

## Supporting information: "Data-driven analysis of the electronic-structure factors controlling the work functions of perovskite oxides"

Yihuang Xiong,<sup>a</sup> Weinan Chen,<sup>a</sup> Wenbo Guo,<sup>b</sup> Hua Wei,<sup>b</sup> and Ismaila Dabo<sup>a,c</sup>

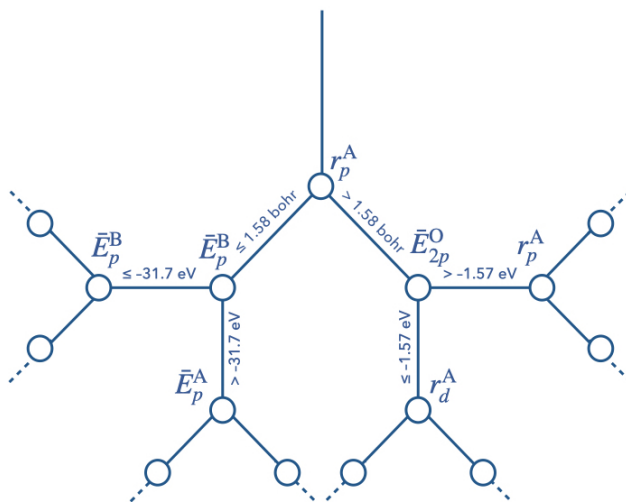
<sup>a</sup> Department of Materials Science and Engineering, and Materials Research Institute, The Pennsylvania State University, University Park, PA, USA

<sup>b</sup> College of Information Sciences and Technology, The Pennsylvania State University, University Park, PA 16802, USA

<sup>c</sup> Institutes of Energy and Environment, The Pennsylvania State University, University Park, PA 16802, USA

### 1 Regression tree

A regression tree is depicted schematically in Fig. S1.



**Fig. S1** Decision tree in the random forest model. Following the decision rule shown at each branching, the dataset is split into subsets.

### 2 Geometric mean of the electronegativity of the perovskite constituents ( $\chi_M^{ABO_3}$ ) compared to the work function

Various methods have been used to predict the absolute band edge energies, including an empirical method based on the geometric mean of the atomic Mulliken electronegativities ( $\chi_M$ )<sup>1</sup>. As previously pointed out<sup>2,3</sup>, this method cannot distinguish between materials that have the same chemical formula but different crystal structures. We plot the work functions with respect to  $\chi_M$  in Fig. S2. It can be seen that for either AO or BO<sub>2</sub> terminations, the correlation is scattered. In fact,  $\chi_M$  is shown to be much less relevant than other features through the recursive reduction process presented in Sec. III of the main text.

### 3 Feature elimination based on correlation matrix

In this work, we start by using a set of 38 features. As discussed in the main text, using highly correlated features would not deteriorate the performance of random forest regressor; however, features that carry similar information would dilute the importance score, which could lead to the misidentification of the predominant features. To overcome this issue, we first perform a Pearson correlation analysis to remove strongly correlated features. The Pearson correlation coefficient measures the linear relationship between variables. The resulting coefficients lie between -1 to 1, indicating a negative or positive correlation. The Pearson correlation coefficient  $r_{xy}$  is calculated as

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}, \quad (1)$$

where  $x_i$ ,  $y_i$  represent the individual point, and  $\bar{x}$ ,  $\bar{y}$  stand for the sample mean of the variables. The computed correlation matrix is shown in Fig. S3.

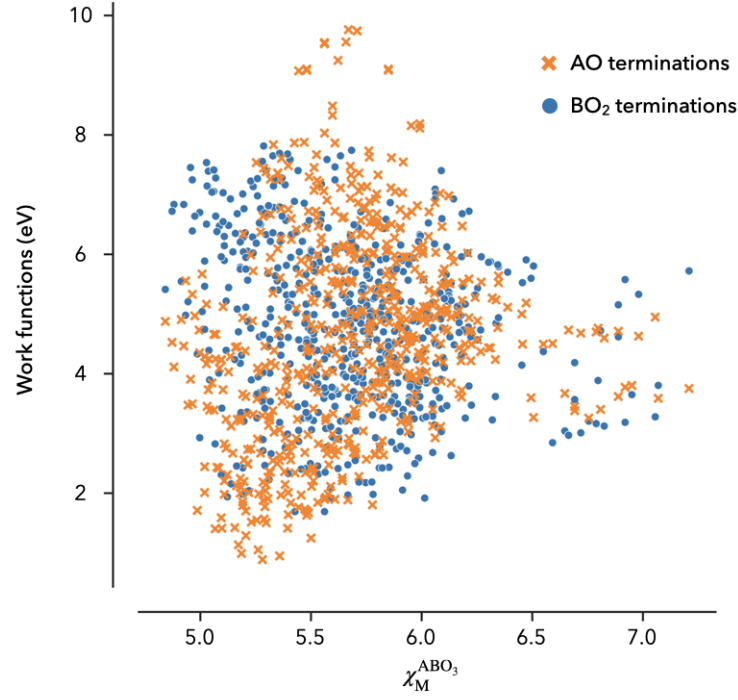


Fig. S2 Work function of AO and BO<sub>2</sub>-terminated perovskites as a function of the geometric mean of the atomic electronegativities  $\chi_M^{ABO_3}$ .

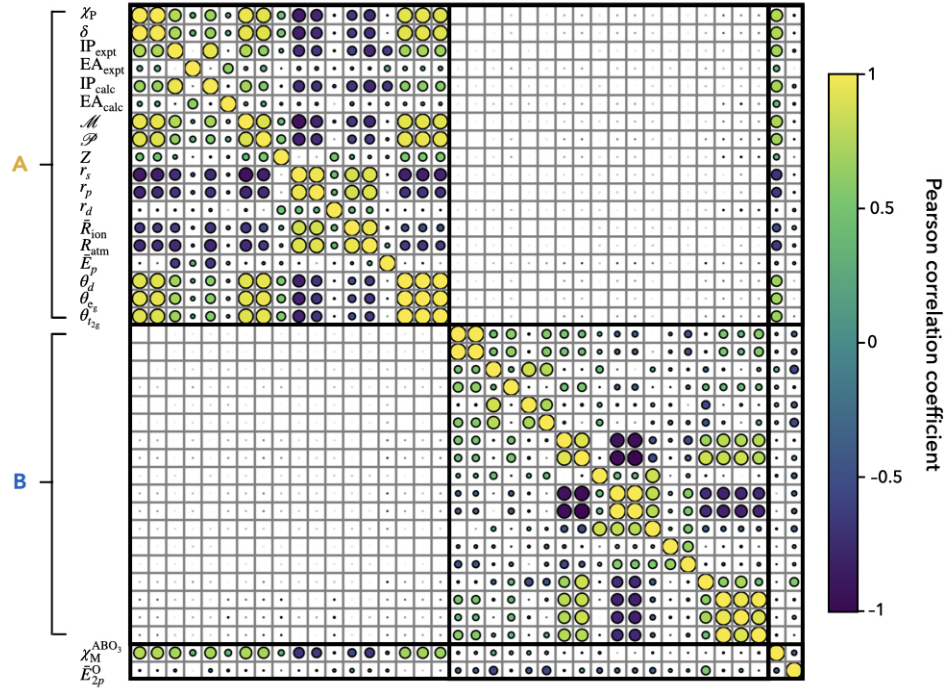


Fig. S3 Pearson matrix of the 38 features. The color indicates whether the correlation is positive or negative. The marker size represents the magnitude of the correlation. Both the A and B elements have 18 features. The geometric mean of the atomic electronegativities ( $\chi_M^{ABO_3}$ ) and the center of the oxygen 2p band ( $\bar{E}_{2p}^O$ ) are also included.

Based on the Pearson correlation coefficient, we removed highly-correlated features if its Pearson correlation coefficient is larger than 0.8. This process reduces the number of features from 38 to 21. The features for A are the experimental and calculated electron affinity  $EA_{expt}^A$ ,  $EA_{calc}^A$ , calculated ionization potential  $IP_{calc}^A$ , radius of s and d orbital  $r_s^A$  and  $r_d^A$ , atomic number  $Z^A$ , band center of p orbital  $\bar{E}_p^A$

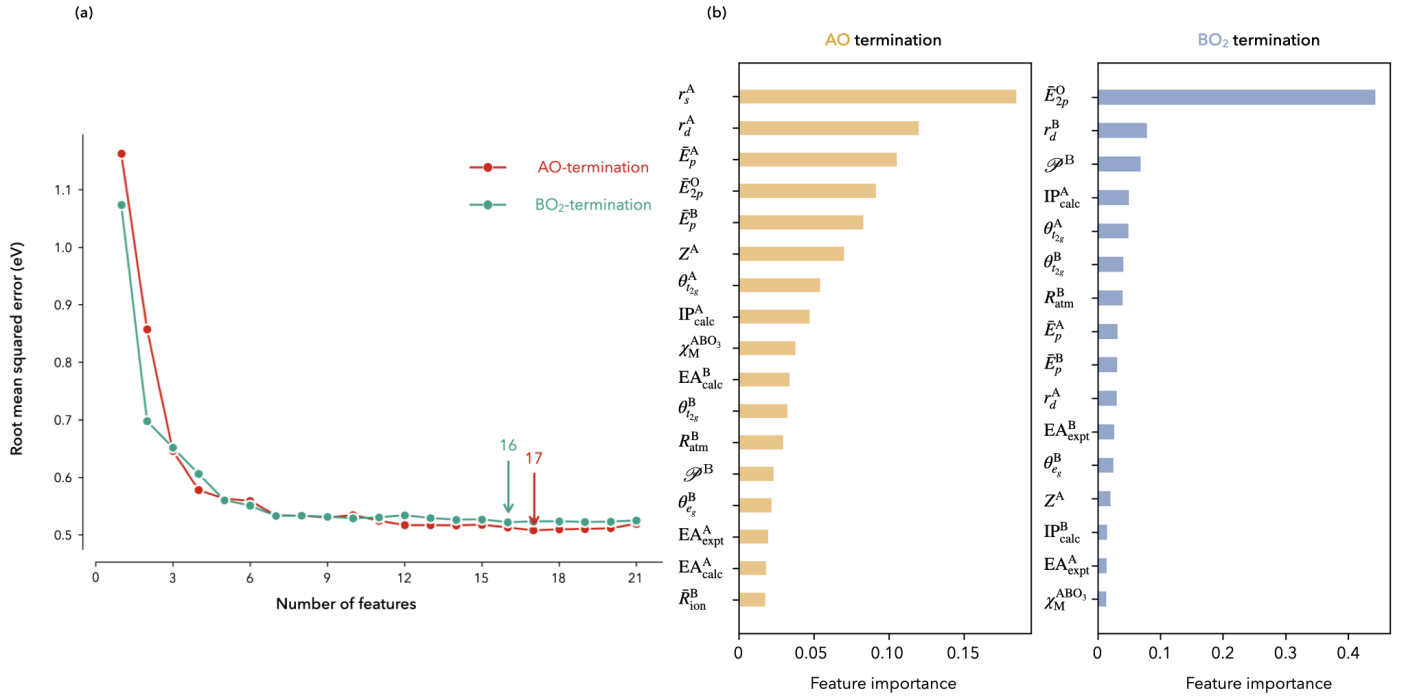
and filling factor  $\theta_{t_{2g}}^A$ . The features for B are the experimental and calculated electron affinity  $EA_{\text{expt}}^B$  and  $EA_{\text{calc}}^B$ , calculated ionization potential  $IP_{\text{calc}}^B$ , atomic radii and averaged ionic radii  $R_{\text{atm}}^B$ ,  $\bar{R}_{\text{ion}}^B$ , radius of  $d$  orbital  $r_d^B$ , Pauling's electronegativity  $\chi_p^B$ , band center of  $p$  orbital  $\bar{E}_p^B$ , filling factor of  $e_g$  band  $\theta_d^B$  and  $t_{2g}$  band  $\theta_{t_{2g}}^B$ . Pettifor's chemical scale  $\mathcal{P}^B$ . The last two are the geometric mean of the Mulliken electronegativity of the bulk perovskite  $\chi_M^{\text{ABO}_3}$  and the  $2p$  band center of oxygen  $\bar{E}_{2p}^O$ .

#### 4 Hyperparameters tuning and recursive feature elimination

The random forest models are developed using SCIKIT-LEARN library<sup>4</sup>. Four primary hyperparameters are considered, including the number of trees, the maximum depth of the tree, the minimum number of samples required to split an internal node, and the number of features to consider when looking for the best split. To select the hyperparameters that optimize the performance of the random forest regressor, a grid of hyperparameter combinations were generated using GridSearch or RandomizedSearch. In addition, we also perform the recursive feature elimination to optimize the feature space. To identify the most important features, we recursively remove the least important one, as we discussed in the main text. The change of root mean square error (RMSE) with respect to the number of features is shown in Fig. S4(a). Fig. S4(b) shows the normalized importance of all selected features for AO and BO<sub>2</sub> termination. The aforementioned hyperparameter tuning was performed when each feature is removed. This process yields the compatible number of features and hyperparameters. To comprehensively evaluating the model's performance, fivefold cross-validation was applied at each step during the recursive feature eliminations. For AO and BO<sub>2</sub> interfaces, the final selected hyperparameters are shown below:

AO-termination: n\_estimators: 130, max\_depth = 14, max\_features = "sqrt", min\_samples\_split = 2

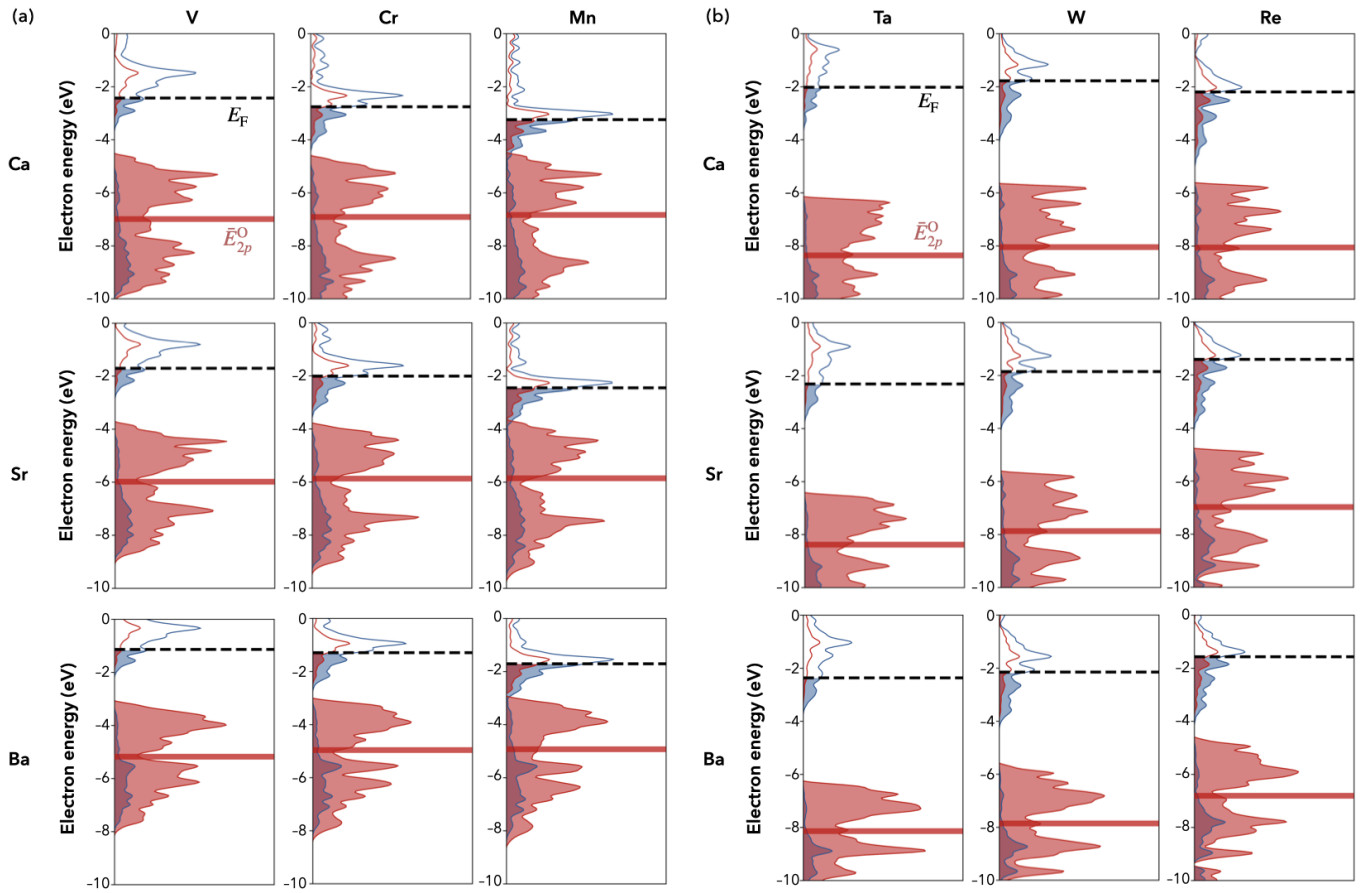
BO<sub>2</sub>-termination: n\_estimators: 100, max\_depth = 16, max\_features = "log2", min\_samples\_split = 2



**Fig. S4** (a) Averaged RMSE with respect to the number of features for the AO and BO<sub>2</sub> terminations. The results show that the best performances are achieved when 17 and 16 number of features are included in the random forest, respectively. (b) The overall importance ranking of AO and BO<sub>2</sub> terminations.

#### 5 Projected density of states of selected perovskites

We examined the partial dependence of the work functions of perovskites with respect to the bulk oxygen  $2p$  band center ( $\bar{E}_{2p}^O$ ). We found that the correlation between  $\bar{E}_{2p}^O$  and the work function is almost linear, except for the ones with deep  $\bar{E}_{2p}^O$  ( $< -4$  eV). We pay specific attention to the compounds that are close to the plateau of the PDP of  $\bar{E}_{2p}^O$  and the work functions. Those perovskites turn out to have relatively low work functions, but the decrease of the work function is limited, especially for AO terminations, as shown in Fig. 7(a) in the main text. To understand this trend, we plot the projected density of states of selected perovskites with AO terminations in Fig. S5. We found that for the perovskites with  $3d$  transition metals, the  $\bar{E}_{2p}^O$  energy level is almost constant with respect to vacuum



**Fig. S5** Projected density of states for perovskite slabs with AO terminations. The representative compositions contains alkaline earth metals (Ca, Sr, Ba) for A site, and (a) 3d (V, Cr, Mn) and 5d (Ta, W, Re) transition metals for B sites. The vacuum energy is referenced to 0, the Fermi energy ( $E_F$ ) and oxygen 2p band center  $\bar{E}_{2p}^O$  are shown with dashed line and red solid line.

(the vacuum energy is referenced to 0 eV). On the other hand, for Ta, W, Re, the hybridization between the  $d$  bands and the oxygen 2p bands involves a noticeable shift of the oxygen 2p bands with respect to vacuum. This shift largely cancels out the decrease in  $\bar{E}_{2p}^O$  and ultimately alters the correlation between  $\bar{E}_{2p}^O$  and the work function.

## Notes and references

- [1] I. E. Castelli, D. D. Landis, K. S. Thygesen, S. Dahl, I. Chorkendorff, T. F. Jaramillo, and K. W. Jacobsen, *Energy Environ. Sci.* **5**, 9034 (2012).
- [2] V. Stevanović, S. Lany, D. S. Ginley, W. Tumas, and A. Zunger, *Phys. Chem. Chem. Phys.* **16**, 3706 (2014).
- [3] Q. Yan, J. Yu, S. K. Suram, L. Zhou, A. Shinde, P. F. Newhouse, W. Chen, G. Li, K. A. Persson, J. M. Gregoire, and J. B. Neaton, *PNAS* **114**, 3040 (2017).
- [4] 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.* **12**, 2825 (2011).