te_rmse, jsc_te_mape = get_metrics(y_jsc.iloc[test_idx],
292 y_jsc_te_pred)
293 ff_tr_r2, ff_tr_rmse, ff_tr_mape = get_metrics(y_ff.iloc[train_idx], y_ff_tr_pred)
294 ff_te_r2, ff_te_rmse, ff_te_mape = get_metrics(y_ff.iloc[test_idx],
295 y_ff_te_pred)
296
297 shap analysis
298 shap.initjs()
299 explainer_pce = shap.TreeExplainer(pce_model); shap_values_pce =
300 explainer_pce.shap_values(X_test)
301 explainer_jsc = shap.TreeExplainer(jsc_model); shap_values_jsc =
302 explainer_jsc.shap_values(X_test)
303 explainer_ff = shap.TreeExplainer(ff_model); shap_values_ff =
304 explainer_ff.shap_values(X_test)
305 explainer_voc = shap.TreeExplainer(voc_model)
306 shap_values_voc_all = explainer_voc.shap_values(X_test)
307 shap_values_voc = shap_values_voc_all[:, :, 2]
308 def plot_shap_summary(shap_values, features, title, filename):
309     plt.figure(figsize=(12, 8))
310     feature_importance = np.abs(shap_values).mean(axis=0)
311     top_indices = np.argsort(feature_importance)[-15:][::-1]
312     shap.summary_plot(shap_values[:, top_indices], features.iloc[:, top_indices],
313 show=False)
314     plt.title(f'{title} - SHAP', fontsize=20, fontname='Times New Roman')
315     plt.tight_layout()
316     plt.savefig(f'shap_summary_{filename}.png', dpi=300, bbox_inches='tight')
317     plt.show()
318 plot_shap_summary(shap_values_pce, X_test, 'PCE', 'pce')
319 plot_shap_summary(shap_values_voc, X_test, 'Voc', 'voc')
320 plot_shap_summary(shap_values_jsc, X_test, 'Jsc', 'jsc')
321 plot_shap_summary(shap_values_ff, X_test, 'FF', 'ff')
322
323 PCE Prediction
324 manual_composition = {
325     'Cs': 0.1,

```

```

326     'FA': 0.888888889,
327     'MA': 0.011111111,
328     'Br': 0.059055118,
329     'I': 2.940944882,
330     'Cl': 0.0,
331 }
332 manual_process_id = 5
333 for comp, value in manual_composition.items():
334     print(f" {comp}: {value}")
335 cation_sum = manual_composition['Cs'] + manual_composition['FA'] +
336 manual_composition['MA']
337 if cation_sum > 0:
338     manual_composition['Cs'] = manual_composition['Cs'] / cation_sum
339     manual_composition['FA'] = manual_composition['FA'] / cation_sum
340     manual_composition['MA'] = manual_composition['MA'] / cation_sum
341 halogen_sum = manual_composition['Br'] + manual_composition['I'] +
342 manual_composition['Cl']
343 if halogen_sum > 0:
344     scale = 3 / halogen_sum
345     manual_composition['Br'] = manual_composition['Br'] * scale
346     manual_composition['I'] = manual_composition['I'] * scale
347     manual_composition['Cl'] = manual_composition['Cl'] * scale
348 for comp, value in manual_composition.items():
349     print(f" {comp}: {value:.4f}")
350 process_names = {
351     0: 'Spin coating',
352     1: 'Blade Coating',
353     2: 'Slot-Die Coating',
354     3: 'Evaporation',
355     4: 'Inkjet printing',
356     5: 'VASP'
357 }
358 fixed_process = manual_process_id
359 other_fixed_params = {
360     'Preparation process': fixed_process,
361     'Substrate function': 4.5,
362     'electrode function': 5.28,
363     'Structure': 1,
364     'ETL_1 CB': 4.5,
365     'ETL_1 VB': 7.5,
366     'ETL_2 CB': 4.5,
367     'ETL_2 VB': 7.5,
368     'HTL_1 CB': 2.2,
369     'HTL_1 VB': 5.22,

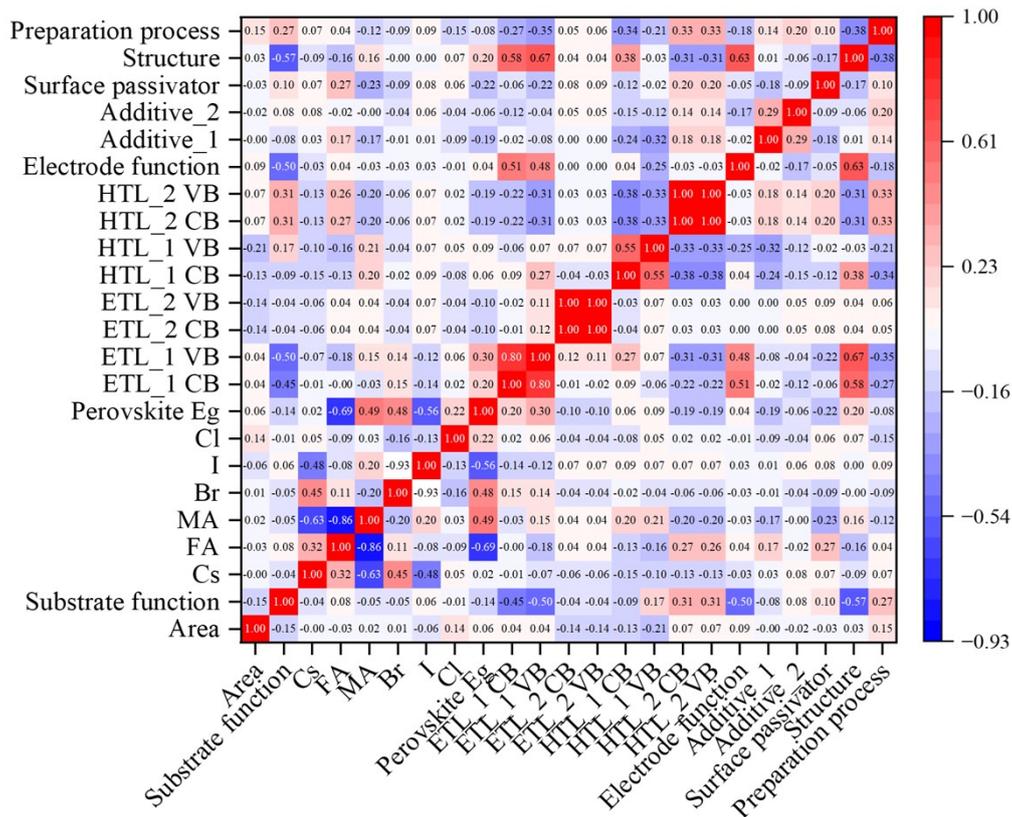
```

```

370     'HTL_2 CB': 0,
371     'HTL_2 VB': 0,
372 }
373 passivator_features = [col for col in X.columns if col.startswith('Passivator_')]
374 for pas in passivator_features:
375     if pas == 'Passivator_Organic_ammonium_salts':
376         other_fixed_params[pas] = 1
377     else:
378         other_fixed_params[pas] = 0
379 additive1_features = [col for col in X.columns if col.startswith('Additive1_')]
380 additive2_features = [col for col in X.columns if col.startswith('Additive2_')]
381 for add in additive1_features:
382     other_fixed_params[add] = 1 if add == 'Additive1_others' else 0
383 for add in additive2_features:
384     other_fixed_params[add] = 1 if add == 'Additive2_none' else 0
385 for feature in X.columns:
386     if feature not in manual_composition and feature not in other_fixed_params and
387 feature != 'Area':
388         if pd.api.types.is_numeric_dtype(X[feature]):
389             other_fixed_params[feature] = high_efficiency_samples[feature].median()
390         else:
391             other_fixed_params[feature] =
392 high_efficiency_samples[feature].value_counts().idxmax()
393 start_area, end_area, step_area = 0.03, 100.03, 0.5
394 area_range = np.arange(start_area, end_area + step_area, step_area)
395 area_results = []
396 for i, area in enumerate(tqdm(area_range, desc=f" {manual_process_id} ")):
397     param_dict = {'Area': area}
398     param_dict.update(manual_composition)
399     param_dict.update(other_fixed_params)
400     param_dict['Area_inv'] = 1 / (area + 0.001)
401     param_dict['log_Area'] = np.log(area + 0.001)
402     input_data = pd.DataFrame([param_dict])
403     for col in X.columns:
404         if col not in input_data.columns:
405             input_data[col] = other_fixed_params.get(col, 0)
406     input_data = input_data[X.columns]
407     try:
408         efficiency = predict_with_iso(pce_model, iso_pce, input_data)[0]
409         adjusted_efficiency = apply_process_specific_adjustment(efficiency, area,
410 manual_process_id)
411         area_results.append((adjusted_efficiency, area))
412     except Exception as e:
413         continue

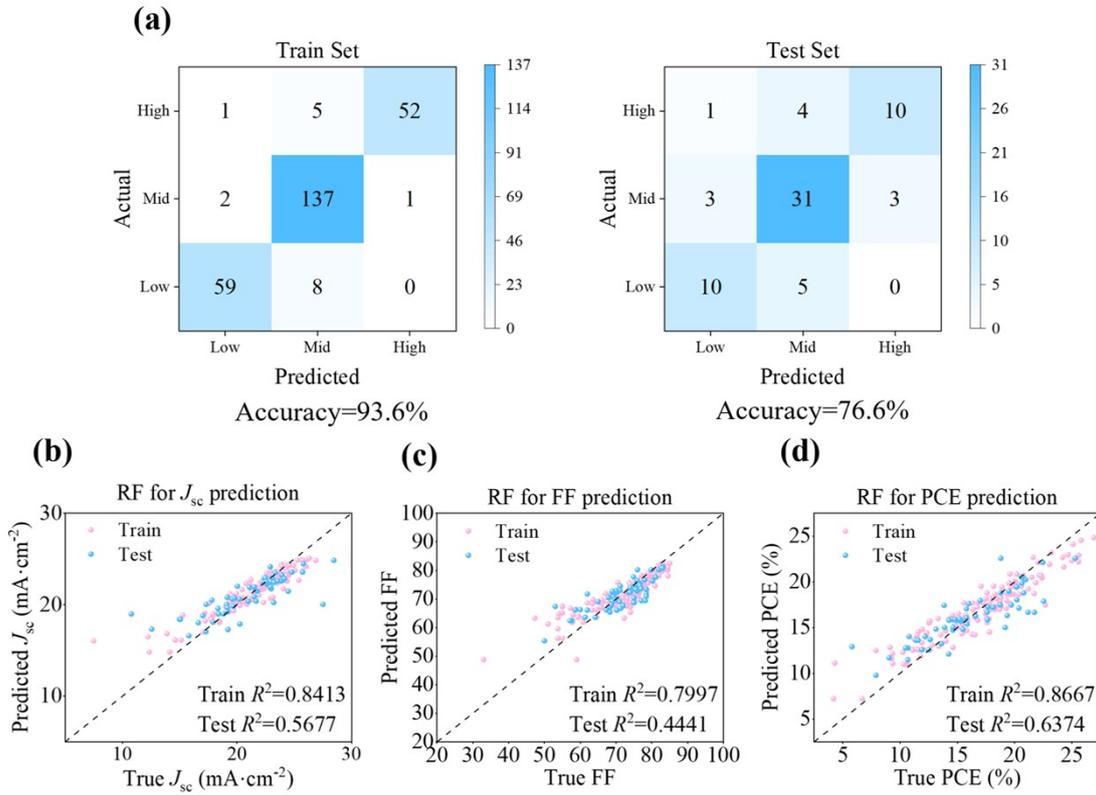
```





415

416 **Fig. S1** The heatmap representing the strength of correlation between input feature  
 417 quantities.

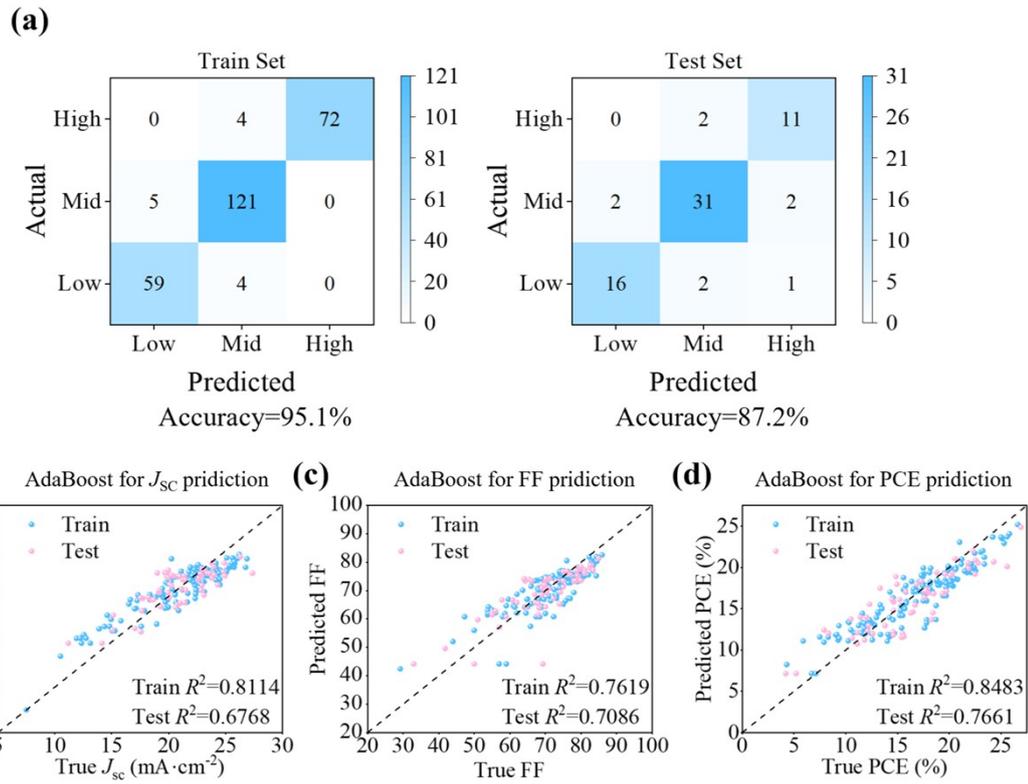


418

419 **Fig. S2** The RF model was used to train  $V_{OC}$ ,  $J_{SC}$ , FF and PCE. (a)  $V_{OC}$  confusion matrix.

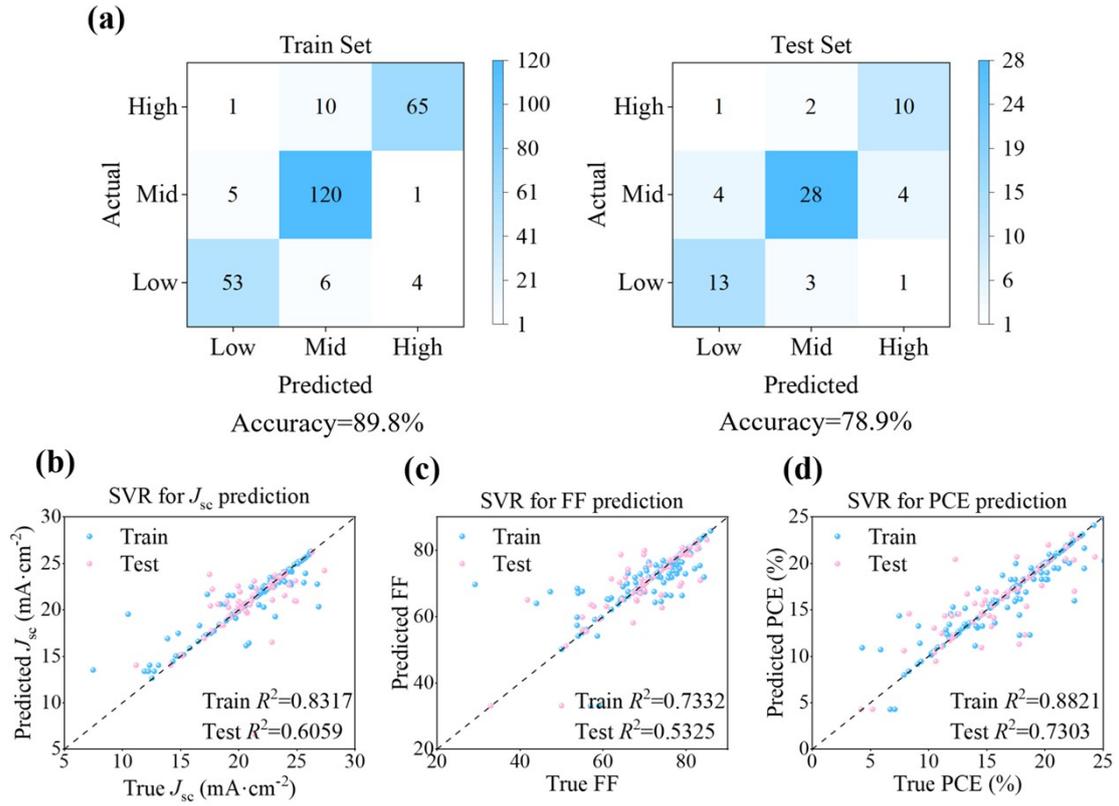
420 In the other figures, red represents the training set and blue represents the testing set:

421 (b)  $J_{SC}$ , (c) FF, (d) PCE.



422

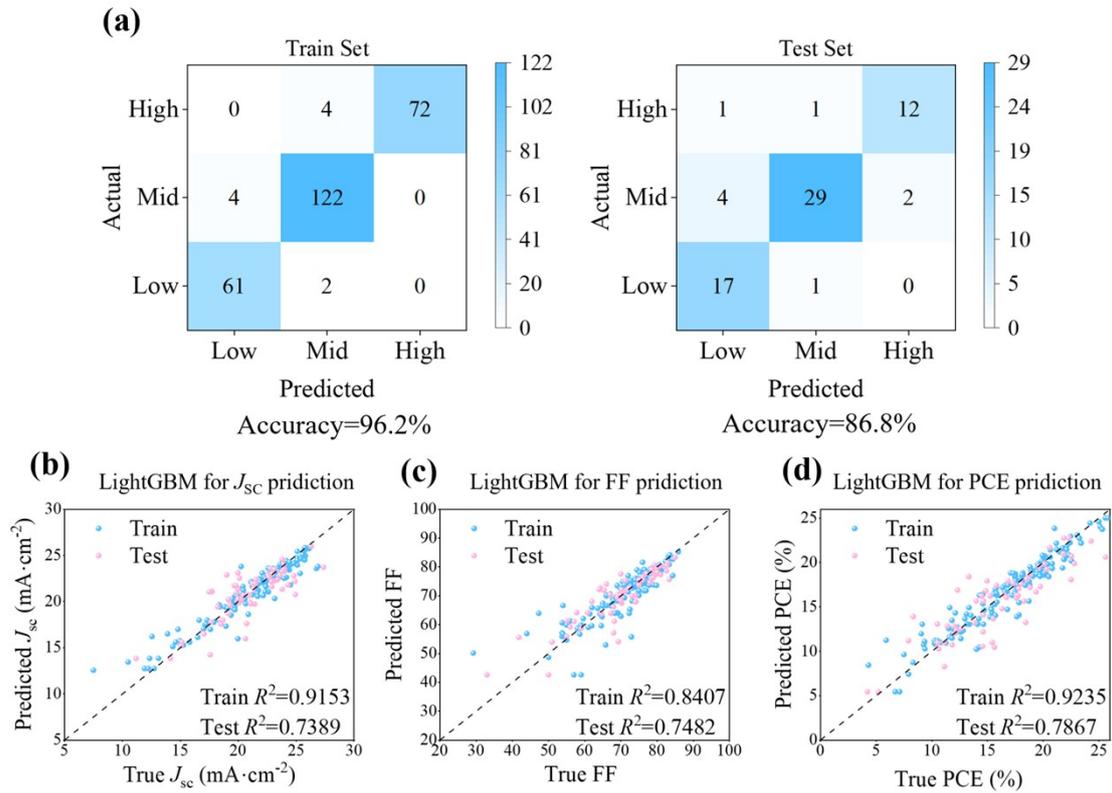
423 **Fig. S3** The AdaBoost model was used to train  $V_{OC}$ ,  $J_{SC}$ , FF and PCE. (a)  $V_{OC}$  confusion  
 424 matrix. In the other figures, red represents the training set and blue represents the testing  
 425 set: (b)  $J_{SC}$ , (c) FF, (d) PCE.



426

427 **Fig. S4** The SVR model was used to train  $V_{OC}$ ,  $J_{SC}$ , FF and PCE. (a)  $V_{OC}$  confusion  
 428 matrix. In the other figures, red represents the training set and blue represents the testing  
 429 set: (b)  $J_{SC}$ , (c) FF, (d) PCE.

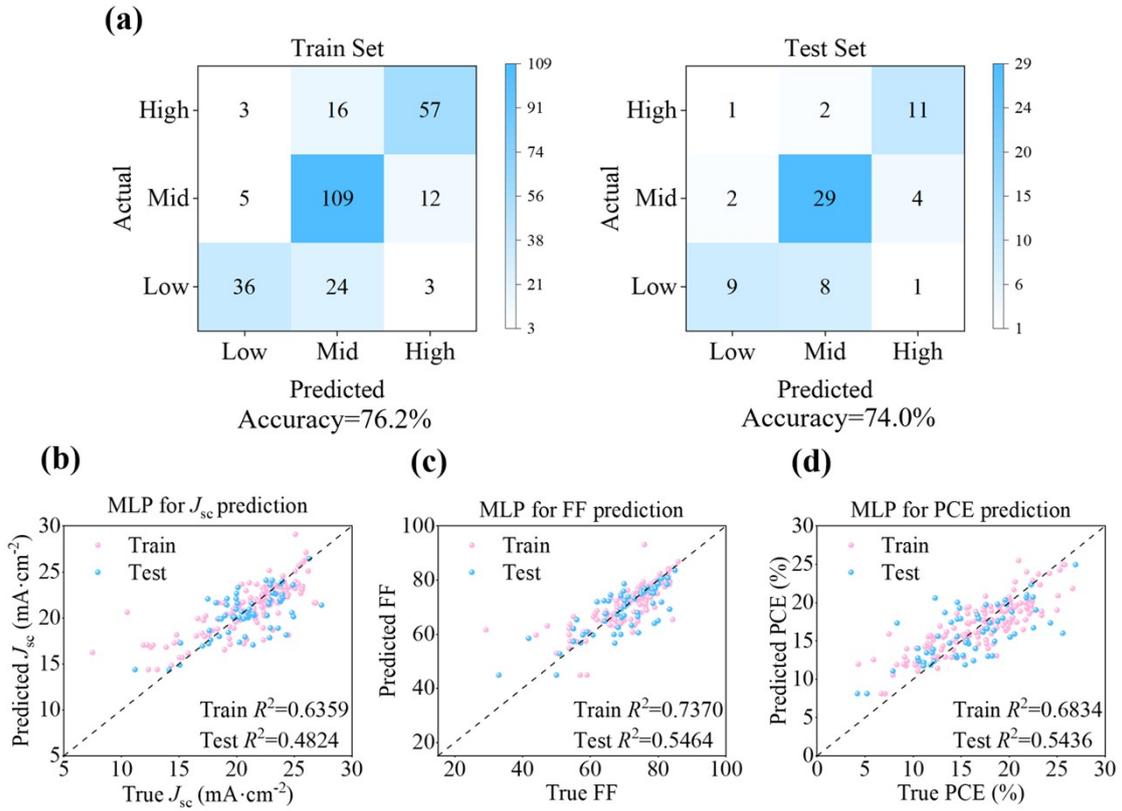
430



431

432 **Fig. S5** The LightGBM model was used to train  $V_{OC}$ ,  $J_{SC}$ , FF and PCE. (a)  $V_{OC}$   
 433 confusion matrix. In the other figures, red represents the training set and blue represents  
 434 the testing set: (b)  $J_{SC}$ , (c) FF, (d) PCE.

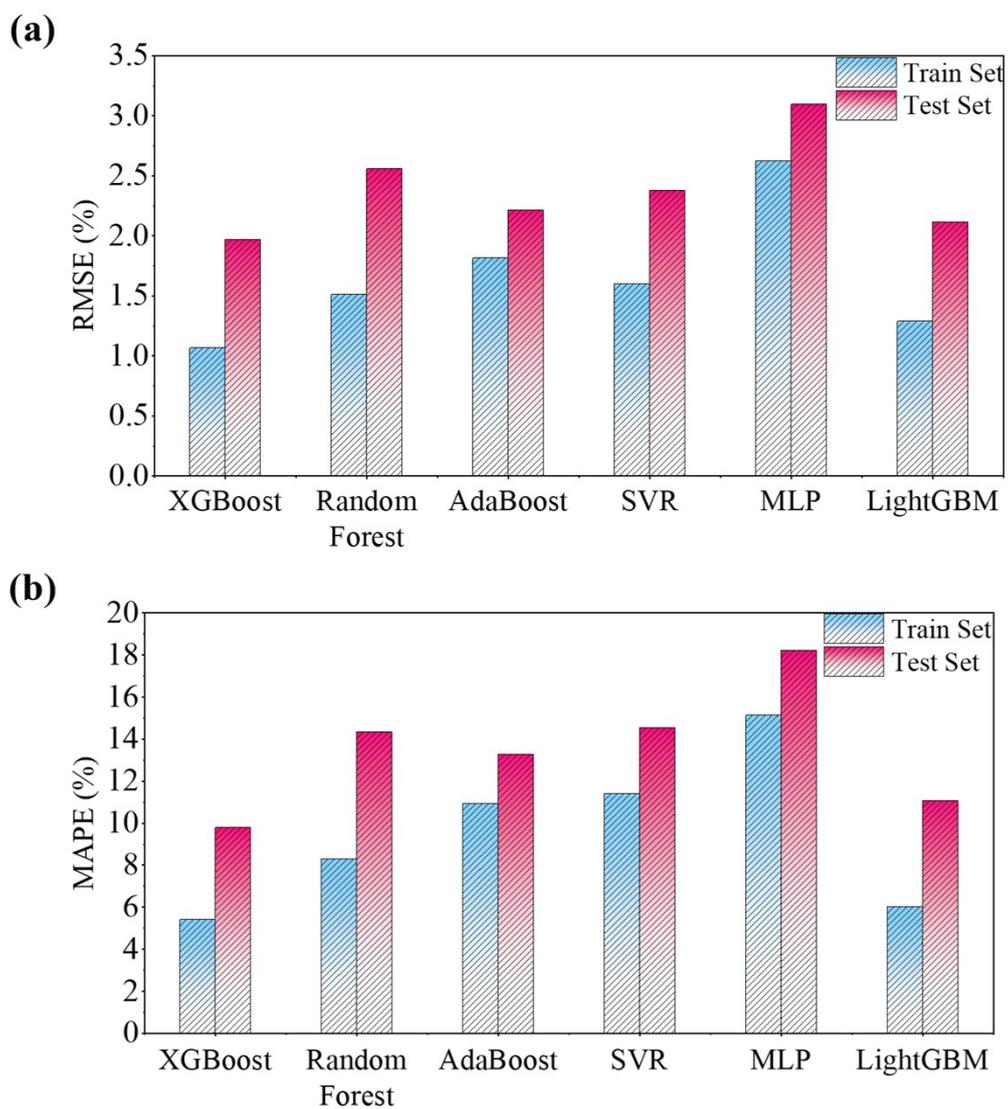
435



436

437 **Fig. S6** The MLP model was used to train  $V_{OC}$ ,  $J_{SC}$ , FF and PCE. (a)  $V_{OC}$  confusion  
 438 matrix. In the other figures, red represents the training set and blue represents the testing  
 439 set: (b)  $J_{SC}$ , (c) FF, (d) PCE.

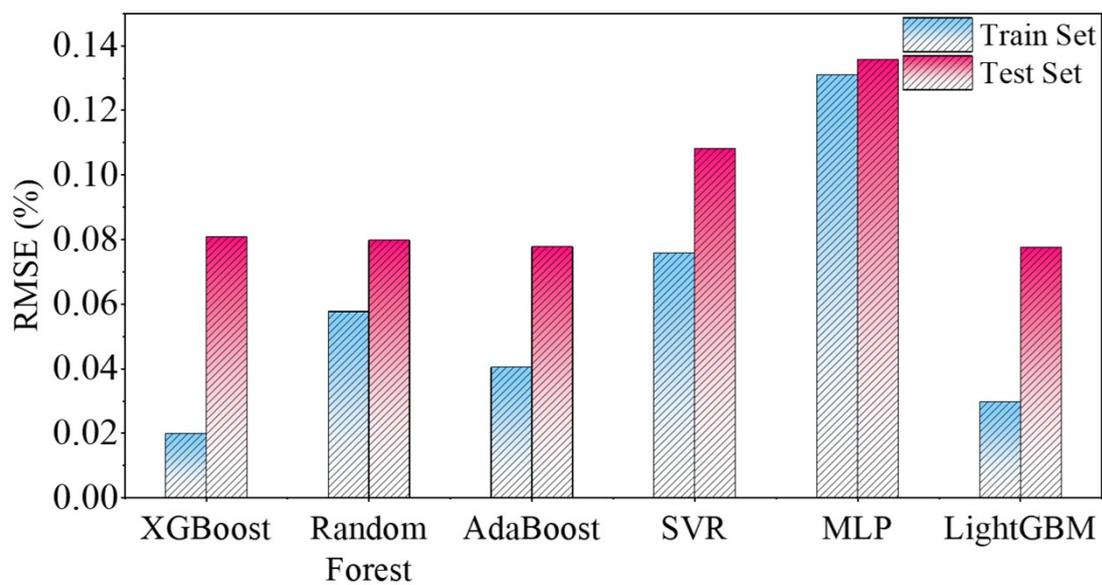
440



441

442 **Fig. S7** Performance of 6 ML models: **(a)** RMSE, **(b)** MAPE.

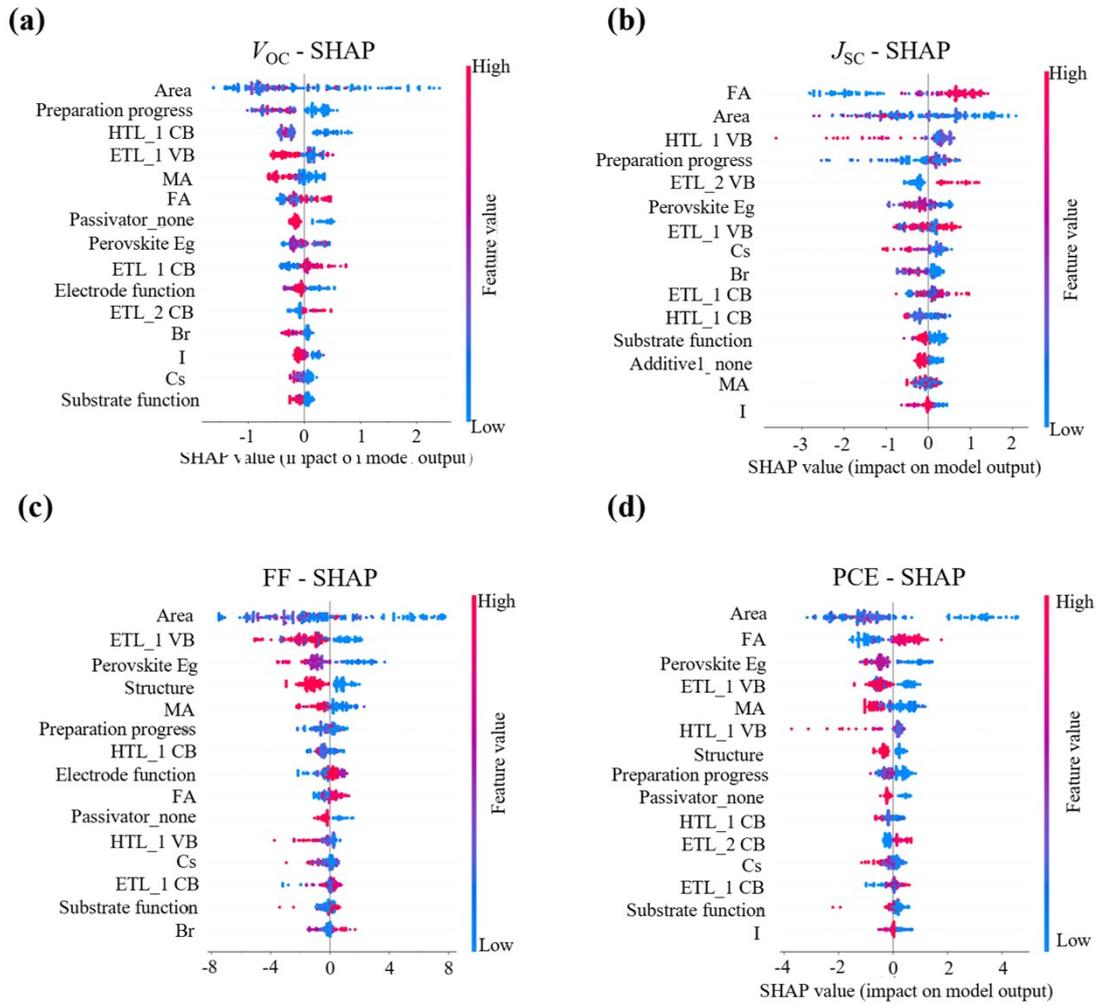
443



444

445 **Fig. S8** RMSE of six models trained on  $V_{oc}$ .

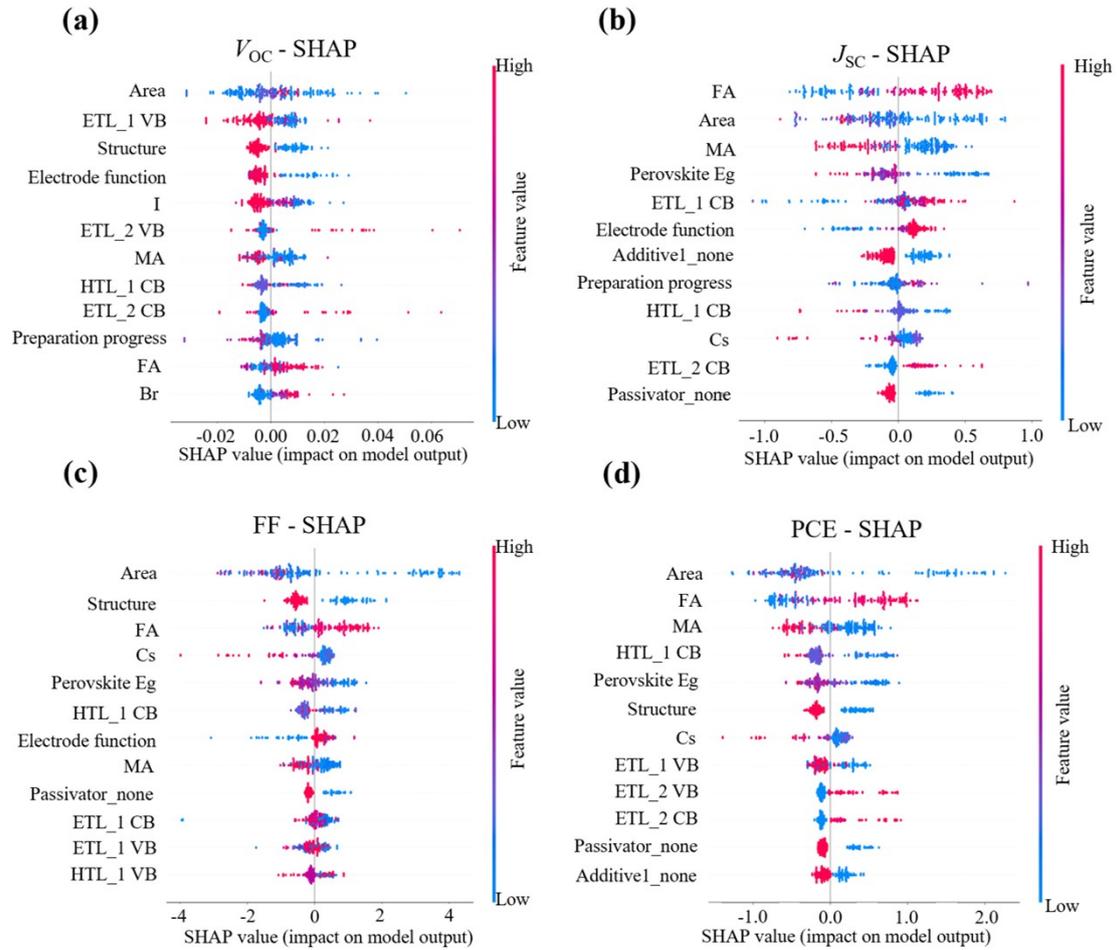
446



447

448 **Fig. S9** The SHAP analysis diagram of characteristic parameters  $V_{OC}$ ,  $J_{SC}$ , FF, PCE  
 449 affecting perovskite solar cells based on XGBoost model: (a)  $V_{OC}$ , (b)  $J_{SC}$ , (c) FF, and  
 450 (d) PCE.

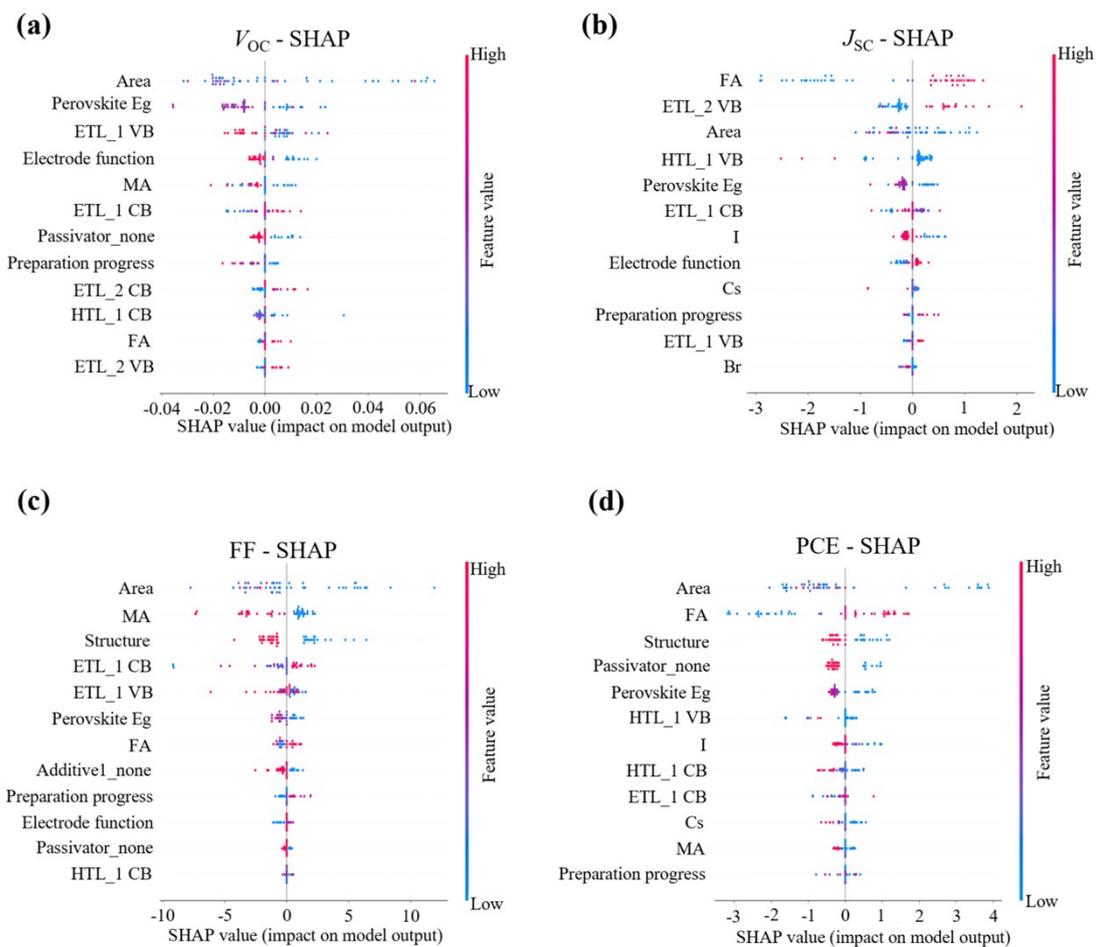
451



452

453 **Fig. S10** The SHAP analysis diagram of characteristic parameters  $V_{OC}$ ,  $J_{SC}$ , FF, PCE  
 454 affecting perovskite solar cells based on RF model: **(a)**  $V_{OC}$ , **(b)**  $J_{SC}$ , **(c)** FF, and **(d)**  
 455 PCE.

456



457

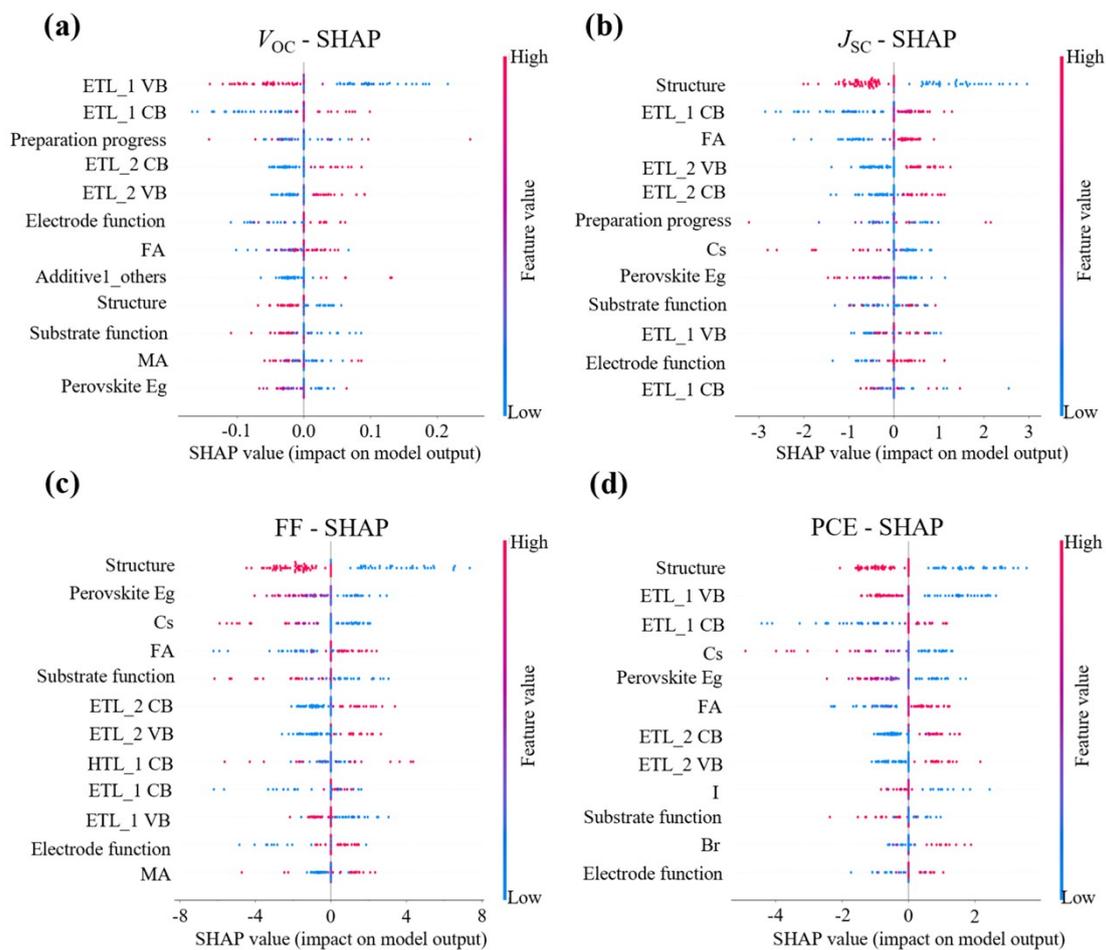
458 **Fig. S11** The SHAP analysis diagram of characteristic parameters  $V_{OC}$ ,  $J_{SC}$ , FF, PCE

459 affecting perovskite solar cells based on AdaBoost model: (a)  $V_{OC}$ , (b)  $J_{SC}$ , (c) FF, and

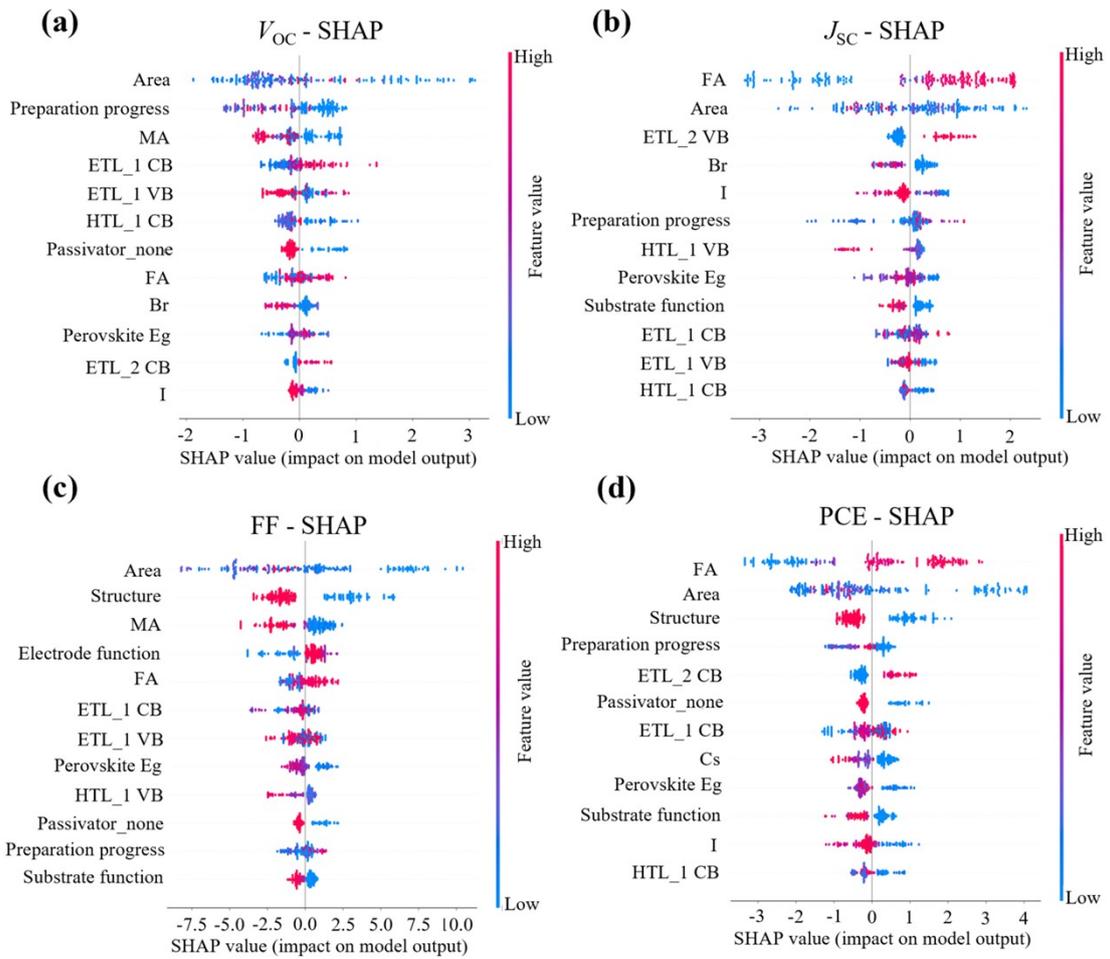
460 (d) PCE.

461

462



464 **Fig. S12** The SHAP analysis diagram of characteristic parameters  $V_{OC}$ ,  $J_{SC}$ , FF, PCE  
 465 affecting perovskite solar cells based on SVR model: (a)  $V_{OC}$ , (b)  $J_{SC}$ , (c) FF, and (d)  
 466 PCE.  
 467



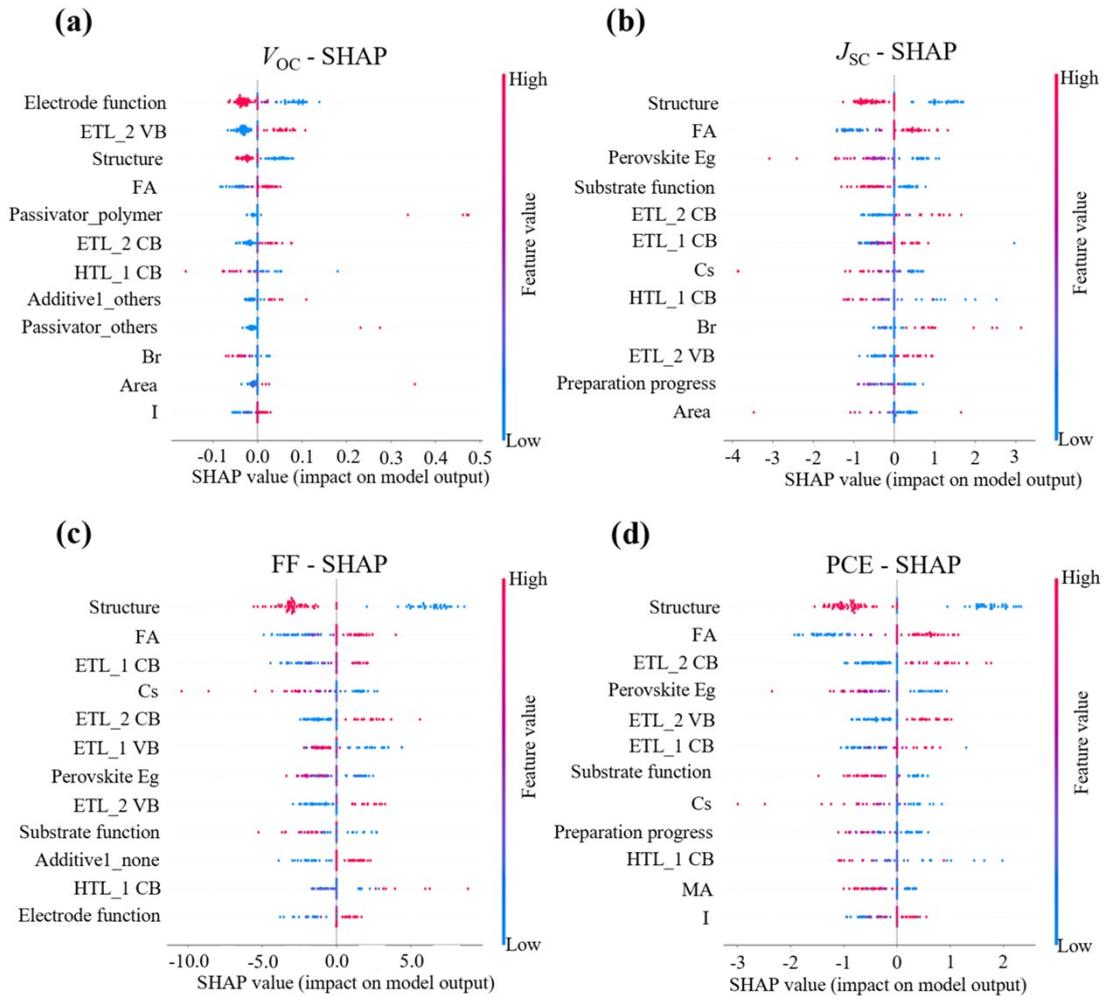
468

469 **Fig. S13** The SHAP analysis diagram of characteristic parameters  $V_{OC}$ ,  $J_{SC}$ , FF, PCE

470 affecting perovskite solar cells based on LightGBM model: (a)  $V_{OC}$ , (b)  $J_{SC}$ , (c) FF, and

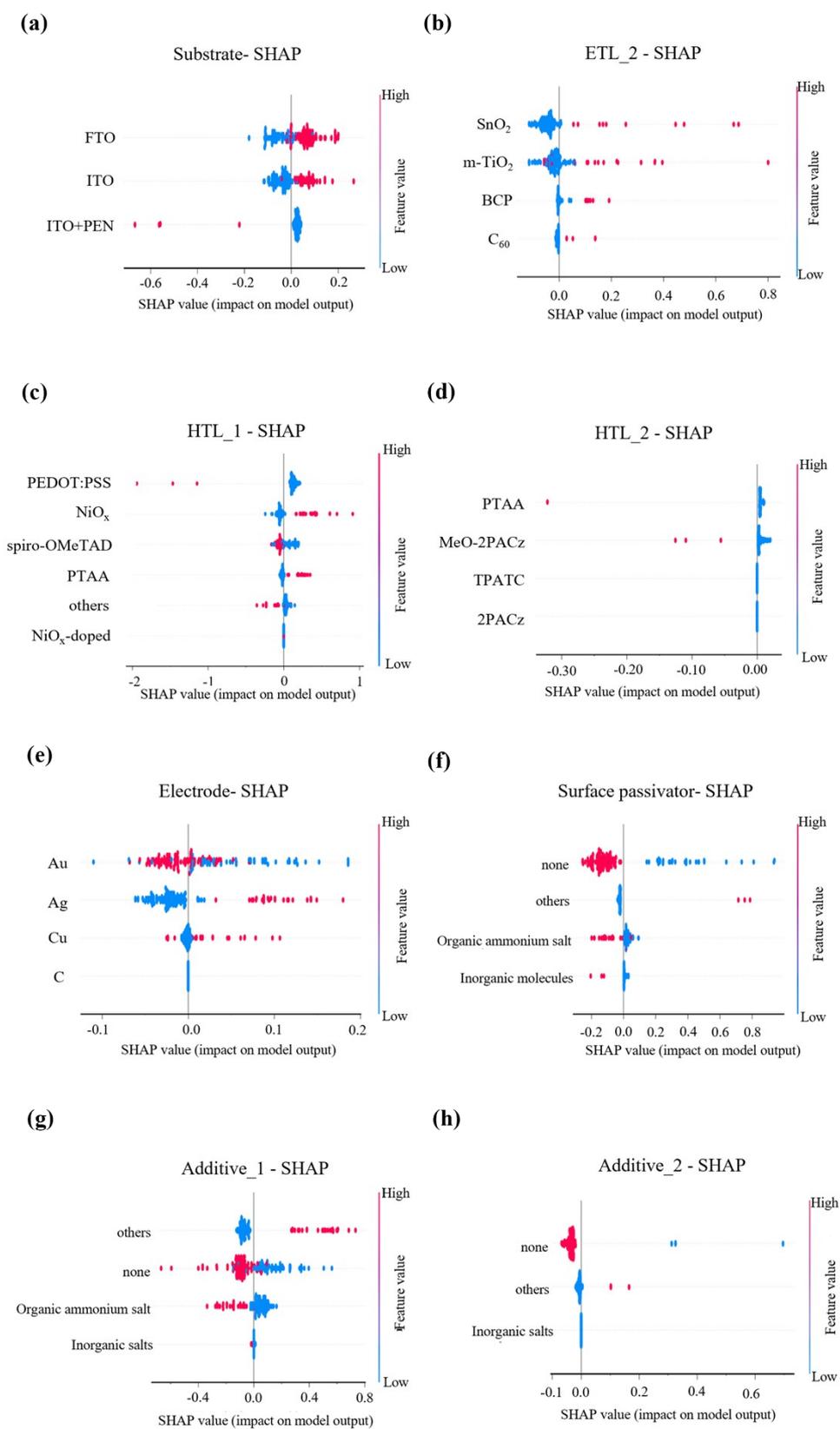
471 (d) PCE.

472



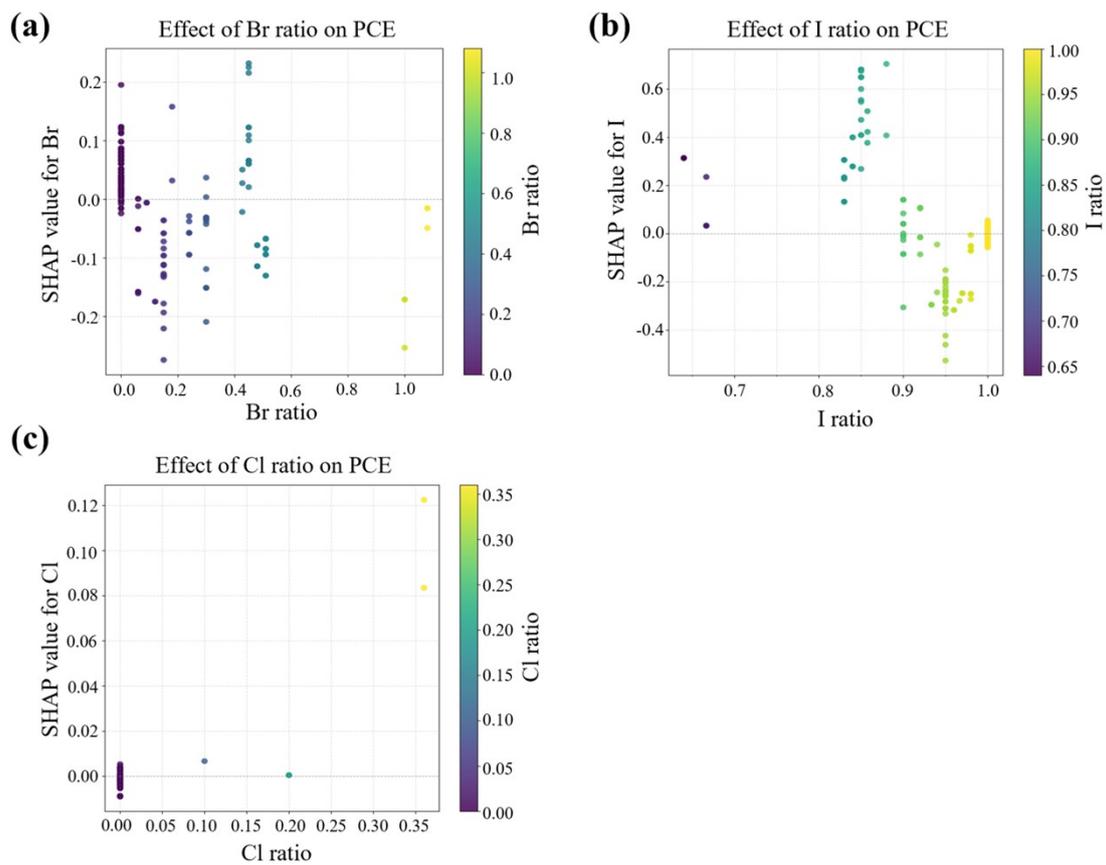
473

474 **Fig. S14** The SHAP analysis diagram of characteristic parameters  $V_{OC}$ ,  $J_{SC}$ , FF, PCE  
 475 affecting perovskite solar cells based on MLP model: **(a)**  $V_{OC}$ , **(b)**  $J_{SC}$ , **(c)** FF, and **(d)**  
 476 PCE.



477

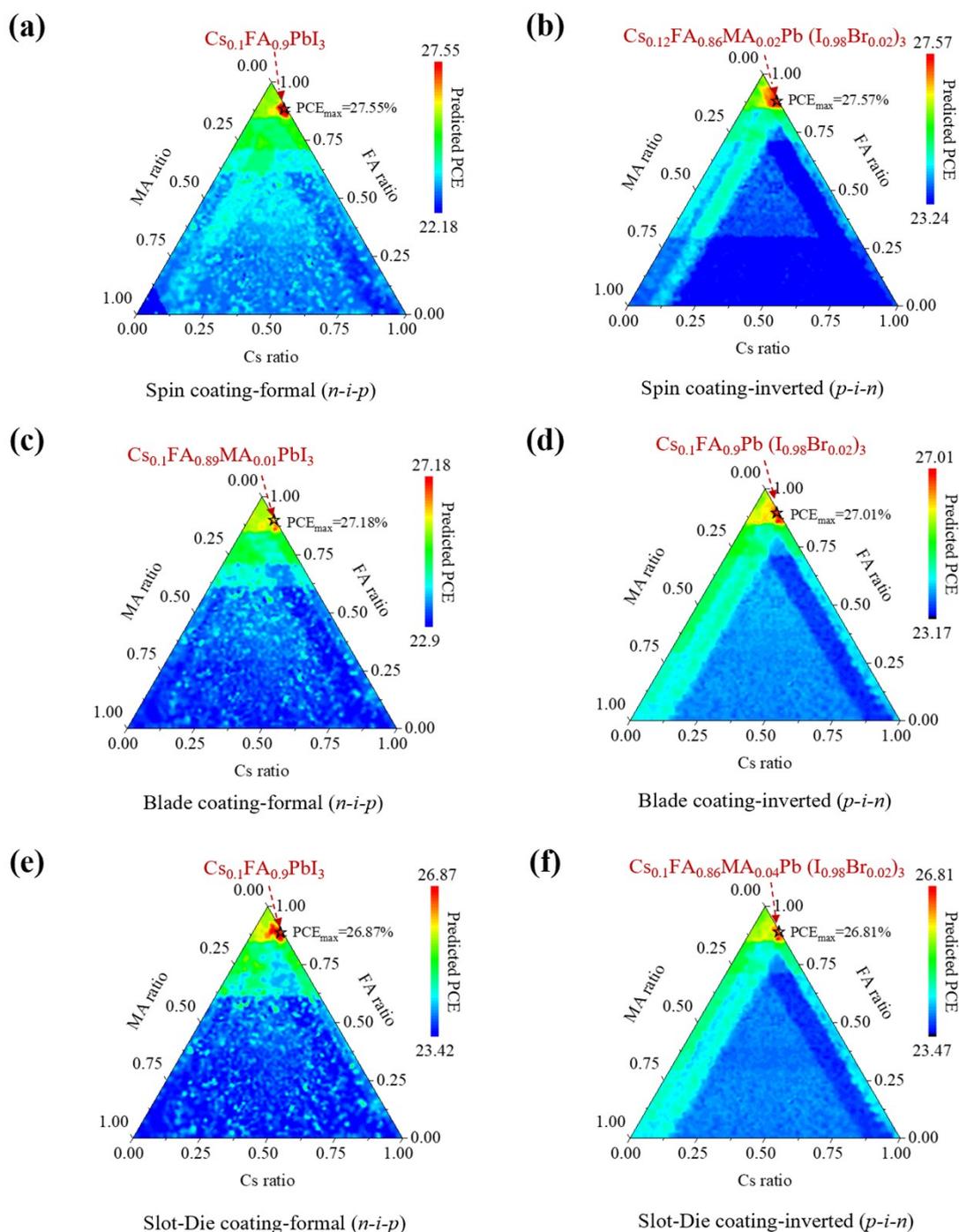
478 **Fig. S15** SHAP analysis diagram of functional layers in perovskite solar cells based on  
 479 the optimal XGBoost model: **(a)** Substrate, **(b)** ETL\_2, **(c)** HTL\_1, **(d)** HTL\_2, **(e)**  
 480 Electrode, **(f)** Surface Passivator, **(g)** Additive\_1, **(h)** Additive\_2.



481

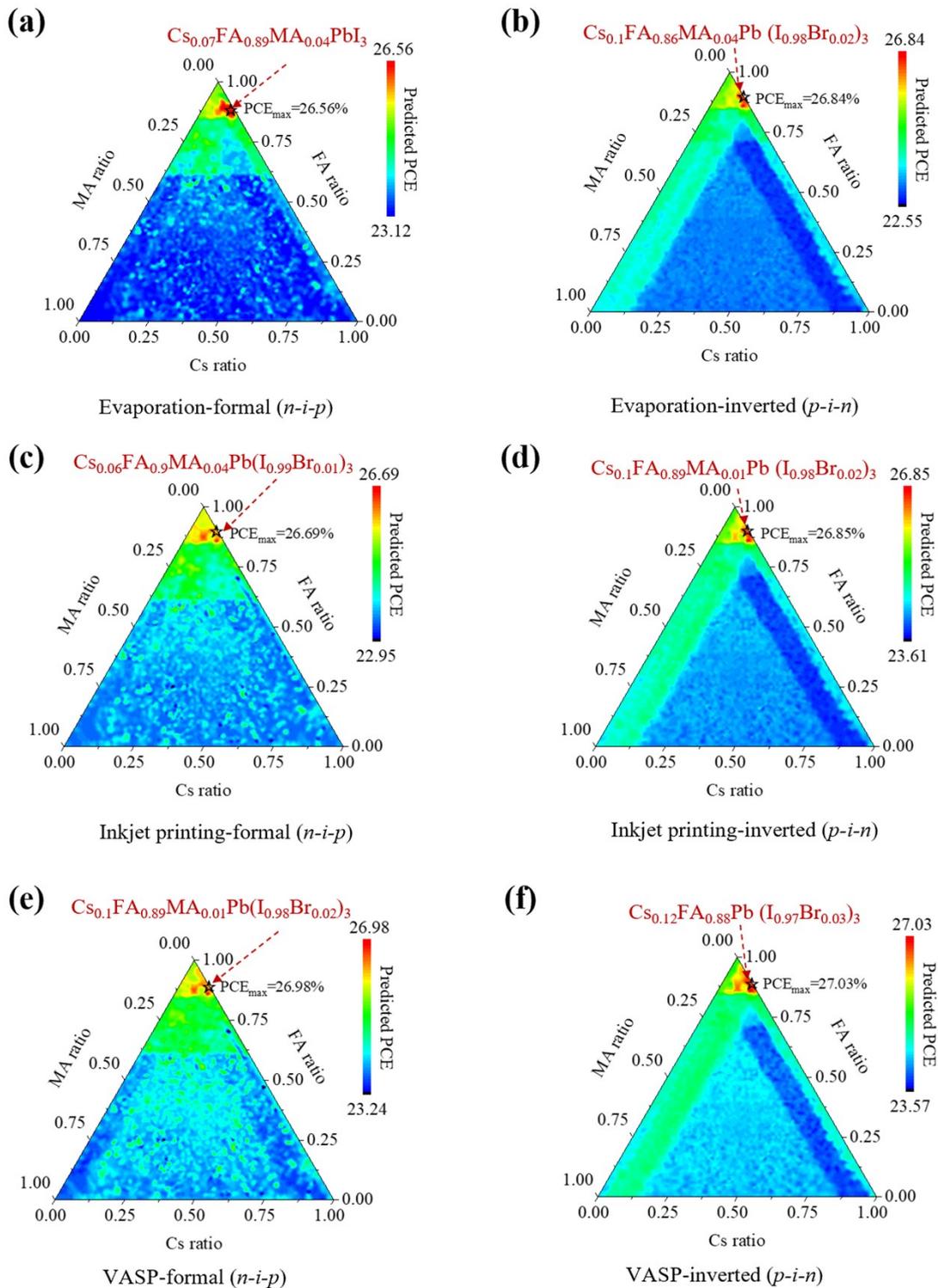
482 **Fig. S16** SHAP analysis diagram of the influence of different anion ratios on PCE based  
 483 on the optimal XGBoost model: **(a)** Br anion, **(b)** I anion, and **(c)** Cl anion.

484



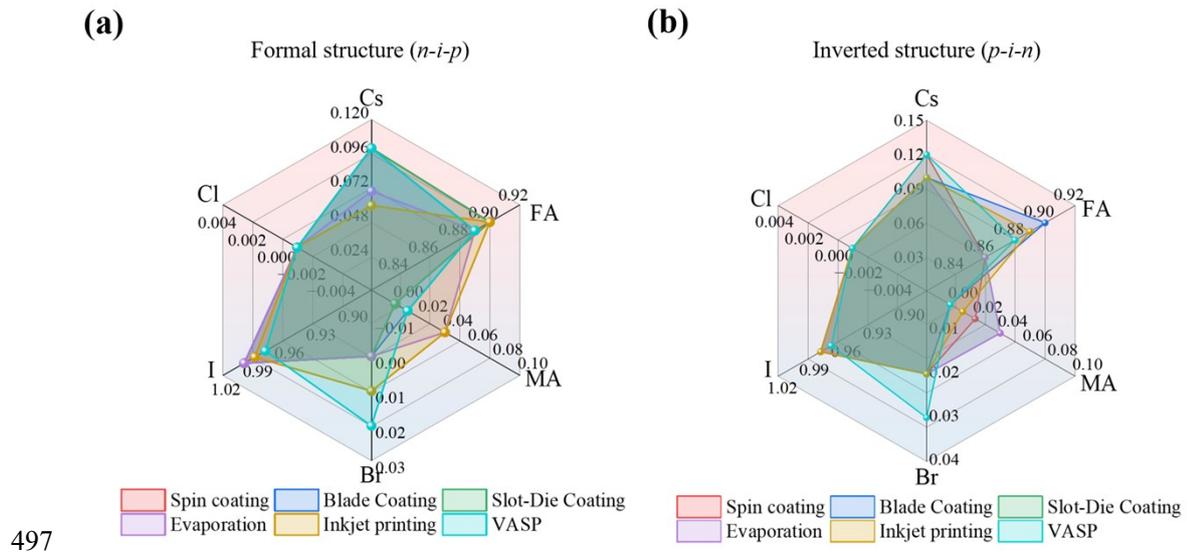
485

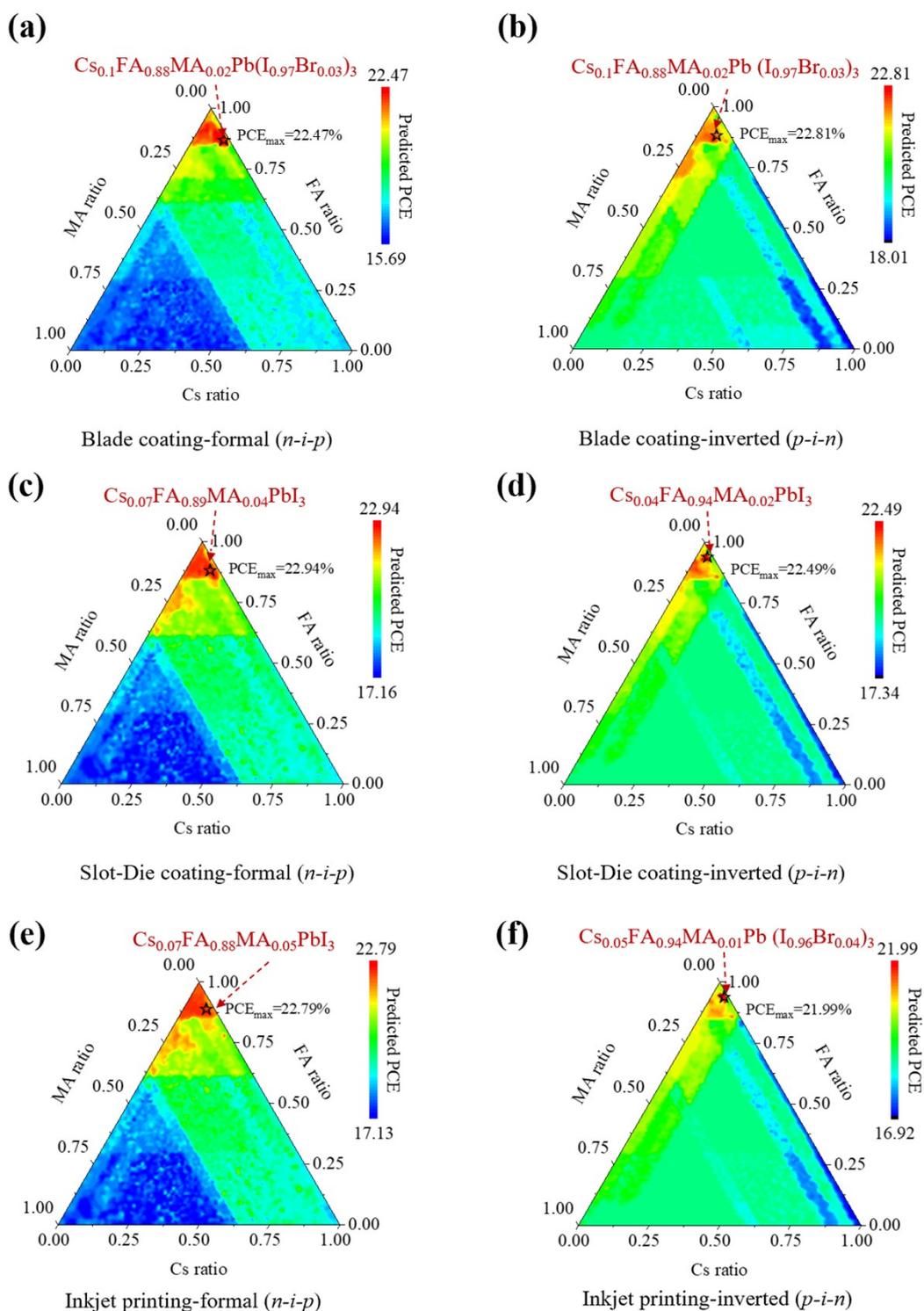
486 **Fig. S17** The triangular heatmap diagram of the effects of Cs, FA, and MA cations on  
 487 PCE for three preparation processes of formal and inverted structures under a small  
 488 area of  $0.03 \text{ cm}^2$  based on the optimal XGBoost model: (a) Spin coating-formal, (b)  
 489 Spin coating-inverted, (c) Blade coating-formal, (d) Blade coating-inverted, (e) Slot-  
 490 Die coating-formal, and (f) Slot-Die coating-inverted.



491

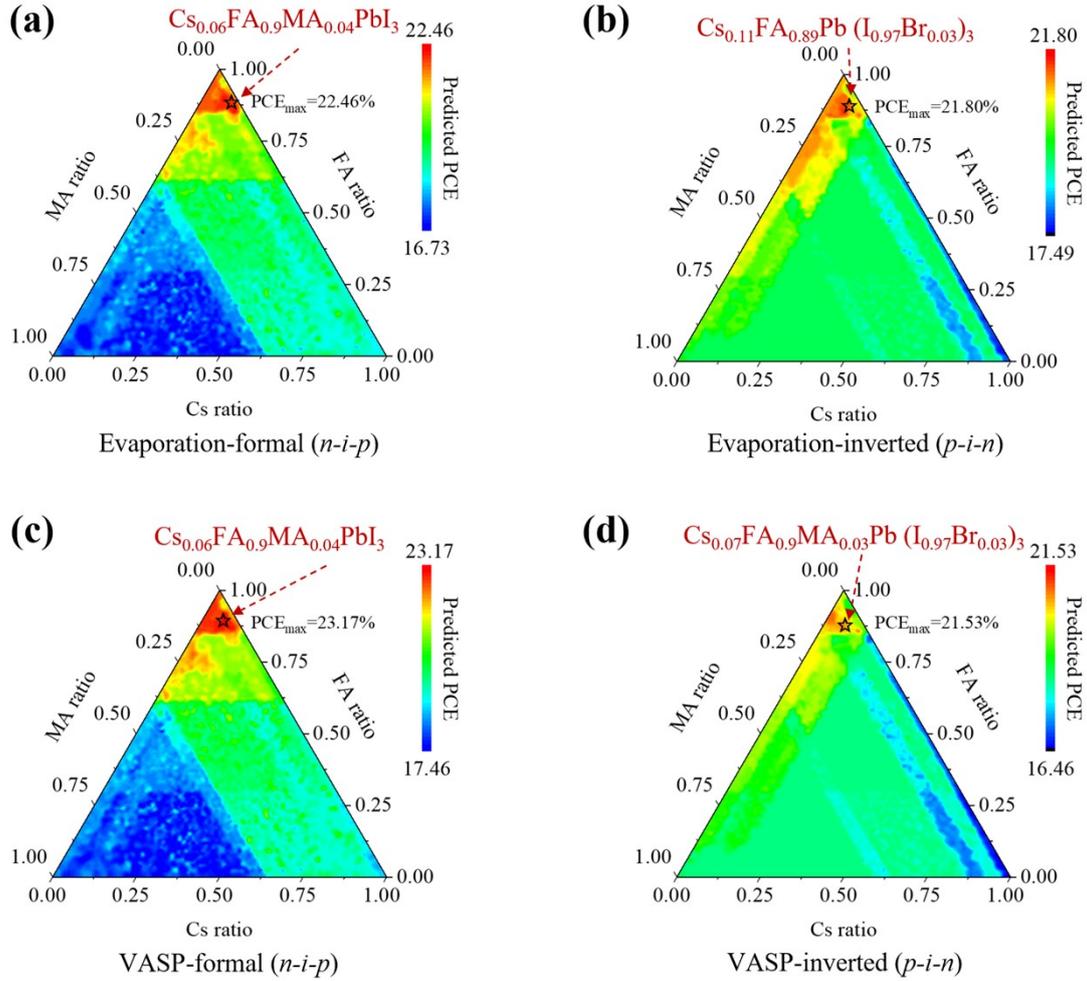
492 **Fig. S18** The triangular heatmap of the effects of Cs, FA, and MA cations on PCE for  
 493 three preparation processes of formal structure and inverted structure under a small area  
 494 of 0.03 cm<sup>2</sup> based on the optimal XGBoost model: (a) Evaporation-formal, (b)  
 495 Evaporation-inverted, (c) Inkjet printing-formal, (d) Inkjet printing-inverted, (e)  
 496 VASP-formal, and (f) VASP-inverted.





502

503 **Fig. S20** The triangular heatmap of the effects of Cs, FA, and MA cations on PCE for  
 504 three preparation processes of formal and inverted structures under a large area of 10  
 505  $\text{cm}^2$  based on the optimal XGBoost model: (a) Blade coating-formal, (b) Blade coating-  
 506 inverted, (c) Slot-Die coating-formal, (d) Slot-Die coating-inverted (e) Inkjet printing-  
 507 formal, and (f) Inkjet printing-inverted.



508

509 **Fig. S21** The ternary phase diagram of the effects of Cs, FA, and MA cations on PCE  
 510 for two preparation processes of formal and inverted structures under a large area of 10  
 511  $\text{cm}^2$  based on the optimal XGBoost model: (a) Evaporation-formal, (b) Evaporation-  
 512 inverted, (c) VASP-formal, and (d) VASP-inverted.

**Table. S1** The 23 input features and corresponding explanations of ML models.

Factors	Description or alternatives for one-hot encoding
ETL_1 VB	VB value of the first ETL layer
ETL_1 CB	CB value of the first ETL layer
ETL_2 VB	VB value of the second ETL layer
ETL_2 CB	CB value of the second ETL layer
HTL_1 VB	VB value of the first HTL layer
HTL_1 CB	CB value of the first HTL layer
HTL_2 VB	VB value of the second HTL layer
HTL_2 CB	CB value of the second HTL layer
Electrode function	Work functions of different types of electrodes
Substrate function	Work functions of different types of substrates
Area	The mask area of the J-V test of the device
Perovskite Eg	Perovskite Eg is given based on empirical formulas
MA	The proportion of MA cations to all cations in perovskite
FA	The proportion of FA cations to all cations in perovskite
Cs	The proportion of Cs cations to all cations in perovskite
Br	The proportion of Br anions to all anions in perovskite
I	The proportion of I anions to all anions in perovskite
Cl	The proportion of Cl anions to all anions in perovskite
Additive_1	Inorganic salts, Organic ammonium salt, others, none
Additive_2	Inorganic salts, Organic ammonium salt, others, none
Surface passivator	Organic ammonium salts, polymer, Inorganic molecules, others, none
Structure	n-i-p, p-i-n
Preparation process	Spin coating, Blade Coating, Slot-Die Coating, Evaporation, Inkjet printing, VASP

**Table S2.** Data distribution of the dataset across different device areas and preparation processes.

Preparation process	Sample size					total
	<1	1 to 15	15 to 60	60 to 100	>100	
Spin Coating	32	44	97	7	17	197
Blade Coating	4	9	23	4	5	45
Slot-Die Coating	8	4	22	2	1	37
Evaporation	4	2	7	1	0	14
Inkjet Printing	5	3	2	1	11	22
VASP	7	5	6	2	2	22
total	60	67	157	17	36	337

**Table S3.** Data distribution ratio before and after weighted processing of different preparation processes.

Preparation process	Spin coating	Blade Coating	Slot-Die Coating	Evaporation	Inkjet printing	VASP
Unweighted processing	58.5%	13.4%	11.0%	4.2%	6.5%	6.5%
Weighted processing	18.0%	18.0%	16.5%	12.9%	16.5%	18.0%

**Table. S4** Training parameters for 6 models used for PCE prediction.

ML algorithms	Training parameters
LightGBM	'regressor__n_estimators': 300, 'regressor__max_depth': 6, 'regressor__learning_rate': 0.05, 'regressor__subsample': 0.8, 'regressor__colsample_bytree': 0.8
SVR	'regressor__C': 100, 'regressor__gamma': 'scale', 'regressor__kernel': 'rbf'
RF	'regressor__n_estimators': 300, 'regressor__max_depth': 15, 'regressor__min_samples_split': 5, 'regressor__min_samples_leaf': 2, 'regressor__max_features': 'sqrt', 'regressor__bootstrap': True, 'regressor__random_state': 42,
MLP	'regressor__activation': 'relu', 'regressor__alpha': 0.001, 'regressor__max_iter': 1000, 'regressor__learning_rate_init': '0.001', 'regressor__solver': 'adam', 'regressor__early_stopping': True
AdaBoost	'regressor__max_depth': 4, 'regressor__min_samples_split': 5, 'regressor__min_samples_leaf': 2
XGBoost	'regressor__colsample_bytree': 0.8, 'regressor__gamma': 0.1, 'regressor__learning_rate': 0.05, 'regressor__max_depth': 6, 'regressor__n_estimators': 300, 'regressor__subsample': 0.8, 'regressor__min_child_weight': 3

**Table. S5** Training metrics of 6 ML models on training and testing sets.

ML algorithms	Train set			Test set		
	$R^2$	RMSE(%)	MAPE(%)	$R^2$	RMSE(%)	MAPE(%)
XGBoost	0.9475	1.0678	5.4055	0.8151	1.9699	9.787
RF	0.8667	1.5102	8.3044	0.63734	2.5585	14.3297
AdaBoost	0.8483	1.8154	10.9258	0.7661	2.2157	13.2769
SVR	0.882	1.6007	11.4056	0.7303	2.3789	14.5416
LightGBM	0.9235	1.2892	6.0133	0.7867	2.1158	11.069
MLP	0.6834	2.6226	15.1249	0.5436	3.0949	18.2154

**Table. S6** The optimal perovskite components and predicted PCE for six preparation processes under small area (0.03 cm<sup>2</sup>) formal and inverted structures obtained from the triangle heatmap.

Structure	Preparation process	Cs	FA	MA	Br	I	Cl	predicted PCE(%)
Formal ( <i>n-i-p</i> )	Spin coating	0.1	0.9	0.0	0	1	0	27.55
	Blade Coating	0.1	0.89	0.01	0	1	0	27.18
	Slot-Die Coating	0.1	0.9	0.0	0	1	0	26.87
	Evaporation	0.07	0.89	0.04	0	1	0	26.56
	Inkjet printing	0.06	0.9	0.04	0.01	0.99	0	26.69
	VASP	0.1	0.89	0.01	0.02	0.98	0	26.98
Inverted ( <i>p-i-n</i> )	Spin coating	0.12	0.86	0.02	0.02	0.98	0	27.57
	Blade Coating	0.1	0.9	0	0.02	0.98	0	27.01
	Slot-Die Coating	0.1	0.86	0.04	0.02	0.98	0	26.81
	Evaporation	0.1	0.86	0.04	0.02	0.98	0	26.84
	Inkjet printing	0.1	0.89	0.01	0.02	0.98	0	26.85
	VASP	0.12	0.88	0	0.03	0.97	0	27.03

**Table. S7** The optimal perovskite components and PCE for six preparation processes under large area (10 cm<sup>2</sup>) formal and inverted structures obtained from the triangle heatmap.

Structure	Preparation process	Cs	FA	MA	Br	I	Cl	PCE(%)
Formal ( <i>n-i-p</i> )	Spin coating	0.08	0.89	0.03	0	1	0	22.99
	Blade Coating	0.1	0.88	0.02	0.03	0.97	0	22.47
	Slot-Die Coating	0.07	0.89	0.04	0	1	0	22.94
	Evaporation	0.06	0.9	0.04	0	1	0	22.46
	Inkjet printing	0.07	0.88	0.05	0	1	0	22.79
	VASP	0.06	0.9	0.04	0	1	0	23.17
Inverted ( <i>p-i-n</i> )	Spin coating	0.07	0.89	0.04	0	1	0	22.70
	Blade Coating	0.1	0.88	0.02	0.03	0.97	0	22.81
	Slot-Die Coating	0.04	0.94	0.02	0.04	0.96	0	22.49
	Evaporation	0.11	0.89	0	0.03	0.97	0	21.80
	Inkjet printing	0.05	0.94	0.01	0.04	0.96	0	21.99
	VASP	0.07	0.9	0.03	0.01	0.99	0	21.53

**Table S8.** Ten additional data points were collected to verify the stability of the model predictions.

Number Factor	1	2	3	4	5	6	7	8	9	10
Voc (V)	1.112	1.0725	1.064	1.14	1.154	1.16	0.993	1.08	1.137	1.172
Jsc (mA/cm <sup>2</sup> )	24.63	19.6	20.8	25.85	25.09	25.02	26.8	20.54	25.39	25.53
FF (%)	77.82	78.31	76	85	77.7	84.69	72.1	78.49	80.75	84.44
Area (cm <sup>2</sup> )	25	16	12.5	0.07	6.25	0.07	14	16	0.09	0.09
Substrate	FTO	ITO	FTO	ITO	ITO	ITO	FTO	ITO	ITO	ITO
Substrate function (eV)	4.5	4.7	4.5	4.7	4.7	4.7	4.5	4.7	4.7	4.7
Cs	0	0.05	0.05	0	0	0.07	0.05	0.05	0.05	0.05
FA	1	0.85	0.95	0.92	1	0.93	0.93	0.95	0.95	0.95
MA	0	0.05	0	0.08	0	0	0.02	0	0	0
Br	0	0.15	0	0	0	0	0.15	0	0	0
I	3	2.85	3	3	3	3	2.585	3	3	3
Cl	0	0	0	0	0	0	0	0	0	0
perovskite Eg (eV)	1.48	1.518	1.492	1.486	1.48	1.498	1.518	1.492	1.492	1.492
ETL_1	PCBM	SnO <sub>2</sub>	SnO <sub>2</sub>	C <sub>60</sub>						
ETL_1 CB (eV)	3.8	4.5	4.5	3.9	3.9	3.9	3.9	3.9	3.9	3.9
ETL_1 VB (eV)	6.1	7.5	7.5	5.8	5.8	5.8	5.8	5.8	5.8	5.8
ETL_2	BCP	none	none	BCP						
ETL_2 CB	2.8	0	0	2.8	2.8	2.8	2.8	2.8	2.8	2.8

(eV)										
ETL_2 VB (eV)	6.7	0	0	6.7	6.7	6.7	6.7	6.7	6.7	6.7
HTL_1	NiO <sub>x</sub>	NiO <sub>x</sub>	Spiro-OMeTAD	TP-TPA	NiO <sub>x</sub>	PTAA	NiO <sub>x</sub>	NiO <sub>x</sub>	NiO <sub>x</sub>	NiO <sub>x</sub>
HTL_1 CB (eV)	1.75	1.75	2.2	2.1	1.75	2.1	1.75	1.75	1.75	1.75
HTL_1 VB (eV)	5.1	5.1	5.22	5	5.1	5.3	5.1	5.1	5.1	5.1
HTL_2	MeO-2PACz	NiO <sub>x</sub>	none	none	MeO-4PACz	none	none	none	MeO-4PACz	MeO-4PACz
HTL_2 CB (eV)	2.3	1.75	0	0	2.35	0	0	0	2.35	2.35
HTL_2 VB (eV)	5.3	5.1	0	0	5.35	0	0	0	5.35	5.35
Electrode	Ag	Ag	Au	Ag	Ag	Cu	Cu	Cu	Ag	Ag
Electrode function (eV)	4.3	4.3	5.28	4.3	4.3	4.8	4.8	4.8	4.3	4.3
additive_1	others	Organic ammonium salt	Organic ammonium salt	none	none	others	Organic ammonium salt	Organic ammonium salt	none	Organic ammonium salt
additive_2	none	none	none	none	none	none	none	none	none	none
Surface passivator	Organic ammonium salts	others	polymer	none	Organic ammonium salts	Organic ammonium salts				
Structure	p-i-n	p-i-n	n-i-p	p-i-n	p-i-n	p-i-n	p-i-n	p-i-n	p-i-n	p-i-n
Preparation	Spin coating	Spin coating	Blade Coating	Spin coating	Blade Coating	Slot-Die	Slot-Die	Evaporation	Spin coating	Spin coating

process						Coating	Coating			
DOI	10.1002/anie.2 02516464	10.1002/smt.d. 202501325	10.1039/d5el0 0032g	10.1016/j.cej.2 025.168738	10.1039/d5ee0 4570c	10.1016/j. jechem.20 25.08.028	10.1016/j.sol mat.2025.113	10.1016/j.sol mat.2025.113	10.1016/j.jech em.2025.05.03	10.1016/j.jech em.2025.05.0
Experimental efficiency (%)	21.31	16.46	18.46	24.8	21.3	24.58	19.19	17.58	23.32	25.27
Predicted efficiency (%)	20.67	17.17	17.25	23.41	21.22	23.97	19.96	17.97	23.62	23.63
<b>relative error (%)</b>	<b>2.988</b>	<b>4.333</b>	<b>6.578</b>	<b>5.625</b>	<b>0.368</b>	<b>2.470</b>	<b>3.992</b>	<b>2.237</b>	<b>1.299</b>	<b>6.500</b>

**Table. S9** Decay data of PCE with area variation under the optimal perovskite composition of six preparation processes for formal structure and inverted structure obtained in a small area.

Structure	Area (cm <sup>2</sup> )	Predicted PCE (%)					
		Spin coating	Blade Coating	Slot-Die Coating	Evaporation	Inkjet printing	VASP
Formal ( <i>n-i-p</i> )	0.03	27.55	27.18	26.87	26.56	26.69	26.98
	10.03	22.44	22.39	22.51	21.81	22.09	22.49
	20.03	19.08	18.87	18.68	18.71	18.91	19.14
	30.03	18.09	17.76	18.49	18.37	18.73	18.77
	40.03	17.43	17.67	18.17	18.04	18.62	18.36
	50.03	17.35	17.58	17.83	18	18.47	18.01
	60.03	17.18	17.5	17.61	17.53	18.21	17.73
	70.03	17.14	17.42	17.49	17.39	18.05	17.54
	80.03	17.09	17.34	17.37	17.26	17.69	17.47
	90.03	17.06	17.27	17.33	17.23	17.45	17.42
	100.03	17.04	17.2	17.32	17.21	17.41	17.38
Inverted ( <i>p-i-n</i> )	0.03	27.57	27.01	26.81	26.84	26.85	27.03
	10.03	22.61	22.79	22.49	21.79	21.99	21.49
	20.03	20.09	20.68	20.46	19.57	19.78	19.19
	30.03	18.62	19.76	19.51	18.59	18.76	18.25
	40.03	17.89	19.29	19.06	18.16	18.3	17.86
	50.03	17.52	19.06	18.84	17.98	18.09	17.75
	60.03	17.32	18.94	18.75	17.88	17.99	17.67
	70.03	17.21	18.89	18.71	17.85	17.95	17.64
	80.03	17.15	18.86	18.69	17.84	17.93	17.62
	90.03	17.13	18.85	18.68	17.83	17.92	17.6
	100.03	17.11	18.84	18.67	17.82	17.91	17.59