

## Supplemental Information

### Accelerated design of promising mixed **lead-free** double halide organic-inorganic perovskites for photovoltaics using machine learning

Yilei Wu<sup>†</sup>, Shuaihua Lu<sup>†</sup>, Minggang Ju<sup>\*</sup>, Qionghua Zhou<sup>\*</sup> & Jinlan Wang<sup>\*</sup>

School of Physics, Southeast University, Nanjing 211189, China.

These authors contributed equally: Yilei Wu, Shuaihua Lu.

Correspondence and requests for materials should be addressed to **M. J.**, Q.Z. and J.W.

(email: **juming@seu.edu.cn**; qh.zhou@seu.edu.cn; jlwang@seu.edu.cn)

## Table of Contents

### Supplementary Methods

Model selection	S3
Model evaluation for classification and regression algorithms	S3
Hyper-parameters selection	S6
Last-place elimination feature selection procedure	S6

### Supplementary Notes

The source of initial features	S6
The criteria for perovskite and non-perovskite	S7
Derivation of geometric limits	S7
Performance of regression model on the out-of-sample systems	S8
<b>The decomposition pathways of MA<sub>2</sub>GeSnI<sub>4</sub>Br<sub>2</sub> and MA<sub>2</sub>InB<sup>3+</sup>X<sub>4</sub>X'<sub>2</sub></b>	<b>S9</b>

### Supplementary Figures

Figure S1 Schematic diagram for replacing X-site ions of MA <sub>2</sub> AgBiI <sub>6</sub>	S10
Figure S2 Flowchart and results of last-place elimination feature selection procedure	S11

Figure S3 Representation of geometric limits of double perovskites	S12
Figure S4 Results of model selection	S13
Figure S5 Proportion of perovskites in candidates	S14
Figure S6 Selection of training & test sets for different ML models	S15
Figure S7 Flowchart for predicting candidates	S16
Figure S8 Number of predicted MDHOIPs with appropriate bandgap values	S17
Figure S9 Total energy during 5 ps AIMD simulations for $\text{FA}_2\text{AgBiBr}_4\text{Cl}_2$ , $\text{MA}_2\text{AgSbBr}_4\text{Cl}_2$ and $\text{MA}_2\text{AlInI}_2\text{Br}_4$	S18
Figure S10 Band structures and PDOS of $\text{FA}_2\text{AgBiBr}_4\text{Cl}_2$ , $\text{MA}_2\text{AgBiI}_2\text{Cl}_4$ , $\text{MA}_2\text{AgBiBr}_2\text{I}_4$ , $\text{MA}_2\text{AgSbBr}_4\text{Cl}_2$ , $\text{MA}_2\text{AlInI}_2\text{Br}_4$ , $\text{MA}_2\text{AuInBr}_4\text{Cl}_2$	S19
Figure S11 Band decomposed charge density of $\text{MA}_2\text{GeSnI}_4\text{Br}_2$ , $\text{MA}_2\text{InBiI}_2\text{Br}_4$ , $\text{FA}_2\text{InSbBr}_2\text{Cl}_4$ , $\text{MA}_2\text{AgInBr}_4\text{Cl}_2$	S20
Figure S12. DFT-calculated decomposition energies of $\text{MA}_2\text{GeSnI}_4\text{Br}_2$ , $\text{MA}_2\text{InBiI}_2\text{Br}_4$ , and $\text{FA}_2\text{InSbBr}_2\text{Cl}_4$	S21
<b>Supplementary Tables</b>	
Table S1 Different elements with common valence states	S22
Table S2 Eighty-seven initial features with description	S23
Table S3 Comparison between DFT-calculated and ML-predicted results of out-of-sample systems	S24
Table S4 Comparison between bandgaps from the database and our DFT results	S25
<b>Supplementary References</b>	S26

## Supplementary Methods

**Model selection.** The accuracy of machine learning (ML) models relies on appropriate algorithm, thus selecting out the best model from a series of available ML models is necessary. For ML classification, six classifiers are considered: gradient boosting classifier (GBC), support vector machine (SVM), AdaBoost classifier, random forest classifier (RFC), stochastic gradient descent classifier (SGDC), and decision trees classifier (DTC). For ML regression, six regressors are considered as well: gradient boosting regressor (GBR), kernel ridge regressor (KRR), bagging regressor, random forest regressor (RFR), kernel neighbors' regressor (KNR) and decision trees regressor (DTR).<sup>1</sup> The same training & test sets are applied to train each classification (regression) model. The same model evaluation indexes are utilized to evaluate fairly the performance of classification (regression) model. The GBC model shows the best performance among these six classification models, and the most appropriate regression model is GBR model.

**Model evaluation for classification and regression algorithms.** Model evaluation indexes are essential for measuring the performance of ML models. ML classification and regression models correspond to different model evaluation indexes. Four evaluation indexes (area under curve (AUC), accuracy, precision and recall) are applied for classification models and three different evaluation indexes (coefficient of determination ( $R^2$ ), mean square error (MSE) and mean absolute error (MAE)) are used for regression models.

The classification models produce the prediction probability for samples. The classification threshold is set to 0.5 in this work, and the prediction probability results of samples are compared to the pre-defined threshold. The prediction probability results correspond to the probabilities that samples belong to positive class (i.e. perovskite and perovskite with bandgaps smaller than 0.2 eV) or negative class (i.e. non-perovskite and perovskite with bandgaps larger than 0.2 eV). Therefore, samples with prediction probability results larger than 0.5 are classified into positive class, and

samples with prediction probability results less than 0.5 are classified into negative class.

According to the classification results, the count of positive samples predicted correctly is defined as true positive (TP), and the count of positive samples predicted falsely is defined as false positive (FP). The count of negative samples predicted correctly is defined as true negative (TN), and the count of negative samples predicted falsely is defined as false negative (FN). By calculating the values of TP, FP, TN and FN, we obtained the confusion matrix, which represents the counts of the predicted classes versus the true classes of test set. The confusion matrix is shown as follows:

	Predicted positive	Predicted negative
True positive	TP	FN
True negative	FP	TN

The true positive rate (TPR) and false positive rate (FPR) are calculated based on TP, TN, FP and FN.

$$TPR = \frac{TP}{TP+FN} \quad (1)$$

$$FPR = \frac{FP}{TN+FP} \quad (2)$$

The receiver operating characteristic (ROC) curve can be drawn using TPR and FPR as coordinates, and is often used to measure the performance of classification models. For ROC curves in this work, the probability of positive prediction is the number of positive prediction times derived by total number of times after 100 executions. When comparing the performance of different classification models, if the ROC curve of one model is completely below the ROC curve of the other, it means that the performance of the latter model is better than the former. If the ROC curves of two classification models intersect, the comparison is difficult. Thus the more appropriate model evaluation index is the AUC value. The higher AUC value correspond to the better performance of classification model. The AUC value of a classification model without learning algorithm is equal to 0.5, and the AUC value of the perfect classification model is equal to 1.

Accuracy is the proportion of correctly classified samples among all samples. In

general, the better classification model has the higher accuracy. Accuracy is defined as follows:

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (3)$$

Precision is the proportion of true positive samples in predicted positive class.

Precision is defined as follows:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (4)$$

Recall represents the ability for identifying positive samples of classification models. Recall is defined as follows:

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (5)$$

The above indexes are used to evaluate the performance of classification models, and three different indexes are chosen for evaluating the performance of regression models.

$R^2$  is the proportion of the variance in the dependent variable that is predictable from the independent variable.  $R^2$  is defined as follows:

$$R^2 = 1 - \frac{\sum_i (y_i^{\text{true}} - y_i^{\text{pred}})^2}{\sum_i (y_i^{\text{true}} - \bar{y}_i^{\text{true}})^2} \quad (6)$$

Where  $y_i^{\text{true}}$  are the true values, and  $y_i^{\text{pred}}$  are the predicted values. The predicted values of the perfect regression model are equal to true values, thus the value of  $R^2$  is equal to 1.

Mean square error (MSE) represents the average squared difference between the predicted values and true values. MSE is defined as follows:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i^{\text{true}} - y_i^{\text{pred}})^2 \quad (7)$$

Mean absolute error (MAE) represents the arithmetic average of the absolute errors between predicted values and true values. MAE is defined as follows:

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |y_i^{\text{true}} - y_i^{\text{pred}}| \quad (8)$$

**Hyper-parameters selection.** After model selection, hyper-parameters are optimized by applying a global search algorithm based on the simulated annealing algorithm.<sup>2</sup> Hyper-parameters are optimized before training ML models, and optimized hyper-parameters can improve the performance of ML models. We selected 7 initial hyper-parameters for GBC and GBR models, respectively. For GBC model, max\_depth is in the range of 2~30, n\_estimators is in the range of 40~200, learning\_rate is in the range of 0.005~0.5, subsample is in the range of 0.5~1, min\_child\_weight is in the range of 0.5~10, the random\_state\_seed is set to 42, and the loss function is set to logistic regression. For GBR model, max\_depth is in the range of 1~200, n\_estimators is in the range of 100~500, learning\_rate is in the range of 0.001~0.5, min\_sample\_leaf is in the range of 1~50, max\_features is in the range of 0.01~1, the random\_state\_seed is set to 42, the range of loss function is set to [least squares (ls), least absolute deviation (lad), the combination of ls and lad (huber), quantile regression (quantile)]. During the optimization process, ML models with the best hyper-parameters can achieve the maximum accuracy.

**Last-place elimination feature selection procedure.** To search the most relevant features, a “last-place elimination” feature selection procedure<sup>2, 3</sup> is introduced to GBC algorithm and GBR algorithm. In this work, the initial feature set is consisted of 87 features relating to polarizability, ionic radii and electronegativity. In the first step, 87 features are ranked according to the relative importance after training model. Then the feature at the last position is removed, and the remaining 86 features consist a new feature set. Next, the performance of model with new feature set is evaluated, and the above step is repeated until only two features left. Finally, the accuracy of the model at each step is analyzed, and the feature set corresponds to the inflection point is selected as the optimized feature set.

## Supplementary Notes

**Note S1.** The ionic polarizability of A-site ions in Table S1 is obtained from

Amsterdam Density Functional program package (ADF2013).<sup>3,4</sup> The ionic radii of A-, B-, B'- and X-site ions and electronegativity of B-, B'- and X-site ions are obtained from the python Mendeleev package 0.5.1.<sup>5</sup>

**Note S2.** Among all DHOIPs data with the chemical formula  $A_2BB'X_6$  in literature, we only selected compositions that satisfy the charge neutrality condition and Pauling's valence rule, resulting in the number of DHOIPs in training & test sets reduced to 2274. Then all of them are labeled 'perovskite' or 'non-perovskite'. The criterion is based on (i)  $\theta_{X-B-X} > 160^\circ$ ; (ii)  $R_{B-X}^{\min}/R_{B-X}^{\max} > 2/3$ . Where  $\theta_{X-B-X}$  represents the angle of the X-B-X bonds,  $R_{B-X}^{\min}$  and  $R_{B-X}^{\max}$  represent the minimum length and the maximum length of B-X bonds, respectively.

**Note S3.** To obtain geometric boundaries for double perovskite structure, the geometric limits using rigid sphere model are derived. The generalized Goldschmidt's parameters are introduced as follows: (i) the average octahedron factor:  $\bar{\mu} = (IR_B + IR_{B'})/2IR_X$ , (ii) the octahedron mismatch:  $\Delta\mu = (IR_B - IR_{B'})/2IR_X$ , (iii) the generalized tolerance factor:  $t = \frac{IR_A + IR_X}{\sqrt{2}\{[(IR_B + IR_{B'})/2 + IR_X]^2 + (IR_B - IR_{B'})^2/4\}^{1/2}}$ .<sup>6</sup> Where  $IR_A$ ,  $IR_B$ ,  $IR_{B'}$  and  $IR_X$  represent the ionic radii of A-, B-, B'- and X-site ions, respectively.

When  $IR_B$  is equal to  $IR_{B'}$ , X-site ions sit at the midway between B- and B'-site ions (Figure S3a). When  $IR_B$  is different from  $IR_{B'}$ , X-site ions shift the distance of  $1/2 |IR_B - IR_{B'}|$  from midway O toward large ions between B- and B'-site ions. This offset helps to relieve the lattice strain caused by size mismatch and to reach the overall electrostatic energy of double perovskite (Figure S3b).<sup>7</sup> The larger difference between  $IR_B$  and  $IR_{B'}$ , the larger the distance of offset. The octahedron limit corresponds to the extremal situation wherein two adjacent X-site ions in the same octahedron are tangent to each other (Figure S3c). In this situation the distance between centers of B'-site ions and X-site ions satisfies the condition  $IR_{B'} = (\sqrt{2} - 1)IR_X$ . Then both sides of equation are divided by  $R_X$  and  $\bar{\mu} - \Delta\mu = (\sqrt{2} - 1)$

is obtained. Therefore, the ionic radii must satisfy the condition  $\bar{\mu} - \Delta\mu \geq (\sqrt{2} - 1)$ .

If  $\bar{\mu} - \Delta\mu < (\sqrt{2} - 1)$ , the B/B'-site ions cannot touch six adjacent X-site ions, bringing the instability from reduced coordination number of B/B'-site ions.

The stretch limit is also considered, which corresponds to the extremal situation wherein A-site ion is so large that it is tangent to all twelve X-site ions around the octahedron cavity (Figure S3d). The distance between A-site ions and midway O is  $\frac{\sqrt{2}}{2}(IR_{B'} + 2IR_X + IR_B)$ , and the distance between X-site ions and midway O is  $\frac{1}{2}(IR_B - IR_{B'})$ . According to Pitagoras' theorem to the triangle, the boundary condition  $(IR_A + IR_X)^2 = \frac{1}{4}(IR_B - IR_{B'})^2 + \frac{1}{2}(IR_{B'} + 2IR_X + IR_B)^2$  is obtained. After combining this condition to the generalized tolerance factor, the geometric boundary  $t = 1$  is obtained. For  $t > 1$ , A-site ions are too large to maintain three-dimensional perovskite structure, thus perovskites are likely to form low dimensional structures.

**Note S4.** To validate the generalization ability of the bandgap regression model, we randomly divided the training & test sets (525 DHOIPs) into two subsets based on combinations of B and B'-site cations. In detail, each DHOIPs corresponding to the same combination appears in the same subset. Then one subset (498 DHOIPs) is used to train the regression model, while another subset (27 DHOIPs) contains 6 combinations, and MDHOIPs corresponding to these combinations are labeled as out-of-sample systems, e.g. Ag&In, In&Sb, Ag&Sb, Au&In, Sn&Pb, and As&In. Subsequently, two models are utilized to predict the bandgap values of these out-of-sample systems: ML model trained by 525 DHOIPs (model-1), and ML model trained by 498 DHOIPs (model-2). The comparison between ML-predicted and DFT-calculated results is listed in Table R1. For model-2, the maximum error between ML-predicted and DFT-calculated bandgap values is 0.228 eV, and most of MDHOIPs have errors within 0.15 eV, which is slightly higher than that of ML model-1. Overall, our ML model shows well generalization ability on these

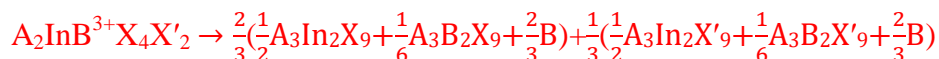
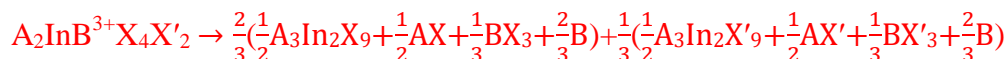
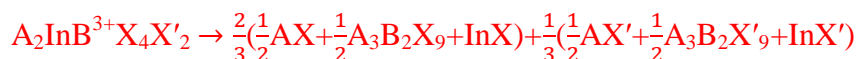
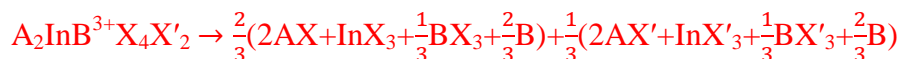
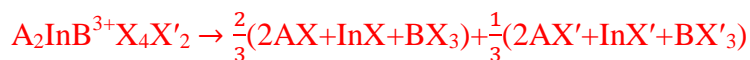


out-of-sample systems.

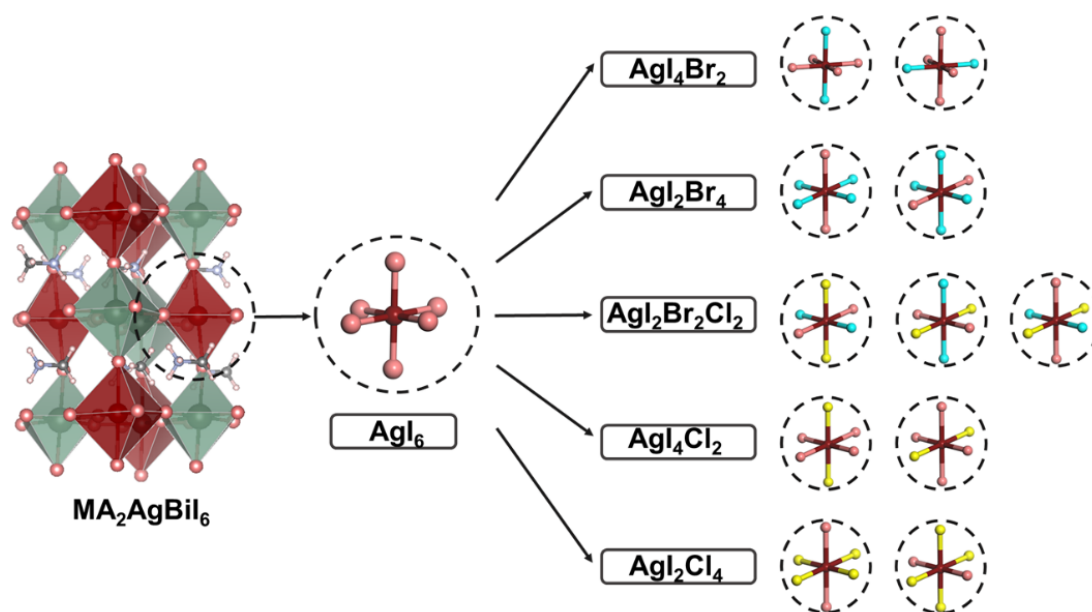
**Note S5.** In this work, the decomposition energy ( $\Delta H$ ) of  $\text{MA}_2\text{GeSnI}_4\text{Br}_2$  is calculated through the pathways as follows:



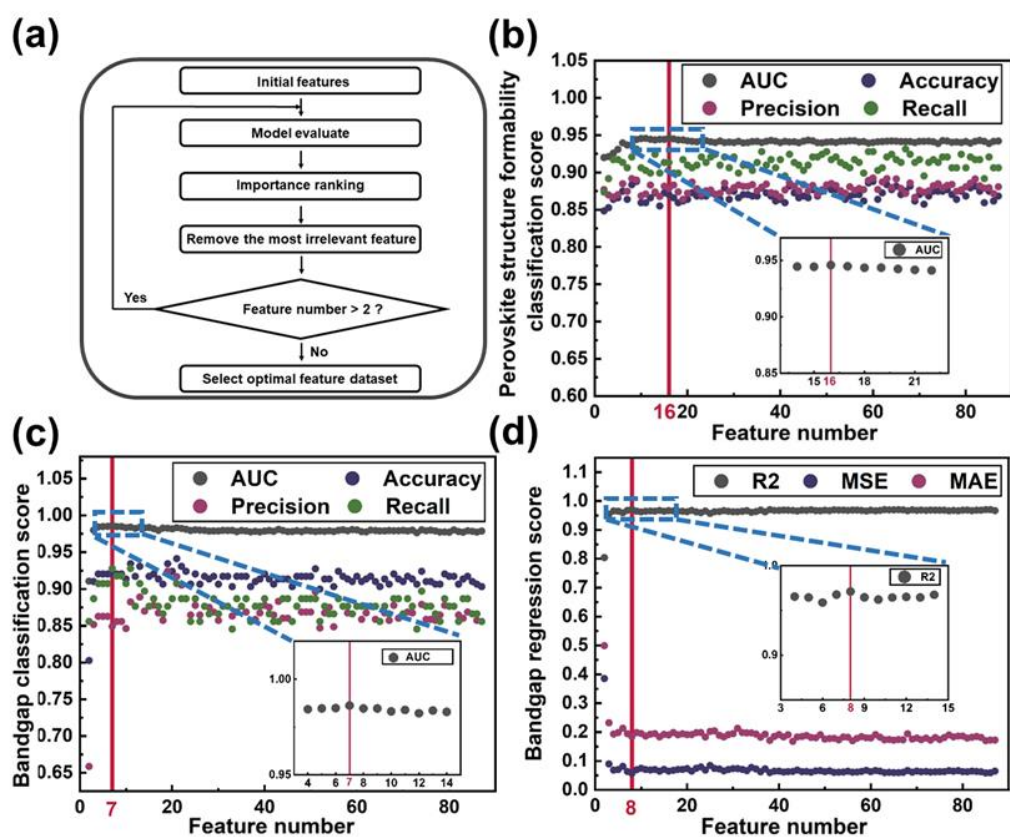
The  $\Delta H$  of  $\text{A}_2\text{InB}^{3+}\text{X}_4\text{X}'_2$  (i.e.,  $\text{MA}_2\text{InBiI}_2\text{Br}_4$  and  $\text{FA}_2\text{InSbBr}_2\text{Cl}_4$ ) is calculated through the pathways as follows:



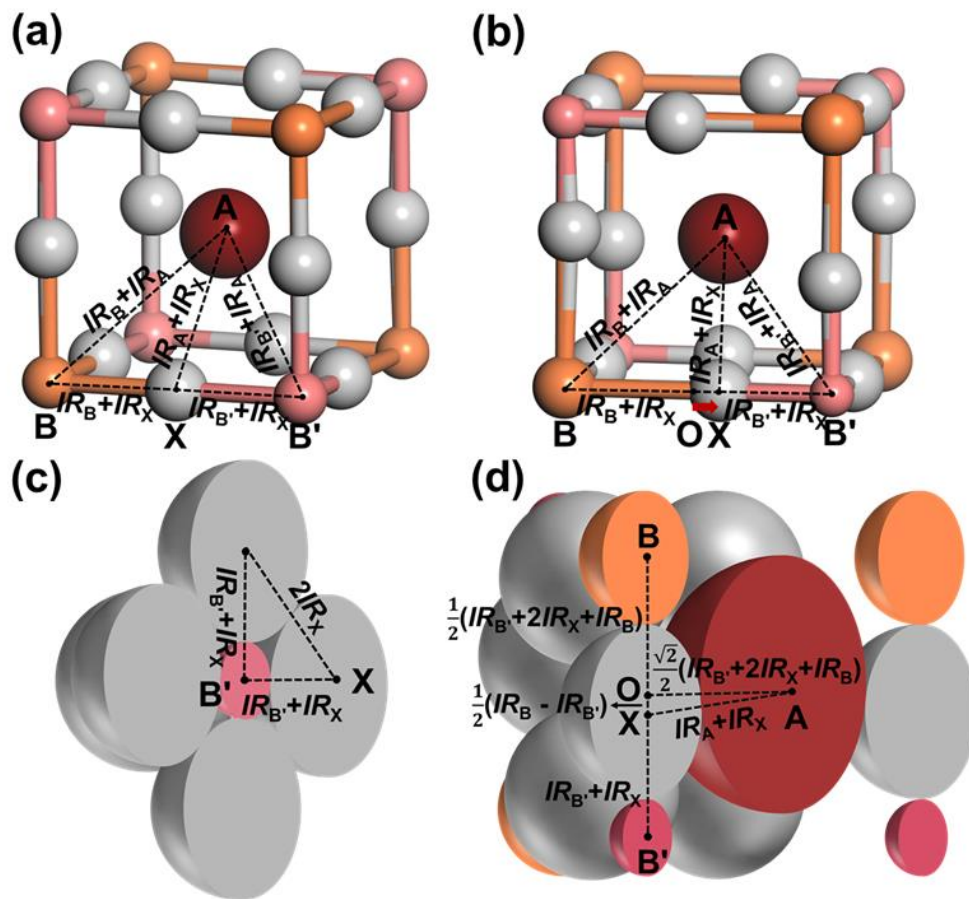
## Supplementary Figures



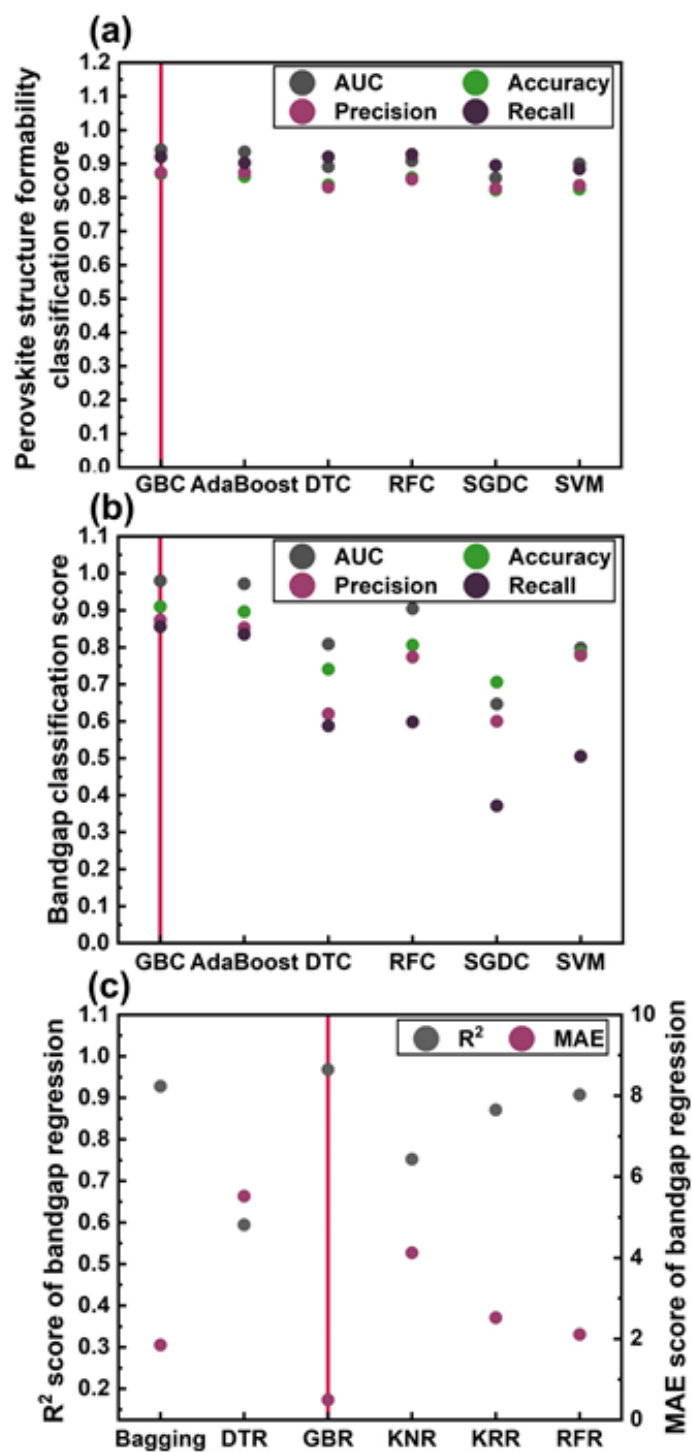
**Figure S1.** 11 candidates in the prediction set are obtained based on  $\text{MA}_2\text{AgBiI}_6$  in training&test sets.



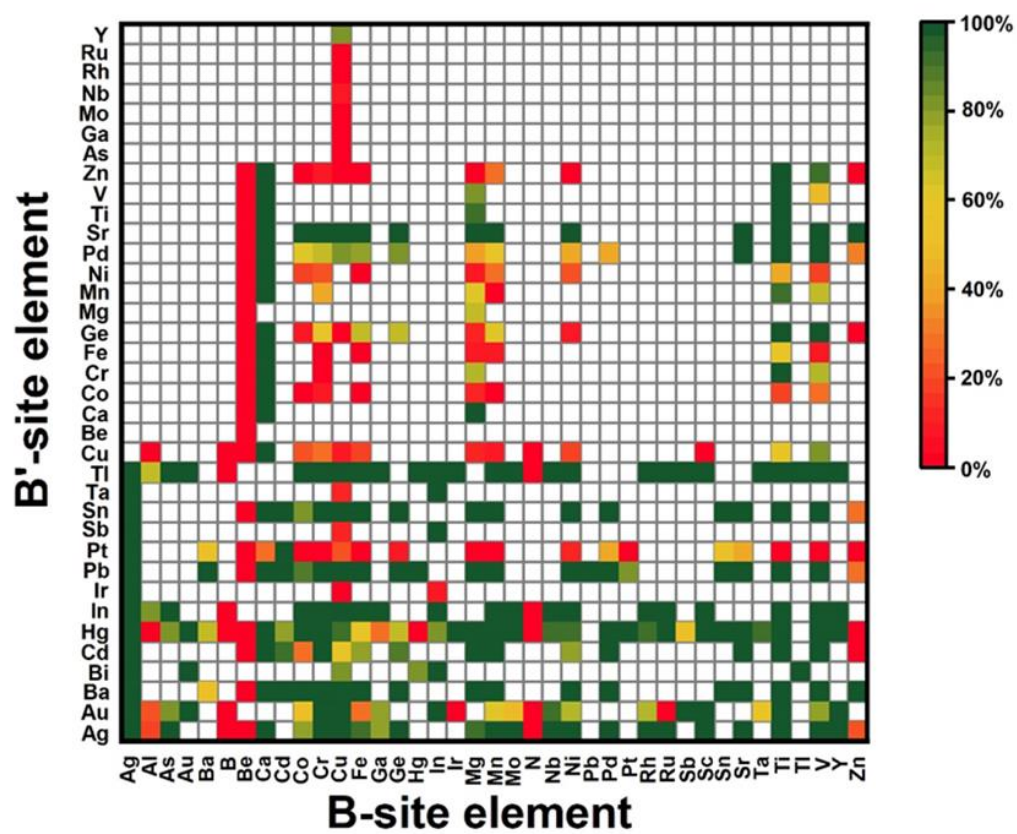
**Figure S2.** Flowchart and results of last-place elimination feature selection procedure. (a) Feature engineering framework combined with "last-place elimination" method. Optimized feature set of (b) perovskite structure formability classification, (c) bandgap classification, and (d) bandgap regression.



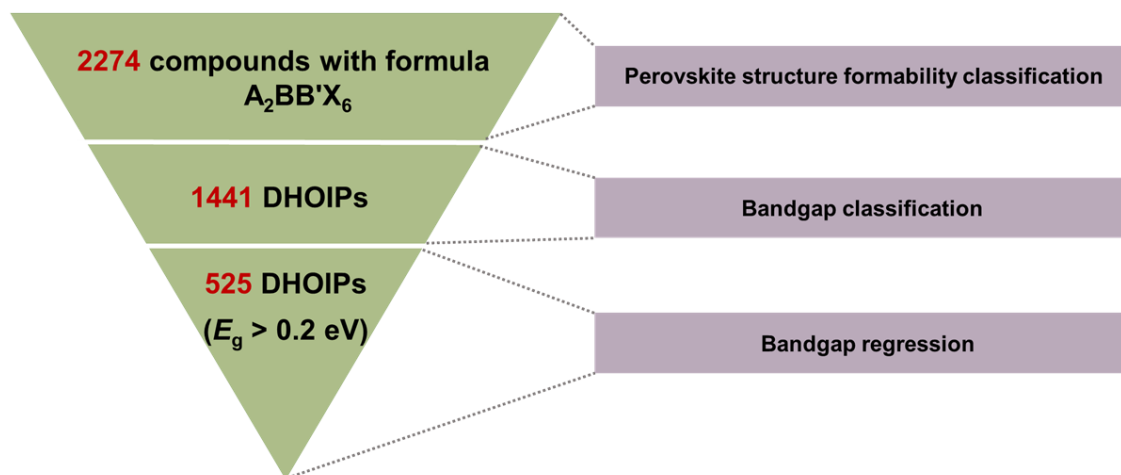
**Figure S3.** Representation for conventional cell of double perovskites. Red, orange, pink and grey spheres represent A-site ions, spheres represent B-site ions, B'-site ions and X-site ions, respectively. Schematic diagram corresponds to the situation wherein B- and B'-site ions possess (a) same ionic radii and (b) different ionic radii. Schematic representations for (c) the octahedron limit and (d) the stretch limit of double perovskites.



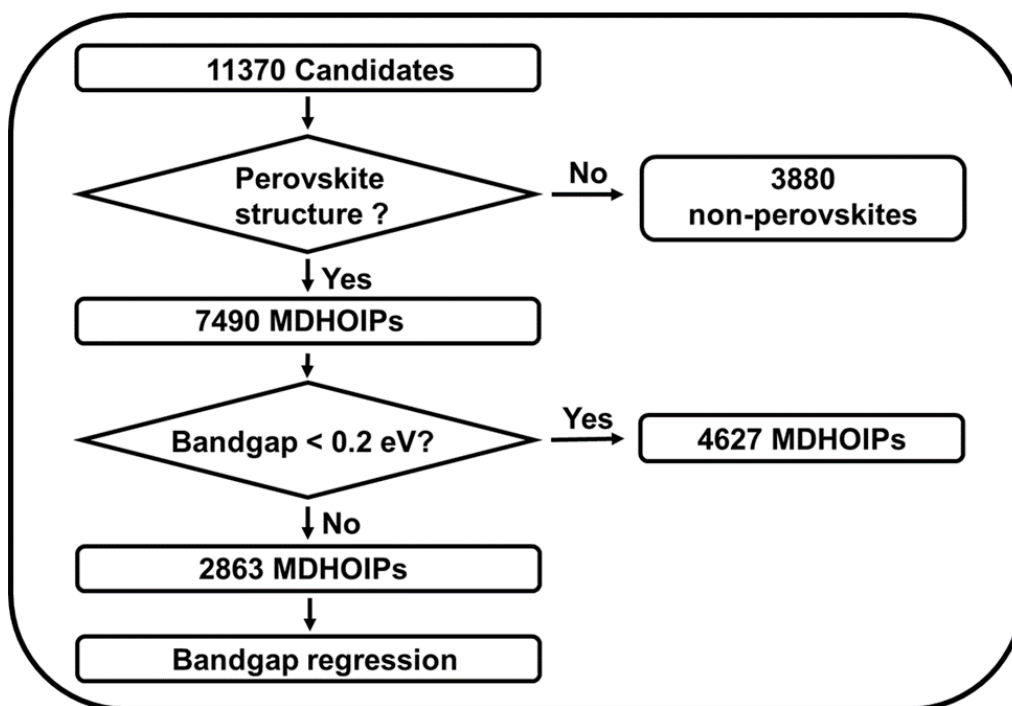
**Figure S4.** Classification model score of (a) perovskite structure formability classification and (b) bandgap classification. (c) Regression model score of bandgap regression.



**Figure S5.** Proportion of perovskites in candidates based on different combinations of B- and B'-site ions.

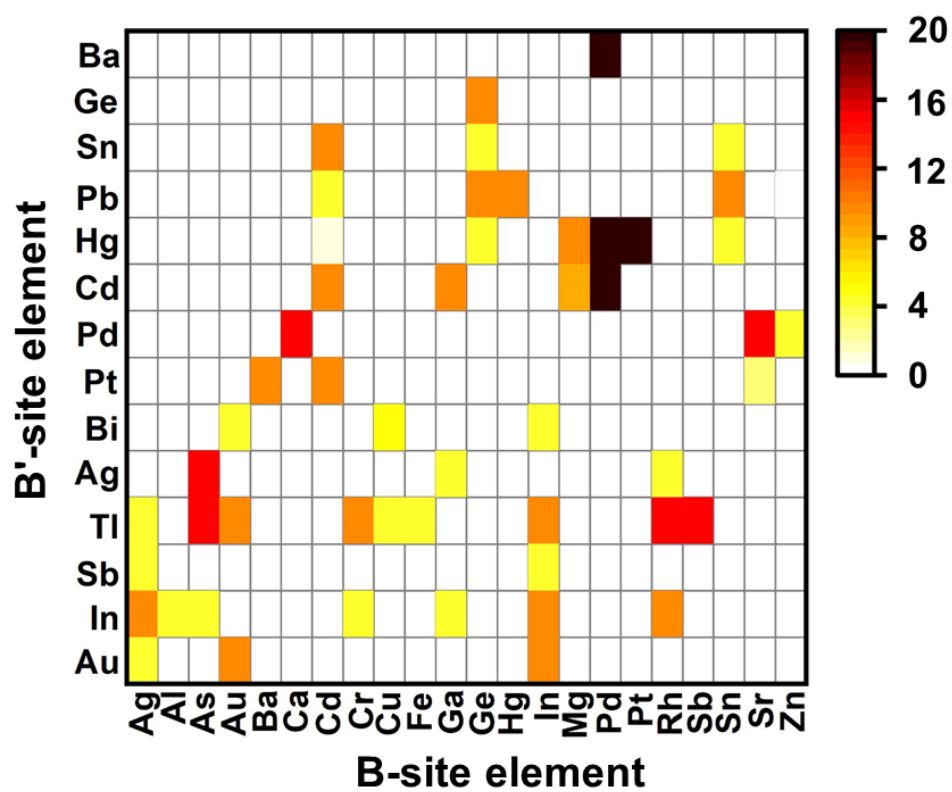


**Figure S6.** Selection of training & test sets corresponds to three ML models with different target properties. Preparing for model training and test, the training set and test set for each model are divided according to the proportion of 80% and 20%.

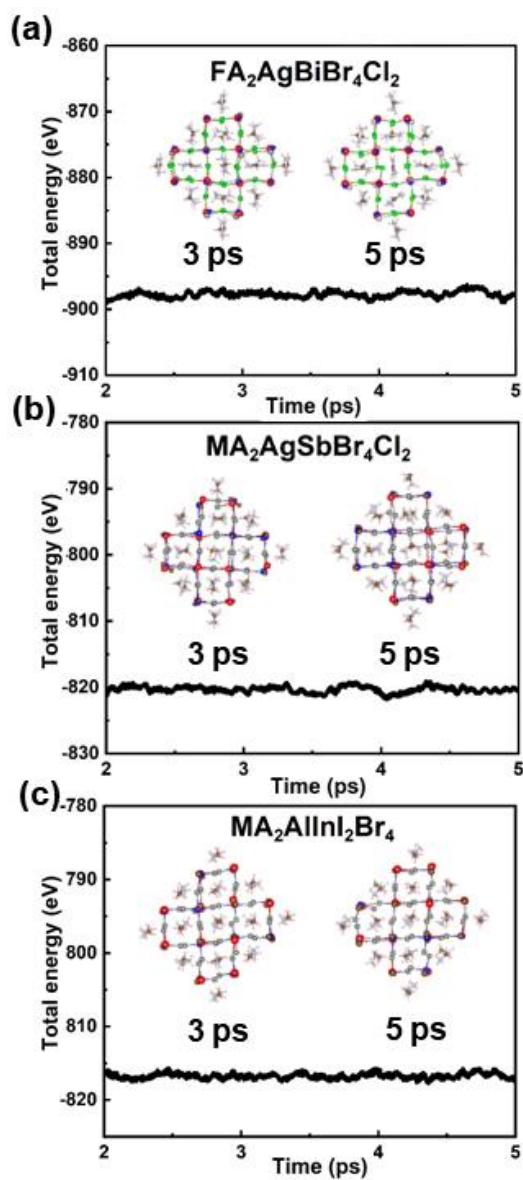


**Figure S7.** Flowchart for predicting candidates using trained ML models.

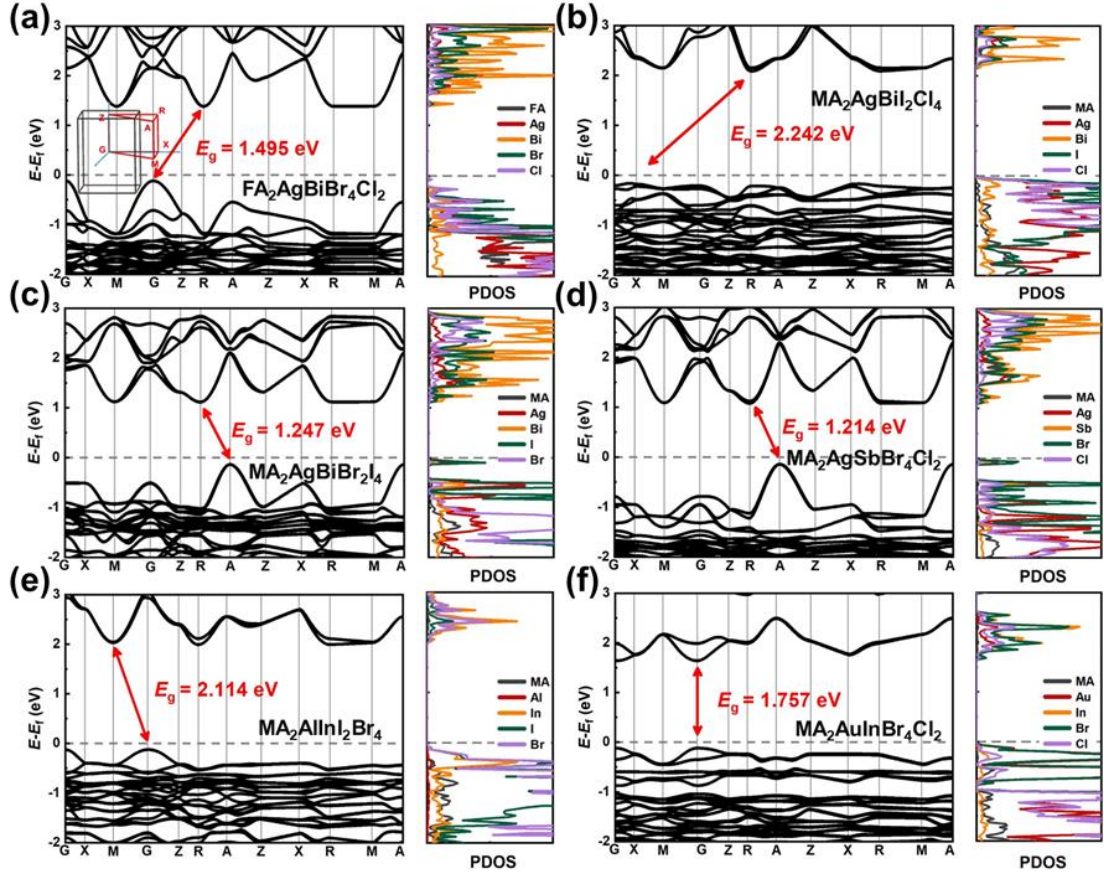




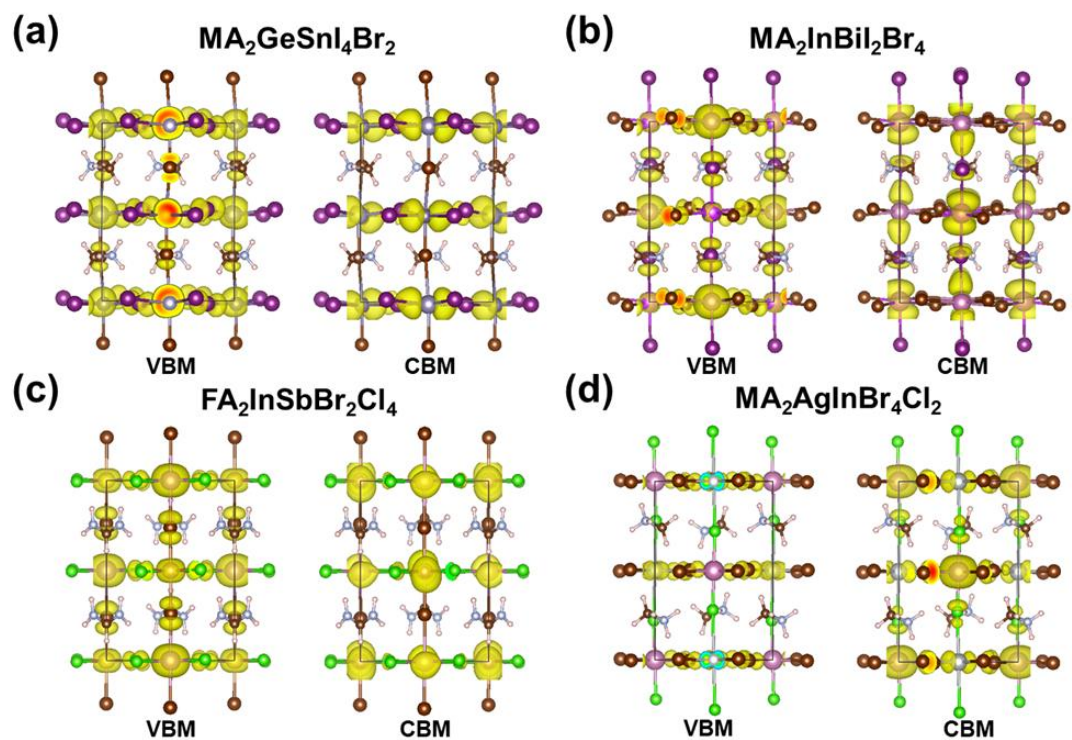
**Figure S8.** Number of optimal MDHOIPs based on different combinations of B- and B'-site ions.



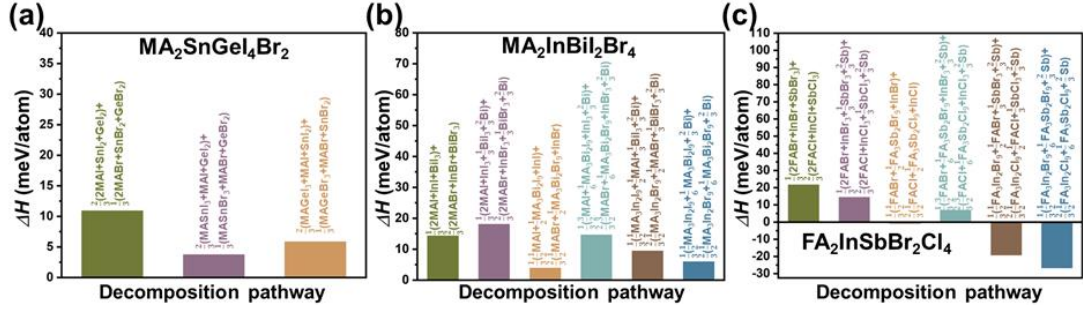
**Figure S9.** Total energy during 5 ps ab initio molecular dynamics (AIMD) simulations for (a)  $\text{FA}_2\text{AgBiBr}_4\text{Cl}_2$ , (b)  $\text{MA}_2\text{AgSbBr}_4\text{Cl}_2$ , and (c)  $\text{MA}_2\text{AlInI}_2\text{Br}_4$ .



**Figure S10.** Calculated band structures and PDOS of (a) FA<sub>2</sub>AgBiBr<sub>4</sub>Cl<sub>2</sub>, (b) MA<sub>2</sub>AgBiI<sub>2</sub>Cl<sub>4</sub>, (c) MA<sub>2</sub>AgBiBr<sub>2</sub>I<sub>4</sub>, (d) MA<sub>2</sub>AgSbBr<sub>4</sub>Cl<sub>2</sub>, (e) MA<sub>2</sub>AlInI<sub>2</sub>Br<sub>4</sub>, and (f) MA<sub>2</sub>AuInBr<sub>4</sub>Cl<sub>2</sub>. Subgraph represents the Brillouin zone for tetragonal lattice.



**Figure S11.** Band decomposed charge density of (a)  $\text{MA}_2\text{GeSnI}_4\text{Br}_2$ , (b)  $\text{MA}_2\text{InBiI}_2\text{Br}_4$ , (c)  $\text{FA}_2\text{InSbBr}_2\text{Cl}_4$ , and (d)  $\text{MA}_2\text{AgInBr}_4\text{Cl}_2$ .



**Figure S12.** DFT-calculated decomposition energies of (a)  $MA_2GeSn_4Br_2$ , (b)  $MA_2InBi_2Br_4$ , and (c)  $FA_2InSbBr_2Cl_4$ .

## Supplementary Tables

**Table S1.** Different elements with common valence states.

Valence	elements
+1	Ag, Au, Cu, Hg, In, Tl
+2	Ag, Ba, Be, Ca, Cd, Co, Cr, Cu, Fe, Ge, Hg, Mg, Mn, Ni, Pd, Pb, Pt, Sn, Sr, Ti, V, Zn
+3	Al, As, Au, B, Bi, Co, Cr, Fe, Ga, In, Ir, Mn, Mo, N, Nb, Ni, Rh, Ru, Sb, Sc, Ta, Ti, V, Y

**Table S2.** Eighty-seven initial features with description.

Feature	Description
$P_A$	Ionic polarizability of the A-site cations
$\chi_B, \chi_{B'}, \chi_{X1}, \chi_{X2}$ and $\chi_{X3}$	Electronegativity of the B-, B'-, X <sub>1</sub> , X <sub>2</sub> and X <sub>3</sub> -site ions
$IR_A, IR_B, IR_{B'}, IR_{X1}, IR_{X2}$ and $IR_{X3}$	Ionic radii of the A-, B-, B'-, X <sub>1</sub> , X <sub>2</sub> and X <sub>3</sub> -site ions
$IR_{i+j}$	Sum of two ionic radii
$IR_{i-j}$	Difference between two ionic radii
$IR_{i/j}$	Ratio between two ionic radii
$\chi_{i+j}$	Sum of electronegativity of two ions
$\chi_{i-j}$	Difference between electronegativity of two ions
$\chi_{i/j}$	Ratio between electronegativity of two ions

**Table S3.** Comparison between DFT-calculated and ML-predicted results.

System	$E_g^{\text{DFT}}$ (eV) <sup>a</sup>	$E_g^{\text{ML1}}$ (eV) <sup>a</sup>	$E_g^{\text{ML2}}$ (eV) <sup>a</sup>
FA <sub>2</sub> InSbBr <sub>2</sub> Cl <sub>4</sub>	0.985	1.013	1.011
MA <sub>2</sub> AgInBr <sub>4</sub> Cl <sub>2</sub>	0.949	0.860	1.177
MA <sub>2</sub> AgSbBr <sub>4</sub> Cl <sub>2</sub>	1.214	1.316	1.096
MA <sub>2</sub> AuInBr <sub>4</sub> Cl <sub>2</sub>	1.757	1.609	1.613
FA <sub>2</sub> SnPbBr <sub>2</sub> Cl <sub>4</sub>	1.507	1.442	1.401
MA <sub>2</sub> AsInCl <sub>4</sub> I <sub>2</sub>	0.712	0.786	0.826

<sup>a</sup>  $E_g^{\text{DFT}}$  represents DFT-calculated bandgap values, and  $E_g^{\text{ML1}}$  and  $E_g^{\text{ML2}}$  represent ML-predicted bandgap values obtained by model-1 and model-2, respectively. In which model-1 is trained by all 525 DHOIPs, and model-2 is trained by 498 DHOIPs.



**Table S4.** Comparison between bandgaps from the database and our DFT results.

System	Database (eV)	Our PBE (eV)
MAPbI <sub>3</sub>	1.74	1.74
MA <sub>2</sub> AgBiI <sub>6</sub>	1.34	1.23
Cs <sub>2</sub> AgBiCl <sub>6</sub>	1.80	1.82
MA <sub>2</sub> AgSbI <sub>6</sub>	1.04	0.92
Cs <sub>2</sub> AgInCl <sub>6</sub>	1.01	1.00

## Supplementary Reference

- 1 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, *J. Mach. Learn. Res.*, 2011, **12**, 2825–2830.
- 2 S.; Lu, Q.; Zhou, L.; Ma, Y.; Guo, and J. Wang, *Small Methods*, 2019, **3**, 1900360.
- 3 S. Lu, Q. Zhou, Y. Ouyang, Y. Guo, Q. Li, and J. Wang, *Nat. Commun.*, 2018, **9**, 3405.
- 4 G. Te Velde, F. M. Bickelhaupt, E. J. Baerends, C. Fonseca Guerra, S. J. A. van Gisbergen, J. G. Snijders, and T. Ziegler, *J. Comput. Chem.*, 2001, **22**, 931-967.
- 5 L. M. Mentel, <https://bitbucket.org/lukaszmentel/mendeleev>, 2014.
- 6 M. R. Filip and F. Giustino, *Proc. Natl. Acad. Sci. U. S. A.*, 2018, **115**, 5397-5402.
- 7 G. Pilania, A. Ghosh, S. T. Hartman, R. Mishra, C. R. Stanek, and B. P. Uberuaga, *npj Comput. Mater.*, 2020, **6**, 71.