

# 1 Input columns used in data cleaning

- Cell\_area\_measured
- ETL\_stack\_sequence
- HTL\_stack\_sequence
- Substrate\_stack\_sequence
- Backcontact\_stack\_sequence
- Backcontact\_thickness\_list
- Perovskite\_composition\_a\_ions
- Perovskite\_composition\_a\_ions\_coefficients
- Perovskite\_composition\_b\_ions
- Perovskite\_composition\_b\_ions\_coefficients
- Perovskite\_composition\_c\_ions
- Perovskite\_composition\_c\_ions\_coefficients
- Perovskite\_additives\_compounds
- Perovskite\_single\_crystal
- Perovskite\_composition\_perovskite\_ABC3\_structure
- Perovskite\_composition\_perovskite\_inspired\_structure
- Perovskite\_band\_gap
- Perovskite\_deposition\_quenching\_induced\_crystallisation
- Perovskite\_deposition\_quenching\_media
- Perovskite\_deposition\_solvent\_annealing
- Perovskite\_deposition\_number\_of\_deposition\_steps
- Perovskite\_deposition\_solvents
- Perovskite\_deposition\_solvents\_mixing\_ratios
- Perovskite\_deposition\_thermal\_annealing\_temperature
- Perovskite\_deposition\_thermal\_annealing\_time
- JV\_reverse\_scan\_Jsc

## 2 Columns removed for duplicate check

- Ref\_ID\_temp
- Ref\_ID
- Ref\_name\_of\_person\_entering\_the\_data
- Ref\_date\_entered\_by\_author

## 3 Device configurations for distribution plots

Distribution A

| <i>Column</i>                                    | <i>Value</i> | <i>Unit</i>       |
|--|--------------|-------------------|
| a_MA_L0  | 1.0          |                   |
| b_Pb_L0  | 1.0          |                   |
| c_I_L0   | 3.0          |                   |
| bandgap_L0                                       | 1.6          | eV                |
| Back_contact_Ag_L0                               | 80.0         | nm                |
| Depo_solvent_GBL_L0                              | 0.7          |                   |
| Depo_solvent_DMSO_L0                             | 0.3          |                   |
| Temp_100_L0                                      | 10.0         | min               |
| Substrate_stack_SLG_L0                           | 1.0          |                   |
| Substrate_stack_FTO_L1                           | 1.0          |                   |
| ETL_stack_sequence_PCBM-60_L0                    | 1.0          |                   |
| ETL_stack_sequence_BCP_L1                        | 1.0          |                   |
| HTL_stack_sequence_NiMgLiO_L0                    | 1.0          |                   |
| Quenching_media_Toluene                          | 1.0          |                   |
| Cell_area_measured                               | 0.09         | cm <sup>2</sup>   |
| Perovskite_quenching_induced_crystallisation     | 1.0          |                   |
| Perovskite_deposition_number_of_steps            | 1.0          |                   |
| Perovskite_composition_perovskite_ABC3_structure | 1.0          |                   |
| JV_reverse_scan_Jsc                              | 18.32        | $\frac{mA}{cm^2}$ |

### Distribution B

| <i>Column</i>                                    | <i>Value</i> | <i>Unit</i>       |
|--|--------------|-------------------|
| a_MA_L0  | 1.0          |                   |
| b_Pb_L0  | 1.0          |                   |
| c_IL0  | 3.0          |                   |
| bandgap_L0                                       | 1.6          | eV                |
| Back_contact_Au_L0                               | 80.0         | nm                |
| Depo_solvent_DMF_L0                              | 08           |                   |
| Depo_solvent_DMSO_L0                             | 0.2          |                   |
| Temp_100_L0                                      | 10.0         | min               |
| Substrate_stack_SLG_L0                           | 1.0          |                   |
| Substrate_stack_FTO_L1                           | 1.0          |                   |
| ETL_stack_sequence_TiO2-c_L0                     | 1.0          |                   |
| ETL_stack_sequence_TiO2-mp_L1                    | 1.0          |                   |
| HTL_stack_sequence_Spiro-MeOTAD_L0               | 1.0          |                   |
| Quenching_media_Ethyl.acetate                    | 1.0          |                   |
| Cell_area_measured                               | 0.1          | cm <sup>2</sup>   |
| Perovskite_quenching_induced_crystallisation     | 1.0          |                   |
| Perovskite_deposition_number_of_steps            | 1.0          |                   |
| Perovskite_composition_perovskite_ABC3_structure | 1.0          |                   |
| JV_reverse_scan_Jsc                              | 21.12        | $\frac{mA}{cm^2}$ |

### Distribution C

| <i>Column</i>                                    | <i>Value</i> | <i>Unit</i>       |
|--|--------------|-------------------|
| a_MA_L0  | 1.0          |                   |
| b_Pb_L0  | 1.0          |                   |
| c_IL0  | 3.0          |                   |
| bandgap_L0                                       | 1.6          | eV                |
| Back_contact_Au_L0                               | 80.0         | nm                |
| Depo_solvent_DMSO_L0                             | 1.0          |                   |
| Temp_100_L0                                      | 20.0         | min               |
| Substrate_stack_SLG_L0                           | 1.0          |                   |
| Substrate_stack_FTO_L1                           | 1.0          |                   |
| ETL_stack_sequence_TiO2-c_L0                     | 1.0          |                   |
| ETL_stack_sequence_TiO2-mp_L1                    | 1.0          |                   |
| HTL_stack_sequence_PTAA_L0                       | 1.0          |                   |
| Quenching_media_Chlorobenzene                    | 1.0          |                   |
| Cell_area_measured                               | 0.16         | cm <sup>2</sup>   |
| Perovskite_quenching_induced_crystallisation     | 1.0          |                   |
| Perovskite_deposition_number_of_steps            | 1.0          |                   |
| Perovskite_composition_perovskite_ABC3_structure | 1.0          |                   |
| JV_reverse_scan_Jsc                              | 22.41        | $\frac{mA}{cm^2}$ |

**Distribution D**

| <i>Column</i>                                    | <i>Value</i> | <i>Unit</i>       |
|--|--------------|-------------------|
| a_MA_L0  | 1.0          |                   |
| b_Pb_L0  | 1.0          |                   |
| c_I_L0   | 3.0          |                   |
| bandgap_L0                                       | 1.6          | eV                |
| Back_contact_Au_L0                               | 60.0         | nm                |
| Depo_solvent_IPA_L0                              | 1.0          |                   |
| Temp_70_L0                                       | 30.0         | min               |
| Substrate_stack_SLG_L0                           | 1.0          |                   |
| Substrate_stack_FTO_L1                           | 1.0          |                   |
| ETL_stack_sequence_TiO2-c_L0                     | 1.0          |                   |
| ETL_stack_sequence_TiO2-mp_L1                    | 1.0          |                   |
| HTL_stack_sequence_Spiro-MeOTAD_L0               | 1.0          |                   |
| Cell_area_measured                               | 0.04         | cm <sup>2</sup>   |
| Perovskite_deposition_number_of_steps            | 2.0          |                   |
| Perovskite_composition_perovskite_ABC3_structure | 1.0          |                   |
| JV_reverse_scan_Jsc                              | 16.42        | $\frac{mA}{cm^2}$ |

## 4 Optuna search space

### XGBoost

| <i>Parameter</i>     | <i>Value Range</i> |
|----------------------|--------------------|
| Max depth            | 1 - 15             |
| Learning rate        | 0.01 - 1.0         |
| Number of estimators | 50 - 1000          |
| Min child weight     | 1 - 10             |
| Gamma                | $10^{-8}$ - 1.0    |
| Subsample            | 0.5 - 0.9          |
| colsample_bytree     | 0.5 - 0.9          |

### Neural Network

| <i>Parameter</i>   | <i>Value Range</i>        |
|--------------------|---------------------------|
| Number of layers   | 1 - 3                     |
| Number of neurons  | input size - 2*input size |
| Dropout percentage | 0.2 - 0.5                 |
| Optimizer          | Adam, RMSprop, SGD        |
| Learning Rate      | $10^{-5}$ - 0.1           |

### Gaussian Process Regressor

| <i>Parameter</i>                | <i>Value Range</i>                   |
|---------------------------------|--------------------------------------|
| Number of kernels               | 1 - 3                                |
| Kernel types                    | Matern, Rational<br>Quadratic, White |
| Matern length scale             | $10^{-5}$ - $10^5$                   |
| Matern nu                       | 0.5 - 5                              |
| Rational quadratic length scale | $10^{-5}$ - $10^5$                   |
| Rational quadratic alpha        | $10^{-5}$ - $10^5$                   |
| White noise level               | $10^{-5}$ - $10^5$                   |
| Alpha                           | $10^{-3}$ - $10^3$                   |
| Number of restarts              | 0 - 10                               |

## 5 Final model configurations

### XGBoost

| <i>Parameter</i>     | <i>Value</i>          |
|----------------------|-----------------------|
| Max depth            | 8                     |
| Learning rate        | 0.325                 |
| Number of estimators | 525                   |
| Min child weight     | 6                     |
| Gamma                | $6.15 \times 10^{-5}$ |
| Subsample            | 0.89                  |
| colsample_bytree     | 0.73                  |

### Neural Network

| <i>Parameter</i>      | <i>Value</i>           |
|-----------------------|------------------------|
| Number of layers      | 1                      |
| Number of neurons L1  | 4147                   |
| Dropout percentage L1 | 0.361                  |
| Optimizer             | RMSprop                |
| Learning Rate         | $8.021 \times 10^{-5}$ |

### Gaussian Process Regressor

| <i>Parameter</i>         | <i>Value</i> |
|--------------------------|--------------|
| Number of kernels        | 1            |
| Kernel type L1           | Matern       |
| Length scale L1          | 87087.9      |
| Nu L1                    | 4.475        |
| White kernel noise level | 45167.8      |
| Alpha                    | 988          |
| Number of restarts       | 0            |