Supplementary material

Machine learning facilitates the development of interconnecting layers for perovskite/silicon heterojunction tandem solar cells with proof-of-concept efficiency > 38%

Xuejiao Wang^{1,2,3,4,5,\xi}, Guanlan Chen ^{1,2,3,4,5,ξ}, Ying Liu^{1,2,3,4,5}, Guangyi Wang^{1,2,3,4,5}, Wei Han^{1,2,3,4,5}, Jin Wang^{1,2,3,4,5}, Pengfei Liu^{1,2,3,4,5}, Jilei Wang⁶, Shaojuan Bao⁶, Bo Yu⁷, Ying Liu⁷, Xinliang Chen^{1,2,3,4,5}, Shengzhi Xu^{1,2,3,4,5}, Ying Zhao ^{1,2,3,4,5,u} and Xiaodan Zhang ^{1,2,3,4,5,u}

¹Institute of Photoelectronic Thin Film Devices and Technology, Nankai University, Tianjin 300350, P.

R. China

²Tianjin Key Laboratory of Efficient Utilization of Solar Energy, Tianjin 300350, P. R. China

³State Key Laboratory of Photovoltaic Materials and Cells, Tianjin 300192, P. R. China

⁴Engineering Research Center of Thin Film Photoelectronic Technology of Ministry of Education,

Tianjin 300350, P. R. China

⁵Collaborative Innovation Center of Chemical Science and Engineering (Tianjin), Tianjin 300072, P. R.

China

⁶Anhui Soltrend New Energy Technology Co., Ltd, Anhui 230000, P. R. China

⁷Yingli Energy Development Co., Ltd, Baoding 071002, P. R. China

 $[\]xi$ These authors contributed equally to this work

u Correspondences to Xiaodan Zhang, E-mail: <u>xdzhang@nankai.edu.cn</u>; Ying Zhao: zhaoygds@nankai.edu.cn;

Introduction to machine learning

Machine learning (ML) is a field of study that gives computers the ability to learn without being explicitly programmed, this is Arthur Samuel's definition of ML. With the rapid development of Artificial intelligence (AI) technology, ML as a powerful computer tool has been widely used in all walks of life. The main task of ML is to generate a certain "model" from the data through the computer with the help of an algorithm, which is called an ML algorithm. Through ML algorithms, as long as the data obtained from experience is learned, a certain model will be obtained. In the actual situation, new data will be generated constantly, and at this time, as long as the new data is substituted into the obtained model, the new data can be analyzed, predicted and judged. The two main types of ML are supervised learning and unsupervised learning. Generally speaking, supervised learning is "teaching machine learning", and unsupervised learning is "machine learning by itself". Supervised learning is mainly for regression issues and classification issues, and clustering is a typical representative of unsupervised learning. This paper mainly adopts the way of supervised learning to carry out machine learning.

Pearson correlation coefficient (r)

The Pearson correlation coefficient (r) is a statistical measure used to measure the degree of linear correlation between two variables. It has a value between -1 and 1 and is commonly used to evaluate the correlation between two continuous variables.

$$r = \frac{\sum (x_i - \overline{x})(y_i - \overline{y})}{\sqrt{\sum (x_i - \overline{x})^2 \sum (y_i - \overline{y})^2}}$$

Where x_i and y_i are the ith observations of the variables X and Y. \bar{x} and \bar{y} are the mean of X and Y. The Pearson correlation coefficient r ranges from -1 to 1. An r > 0 indicates that the mean of the continuous variable is higher when the binary variable equals 1, demonstrating a positive correlation. Conversely, an r < 0 suggests that the mean of the continuous variable is lower when the binary variable equals 1, indicating a negative correlation. The closer r is to 1, the stronger the correlation; the closer it is to 0, the weaker the correlation.

Point-biserial correlation coefficient

The Point-biserial correlation coefficient is a measure of the correlation between a continuous variable and a binary variable. It is a special case of Pearson correlation coefficient, which is suitable for the analysis of the relationship between discrete and continuous variables of binary classification (0/1).

$$r = \frac{\overline{X}_1 - \overline{X}_0}{s} \cdot \sqrt{\frac{n_1 - n_0}{n^2}}$$

Where \overline{X}_1 and \overline{X}_0 are the means of the continuous variables when the binary variables are 1 and 0, respectively, s is the population standard deviation of a continuous variable, n_1 and n_0 are the number of samples of 1 and 0 in a binary categorical variable and n is the total number of samples ($n = n_1 + n_0$). The point-biserial correlation coefficient r ranges from -1 to 1. An r > 0 indicates that the mean of the continuous variable is higher when the binary variable equals 1, demonstrating a positive correlation. Conversely, an r < 0 suggests that the mean of the continuous variable is lower when the binary variable equals 1, indicating a negative correlation. The closer r is to 1, the stronger the correlation; the closer it is to 0, the weaker the correlation.

The p-value reflects the statistical significance of the correlation coefficient:

- p < 0.05: The correlation is considered significant, implying it is unlikely to have occurred by chance.
- $p \ge 0.05$: The correlation may not be significant, suggesting it could be random.

One-way ANOVA

One-way ANOVA is a statistical method used to determine whether there are significant differences between the means of three or more independent groups. It compares the variance between groups to the variance within groups to evaluate whether the grouping factor affects the dependent variable. The ratio of these variances is called the F-statistic, calculated as:

$$F = \frac{\text{Between} - \text{Group Variance}}{\text{Within} - \text{Group Variance}}$$

The P-value represents the probability of observing the current F-value or a more

extreme result under the null hypothesis (that is, assuming that all groups have equal means). $P \le 0.05$ means that at least one group has a significantly different mean from the others, while P > 0.05 was considered to have no significant difference in the mean values of all groups.

F-Statistic:

- Represents the ratio of variance between groups to variance within groups.
- A higher F-value indicates greater differences between group means relative to random noise.

p-Value:

- If p < 0.05, reject the null hypothesis. This means at least one group mean is significantly different.
- If $p \ge 0.05$, fail to reject the null hypothesis, indicating no significant difference between group means.

Lasso regression model

Lasso regression, or Least Absolute Shrinkage and Selection Operator Regression, is a type of linear regression that performs both variable selection and regularization to improve the accuracy and interpretability of the statistical model it produces. There are three key characteristics of Lasso regression, the first one is regularization, which adds a penalty term to the loss function to prevent overfitting and handle multicollinearity. The second one is featuring selection, which shrinks coefficients of less important features to exactly zero, effectively removing them from the model and the last one is penalty term, which is based on the L1 norm of the coefficients, encouraging sparsity in the model. The objective function for Lasso regression is:

Minimize:
$$\frac{1}{2n}\sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$

Where y_i is actual value, \hat{y}_i is predicted value, β_j is coefficients of the features and λ is regularization strength (controls the degree shrinkage).

Random forest model

Random Forest (RF) is a versatile ensemble learning method primarily used for classification and regression tasks. It is based on constructing a multitude of decision

trees during training and combining their outputs (majority vote for classification or average for regression) to improve predictive performance and reduce overfitting. There are also three key features of RF. The first one is ensemble of decision tress. A random forest is essentially a collection of decision trees, each tree is trained on a randomly sampled subset of the data (with replacement, i.e., bootstrapping). The second one is featuring randomness, which means during training, only a random subset of features is considered at each split in the tree, this helps reduce correlation among trees and increases diversity. The third one is Aggregation: for classification, the final prediction is made by majority voting among the trees. For regression, the output is the average of the predictions from all trees.

The RF algorithm consists of an ensemble of decision trees. During the tree-building process, each tree selects features to split the data in a way that maximizes the improvement in data purity. Feature importance is calculated based on the average contribution of each variable to the reduction in impurity (e.g., root mean squared error) across all trees. When a feature is frequently used as a splitting variable and significantly reduces prediction error, its importance score increases. In RF, the contributions of each feature across all trees are aggregated and normalized. The resulting values represent the normalized importance of each feature, typically summing to one.

Neural network model

Neural network (NN) model is a machine learning method inspired by the structure and functioning of biological neural networks. It is designed to identify complex patterns and relationships in data by using a layered structure composed of interconnected nodes (neurons). Neural networks are widely used in tasks like regression, classification, image recognition, natural language processing, and more. Its basic structure can be divided into the following five parts, as listed in Table S1.

Table S1. Five parts of NN.

| Structure | Explanation |
|-------------|--|
| Input layer | Accepts the input data, with each feature of the data corresponding to one neuron in this layer. |

| | Does not perform computations—simply passes the data forward. |
|---------------|--|
| | Consist of one or more layers of neurons between |
| | the input and output layers. |
| | Each neuron applies a mathematical function |
| Hidden layers | (activation function) to the inputs and passes the |
| | results to the next layer. |
| | These layers extract and learn features from the |
| | data. |
| | Provides the final predictions or classifications. |
| 0 1 | For regression, the output is usually a single continuous value. |
| Output layer | For classification, the output layer size |
| | corresponds to the number of classes, often with |
| | a SoftMax function for probabilities. |
| | Every neuron is connected to the neurons in |
| Connections | adjacent layers through weights. Each connection |
| Connections | is assigned a weight that determines its |
| | importance. |
| | An additional parameter added to neurons to shift |
| Bias | the activation function, improving the model's |
| | ability to learn |

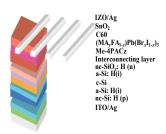


Fig. S1. The schematic diagram of the tandem solar cell structure.

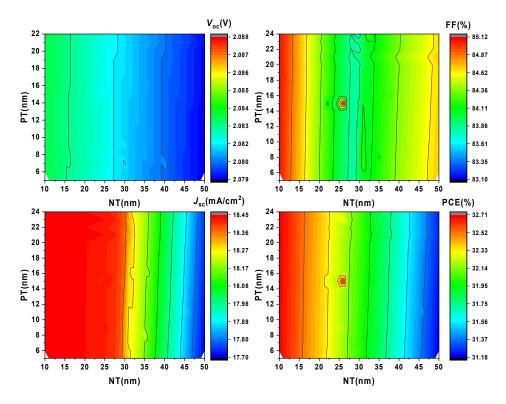


Fig. S2. Simulated device performance as a function of Si-ICL.

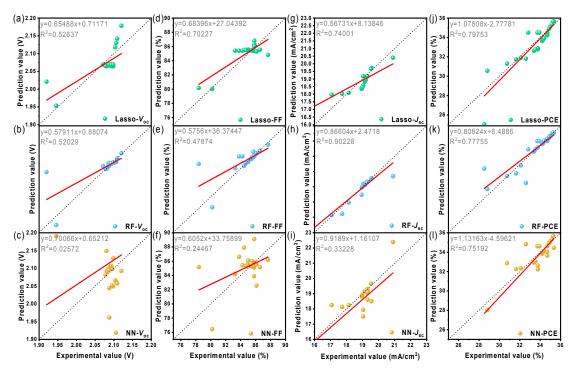


Fig. S3. Tandem solar cell performance predicted by different models based on dataset of TCO-ICL: (a-c) V_{oc} , (d-f) FF, (g-i) J_{sc} , (j-l) PCE. Models used include Lasso, RF and NN. The dashed line represents a perfect fitting, while the solid red line represents the actual fitting result.

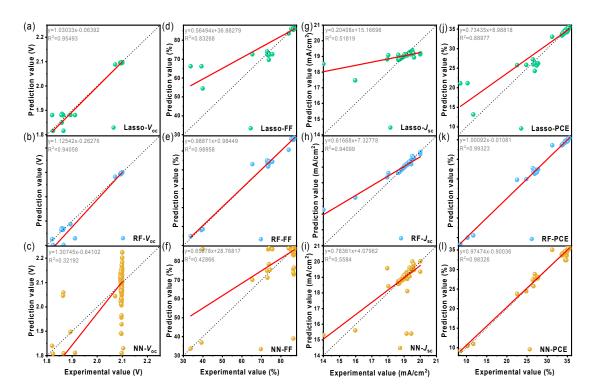


Fig. S4. Tandem solar cell performance predicted by different models based on dataset of Si+TCO-ICL: (a-c) V_{oc} , (d-f) FF, (g-i) J_{sc} , (j-l) PCE. Models used include Lasso, RF and NN. The dashed line represents a perfect fitting, while the solid red line represents the actual fitting result.

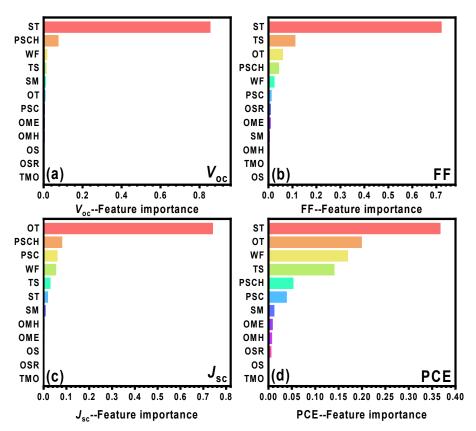


Fig. S5. Feature importance for dataset of TCO-ICL based on the RF model. (a) $V_{\rm oc}$, (b) FF, (c) $J_{\rm sc}$, (d) PCE.

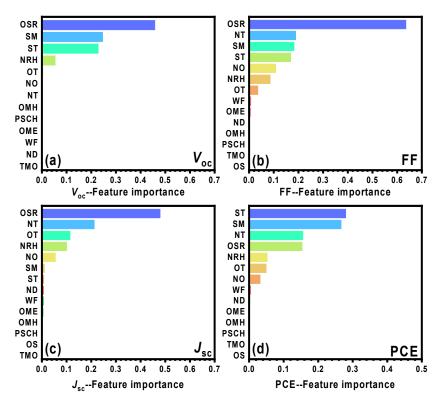


Fig. S6. Feature importance for dataset of Si+TCO-ICL based on the RF model. (a) V_{oc} , (b) FF, (c) J_{sc} , (d) PCE.

Table S2. Point-biserial correlation coefficient results of dataset of Si-ICL.

| Category Variable | Numeric Variables | Correlation coefficient r | P Value |
|-------------------|-------------------|---------------------------|---------------------------|
| SM | $V_{ m oc}$ | -0.9414323515397796 | 9.7483407e ⁻⁹⁹ |
| SM | $J_{ m sc}$ | 0.14571982356195037 | 0.0361671 |
| SM | FF | -0.9125988991043305 | 1.4347159e ⁻⁸¹ |
| SM | PCE | -0.9073200492344362 | 4.4338796e ⁻⁷⁹ |

Table S3. Point-biserial correlation coefficient results of dataset of TCO-ICL.

| Category Variables | Numeric Variables | Correlation coefficient r | P Value |
|--------------------|-------------------|-----------------------------|---------------------|
| SM | $V_{ m oc}$ | -0.579666053 | $4.93e^{-10}$ |
| SM | $J_{ m sc}$ | 0.270775911 | 0.007305854 |
| SM | FF | -0.563437622 | 1.87e ⁻⁹ |
| SM | PCE | -0.375159162 | 1.53e ⁻⁴ |
| PSC | $V_{ m oc}$ | 0.637346795 | $2.23e^{-12}$ |
| PSC | $J_{ m sc}$ | 0.31318726 | $1.79e^{-3}$ |
| PSC | FF | 0.614959877 | $2.07e^{-11}$ |
| PSC | PCE | 0.64375799 | $1.14e^{-12}$ |

Table S4. Point-biserial correlation coefficient results of dataset of Si+TCO-ICL.

| Category Variable | Numeric Variables | Correlation coefficient r | P Value |
|-------------------|-------------------|-----------------------------|----------------------|
| SM | $V_{ m oc}$ | -0.591615133 | 8.82e ⁻⁴⁸ |
| SM | $J_{ m sc}$ | -0.158119715 | 0.00043071 |
| SM | FF | -0.687284807 | $4.91e^{-70}$ |
| SM | PCE | -0.766577847 | 2.59e ⁻⁹⁶ |

Table S5. One-way ANOVA results of dataset of Si+TCO-ICL.

| Categorical Variables | Numerical Variables | F Value | P Value |
|-----------------------|---------------------|----------|----------|
| ТМО | $V_{ m oc}$ | 0.055604 | 0.994236 |
| TMO | $J_{ m sc}$ | 0.027379 | 0.998549 |
| TMO | FF | 0.058606 | 0.993622 |
| TMO | PCE | 0.090014 | 0.985569 |
| PSCH | $V_{ m oc}$ | 0.077154 | 0.998243 |
| PSCH | $J_{ m sc}$ | 0.064096 | 0.998962 |
| PSCH | FF | 0.031515 | 0.999867 |
| PSCH | PCE | 0.059965 | 0.999143 |

 $Table\ S6.\ Details\ of\ relevant\ hyperparameters\ for\ the\ ML\ models\ in\ Si\mbox{-}ICL\ dataset.$

| 37.11 | 0 / / F / | A 1 1 | 3.6 % | T 1 |
|--------|-----------------|-------|----------|--------------------|
| Models | Output Features | Alpha | Max_iter | Tol |
| | $V_{ m oc}$ | 0.001 | 500 | 10 |
| Lasso | $J_{ m sc}$ | 0.01 | 3000 | 0.1 |
| Lasso | FF | 0.1 | 2000 | 0.01 |
| | PCE | 0.001 | 500 | 10 |
| Models | Output Features | Ma | x_depth | N_estimators |
| | $V_{ m oc}$ | | 5 | 400 |
| DE | $J_{ m sc}$ | | 11 | 200 |
| RF | FF | | 11 | 550 |
| | PCE | | 5 | 550 |
| Models | Output Features | Alpha | | Hidden_layer_sizes |
| | $V_{ m oc}$ | 0.001 | | (50,50) |
| NDI | $J_{ m sc}$ | (| 0.001 | (100,100) |
| NN | FF | | 0.01 | (50,50) |
| | PCE | | 0.01 | (100,100) |

 $Table\ S7.\ Details\ of\ relevant\ hyperparameters\ for\ the\ ML\ models\ in\ TCO-ICL\ dataset.$

| | | 1 1 | | |
|--------|-----------------|-------|----------|--------------------|
| Models | Output Features | Alpha | Max_iter | Tol |
| | $V_{ m oc}$ | 0.01 | 8000 | 0.1 |
| Lagge | $J_{ m sc}$ | 0.08 | 8000 | 0.001 |
| Lasso | FF | 0.05 | 8000 | 0.001 |
| | PCE | 0.01 | 8000 | 0.01 |
| Models | Output Features | Ma | x_depth | N_estimators |
| | $V_{ m oc}$ | | 5 | 200 |
| DE | $J_{ m sc}$ | | 8 | 100 |
| RF | FF | | 7 | 400 |
| | PCE | | 13 | 200 |
| Models | Output Features | Alpha | | Hidden_layer_sizes |
| | $V_{ m oc}$ | | 0.1 | (100,100) |
| ND I | $J_{ m sc}$ | 0 | .0001 | (50,50) |
| NN | FF | 0.0 | 000001 | (50,50) |
| | PCE | 0. | 00001 | (100,100) |

Table S8. Details of relevant hyperparameters for the ML models in Si+TCO-ICL dataset.

| Models | Output Features | Alpha | Max_iter | Tol |
|--------|-----------------|----------|----------|--------------------|
| | $V_{ m oc}$ | 0.001 | 8000 | 0.01 |
| Lagge | $J_{ m sc}$ | 0.1 | 8000 | 0.001 |
| Lasso | FF | 1 | 8000 | 0.001 |
| | PCE | 0.1 | 8000 | 0.01 |
| Models | Output Features | Ma | x_depth | N_estimators |
| | $V_{ m oc}$ | | 7 | 500 |
| RF | $J_{ m sc}$ | | 14 | 500 |
| KΓ | FF | | 8 | 500 |
| | PCE | | 9 | 500 |
| Models | Output Features | Alpha | | Hidden_layer_sizes |
| | $V_{ m oc}$ | 0.000001 | | (100,100) |
| NINI | $J_{ m sc}$ | | 0.1 | (100,100) |
| NN | FF | | 0.1 | (50,50) |
| | PCE | (| 0.001 | (50,50) |

Table S9. Performances of different ML model in the prediction of interconnecting layer of the perovskite/silicon heterojunction tandem solar cells in dataset of Si-ICL.

| Photovoltaic Parameters | ML model | RMSE |
|-------------------------|----------|----------------------|
| | Lasso | 0.015446611151380513 |
| $V_{ m oc}$ | RF | 0.015099898358605316 |
| | NN | 0.1602053538591303 |
| | Lasso | 0.286992626872376 |
| $J_{ m sc}$ | RF | 0.27951616820264574 |
| | NN | 1.120907945816788 |
| | Lasso | 0.9567102056018761 |
| FF | RF | 0.8682119604307589 |
| | NN | 3.183064186323462 |
| | Lasso | 0.6176109384517436 |
| PCE | RF | 0.5764471184651523 |
| | NN | 1.9239972952538564 |

Note: RMSE refer to root mean squared error.

Table S10. Performances of different ML model in the prediction of interconnecting layer of the perovskite/silicon heterojunction tandem solar cells in dataset of TCO-ICL.

| Photovoltaic Parameters | ML model | RMSE |
|-------------------------|----------|----------------------|
| | Lasso | 0.042788177574249193 |
| $V_{ m oc}$ | RF | 0.020475879062268873 |
| | NN | 0.12167004946699936 |
| | Lasso | 0.6108993088784077 |
| $J_{ m sc}$ | RF | 0.36550811937001726 |
| | NN | 1.2081575596461083 |
| | Lasso | 1.8687498275760739 |
| FF | RF | 1.3873067486307413 |
| | NN | 4.827784352395491 |
| | Lasso | 1.2841844521563393 |
| PCE | RF | 1.2143686975396648 |
| | NN | 1.9579579108977 |

Note: RMSE refer to root mean squared error.

Table S11. Performances of different ML model in the prediction of interconnecting layer of the perovskite/silicon heterojunction tandem solar cells in dataset of Si+TCO-ICL.

| Photovoltaic Parameters | ML model | RMSE |
|-------------------------|----------|----------------------|
| | Lasso | 0.059321347021387714 |
| $V_{ m oc}$ | RF | 0.0544010083819987 |
| | NN | 0.17480281827306443 |
| | Lasso | 0.9267432287180262 |
| $J_{ m sc}$ | RF | 0.7116399034301324 |
| | NN | 1.271500876890284 |
| | Lasso | 5.419571813587112 |
| FF | RF | 4.16788514484323 |
| | NN | 5.552792000546463 |
| PCE | Lasso | 2.2687260161718403 |
| | RF | 1.7919828243834686 |
| | NN | 2.4296643153898856 |

Note: RMSE refer to root mean squared error.