

# A High-Throughput Computational Dataset of Halide Perovskite Alloys<sup>†</sup>

Jiaqi Yang<sup>a</sup>, Panayotis Manganaris<sup>a</sup>, and Arun Mannodi-Kanakkithodi<sup>a</sup>

## Abstract

Novel halide Perovskites with improved stability and optoelectronic properties can be designed via composition engineering at cation and/or anion sites. Data-driven methods, especially high-throughput first principles computations and subsequent analysis based on unique materials descriptors, are key to achieving this goal. In this work, we report a Density Functional Theory (DFT) based dataset of 495 ABX<sub>3</sub> halide Perovskite compounds, with various atomic and molecular species considered at A, B and X sites, and different amounts of mixing considered at each site generated using the Special Quasirandom Structures (SQS) algorithm for alloys. We perform GGA-PBE calculations on pseudo-cubic Perovskite structures to determine their lattice constants, stability in terms of formation and decomposition energies, electronic band gaps, and properties extracted from optical absorption spectra. To elucidate the importance of the level of theory used, we further perform 299 calculations using the more expensive HSE06 functional and determine lattice constant, stability and band gap, and compare PBE and HSE06 properties with some experimentally measured results. Trends in the datasets are unraveled in terms of the effects of mixing at different sites, the composition in terms of specific atomic or molecular species, and averaged elemental properties of species at different sites. This work presents the most comprehensive DFT perovskite alloy dataset to date and the data, which is open-source, can be exploited to train predictive and optimization models for accelerating the design of completely new compositions that may yield large solar cell efficiencies and improved performance across many optoelectronic applications.

## TODO Introduction

- State "DONE" from [2022-06-05 Sun 18:16]

The perovskite structure has been widely investigated by material scientists in recent decades for its promising industry applications. A perovskite unit cell contains a generalized three-component ABX<sub>3</sub> formula unit wherein A and B are cations with different oxidation states while X is an anion. The symbolic 3D perovskite structure is a network of BX<sub>6</sub> octahedra which are held in place by large A-site cations. This unique structure makes the crystal's electronic properties incredibly tunable. Numerous research efforts are devoted to halide perovskite (HaP), specifically, as a promising solar cell absorption material[XXX]. ABX<sub>3</sub> halide perovskite X-sites anions consist of halogens such as I and Br, B-site cations may be divalent elements such as Pb and Sn, and the A-site is occupied by large monovalent cations. The A-site cation can be so large that either inorganic (e.g Cs, K, Rb) elements or organic ligands (e.g Methylammonium (MA) and Formamidinium (FA)). The most commonly studied hybrid organic-inorganic halide perovskite, MAPbI<sub>3</sub> and FAPbI<sub>3</sub>, have relatively recently demonstrated extraordinary power conversion efficiency (PCE) between 20\single or multi junction

solar cells[XXX]. This is a five fold improvement over prior art and shows off the most attractive feature of halide perovskites, their unique tunability. Experimentally, it has been shown that a perovskite structure is considered stable if A-site, B-site, and X-site ionic radii satisfy the well-known tolerance (t) and octahedral ( $\mu$ ) relations[XXX]. Even under these restrictions, the configuration space of HaP structures, alloy ratios, alloy ordering, and defects still poses a highly multidimensional optimization problem. Generally speaking, halide perovskite design can be divided into several aspects. (1) Compounds: (2) Structure and phases: (3) Defects (4) Polymorphism.

We report a synthetic dataset collected for 495 chemically distinct, pseudo cubic Halide Perovskites. This dataset builds on that of 229 samples which formed the foundation of the prior work by ?. The density functional theory (DFT) computed properties we collected and the levels of theory used are discussed in Methodology. The relatively large size of this dataset is intended to provide an initial sampling suitable for guiding exploration of the alloy space. Structural information is considered constant to better focus on obtaining physically meaningful interpretations of models dealing only with information derived from a sample composition.

literature review about grouping/unsupervised learning

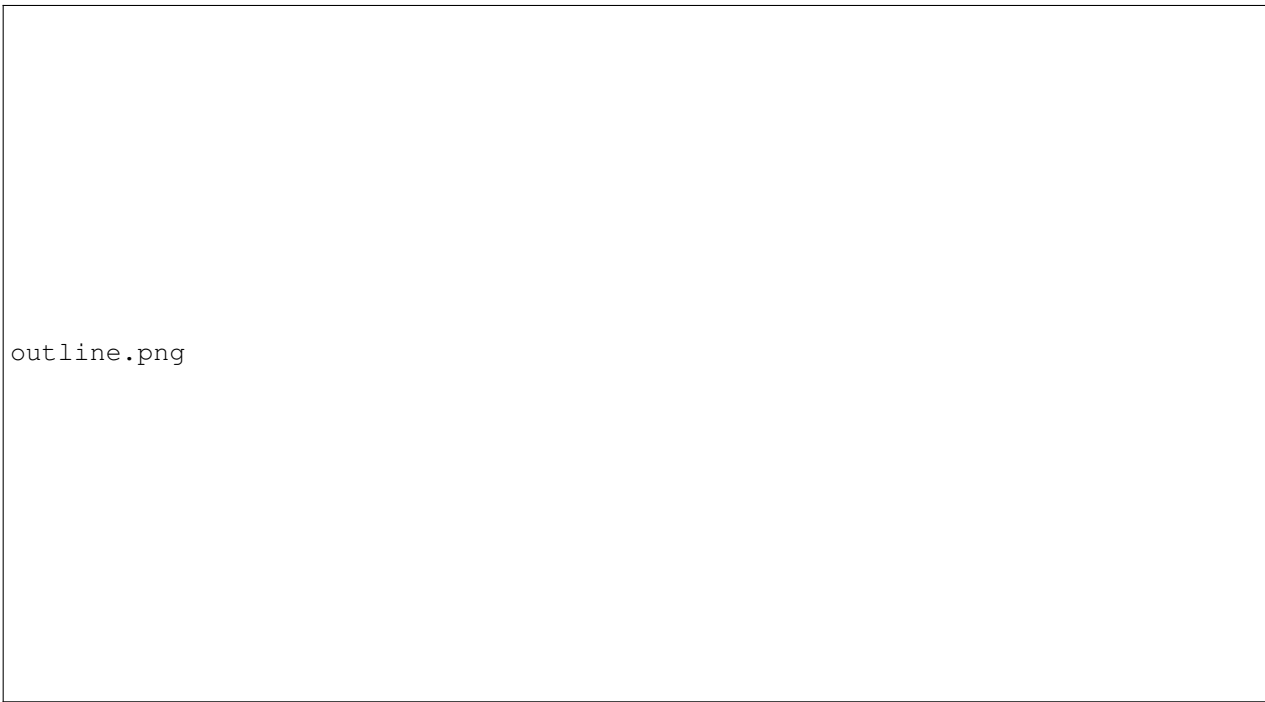
## TODO Methodology

### NEXT Building Perovskite Dataset

- State "NEXT" from "TODO" [2022-06-05 Sun 22:11]

<sup>a</sup>School of Materials Engineering, Purdue University, West Lafayette, IN 47907, USA; E-mail: amannodi@purdue.edu

<sup>†</sup>Electronic Supplementary Information (ESI) available: [https://www.github.com/PanayotisManganaris/REPO\\_ODO.SeeDOI](https://www.github.com/PanayotisManganaris/REPO_ODO.SeeDOI) : 00.0000/00000000.



outline.png

**Fig. 1** (a) Chemical space of  $ABX_3$  perovskites. (b) Number of samples representing each kind of primary alloy. (c) Detailed outline.

The dataset we report is based on standard cubic phase  $ABX_3$  Perovskite structures obtained from public databases. Fourteen common Perovskite constituents are selected to form our Halide Perovskite composition space (Figure 1). Five constituents including MA and FA cations represent the possible A-site occupants. Six metals represent the possible B-site occupants. Three halides represent the possible X-site occupants. In total, these component vectors form a constrained 14 dimensional space (Figure 2) within which all Perovskite compounds consisting of the elements in Figure 1 (a) must exist.

The pure (non-alloyed) possibilities are exhaustively sampled using  $5 * 6 * 3 = 90$  Perovskites. Based on these pure Perovskite structures, we mix candidates for A, B, and X sites systematically. The alloy space sees combinatorial scaling and must be sparsely sampled (Figure ??).

To generate site-mixing structures with high representativeness, special quasi-random structures (SQS) are applied. The SQS method is to build a special periodic structure and make the first nearest-neighbor shells as similar to the target random alloy as possible. The SQS can be considered the best possible periodic unit cell representing a given random alloy.

Each computational run is performed using a  $2x2x2$  supercell, this allows A and B site doping to be modeled in discrete  $1/8^{\text{th}}$  fractions of the total site occupancy, and it allows X site doping to be modeled in  $1/24^{\text{th}}$  fractions. At these mixing levels, it is appropriate to call all these Per-

ovskites alloys.

following the process above, 126 A-site mixing samples, 151 B-site mixing samples and 127 X-site mixing samples are generated. All resulting structures are optimized using a DFT variable-cell relaxation under Perdew-Burke-Ernzerhof (PBE). These same initial structures also underwent a full Heyd-Scuseria-Ernzerhof (HSE) relaxation to help ensure the validity of the PBE relaxations, however this was only computationally tractable for 299 samples.

### DONE Calculation Details

- State "DONE" from "TODO" [2022-06-06 Mon 17:17]

DFT calculations are performed with VASP version 6.2. The projector augmented wave potentials were used as pseudopotentials. The generalized gradient approximation of PBE and the hybrid HSE06 ( $\alpha=0.25$  and  $\omega=0.2$ ) functionals are used as exchange-correlation energy. The energy cutoff for the plane-wave basis is set to 500 eV. For PBE the Brillouin zone was sampled by  $6x6x6$  reciprocal mesh using the Monkhorst-Pack k-point mesh. For HSE the Brillouin zone was sampled by  $2x2x2$  reciprocal mesh using the Monkhorst-Pack k-point mesh. The structural force convergence threshold is set to be 0.005 eV/Å.

### TODO Discussion of DFT Computed Properties

- State "DONE" from "TODO" [2022-06-06 Mon 17:19]

variability\_of\_composition\_vectors.png

**Fig. 2** Plots showing number of Perovskites representing a constituent at a certain atomic fraction of a complete Perovskite.

## DONE Decomposition Energy

- State "DONE" from "TODO" [2022-06-06 Mon 00:43]

The decomposition energy indicates the stability of a compound. To calculate the decomposition energy for an ABX<sub>3</sub> Perovskite, we assume it will decompose to two phases, AX and BX<sub>2</sub>. Using DFT calculations, we can get the optimized energy of a Perovskite and that of its constituent phases. The decomposition energy is calculated using equation (1). This calculation is performed separately for each level of theory.

$$E_{decomp} = E_{opt}(ABX_3) - E_{opt}(AX) - E_{opt}(BX_2) \quad (1)$$

## TODO SLME calculations

- State "DONE" from "TODO" [2022-06-06 Mon 17:30]

The SLME metric developed by Yu and Zunger<sup>1</sup> is used as the primary criterion screening perovskites for their photovoltaic merits. The SLME value is computed considering a 5μm absorption layer for every Perovskite according to equation (2).

(2)

This calculation is performed using Logan William's SL3ME<sup>2</sup>. Based on absorption spectra obtained from VASP. This calculation is performed separately for the PBE band gaps and the HSE band gaps, resulting in two synthetic efficiencies for each record.

## TODO Results and discussion

### TODO Visualization of DFT Data

#### DONE Lattice constant PBE vs HSE

- State "DONE" from [2022-06-06 Mon 01:02]

Fig 4 presents the lattice parameter comparison of PBE calculation and HSE calculations. Full relaxations at both levels of theory mostly agree on the lattice parameters. At least, any deviation appears to be very linearly explained. This suggests the accuracy of PBE relaxation is enough to optimize most Perovskite samples.

Note this also served as a reasonable validation for our results. A few samples did have significantly differing lattice parameters. This prompted checking the optimized structures. We found those Perovskite structures were substantially deformed and no longer had obvious octahedral structures. Thus, we exclude these outliers from any analysis concerned with the dominant pseudo cubic structures which are the focus of this report.

Figure2.png

**Fig. 3** DFT Results: PBE and HSE properties; lattice constants, decomposition energies, band gaps, and converter efficiencies.

pbe\_v\_hse\_LC.png

**Fig. 4** Comparing lattice constants obtain by full PBE and HSE relax calculations

#### TODO Comparing synthetic to physical data

- State "DONE" from [2022-06-06 Mon 01:01]

pbe\_v\_exp\_LC.png

**Fig. 5** Comparison of PBE and HSE computed pseudo-cubic lattice constants with crystallographic measures for lattice constants

hse\_v\_exp\_LC.png

**Fig. 6** Comparison of PBE and HSE computed pseudo-cubic lattice constants with crystallographic measures for lattice constants

The physical data used in the Lattice Constant Comparisons is collected from the works of Briones *et al.*<sup>3</sup>, Jiang *et al.*<sup>4</sup>, Chen *et al.*<sup>5</sup>.

The synthetic lattice constants do mostly agree with experiment. The PBE lattice constants are better than the HSE measures 1.

It should be noted that the physical FASnI<sub>3</sub> measure reported by Chen *et al.*<sup>5</sup> is clearly an orthorhombic phase. The validating lattice constant here is obtained by averaging lattice parameters. However, per the Deviation from Cubicity metric, this phase is still approximately cubic, with angles near enough to 90 degrees that we consider the data point valid. However, this does explain the disagreement in the parity plots.

**Table 1** Root Mean Squared Error

HSE v EXP	PBE v EXP
-----------	-----------

**TODO Decomposition energy vs band gap PBE and HSE**

Fig 2(b) is showing the PBE band gap compared to the PBE decomposition energy. It presents the diversity of our Perovskite dataset. The data sets cover a large range of band gap and decomposition energy. For example, we have Perovskite with low decomposition energy (good stability) and suitable band gap value (between 1 eV to 2.5 eV for PBE calculations). And we can also find samples with low stability and large band gap. The distribution of band gap and decomposition energy shows a great diversity of all Perovskite samples and indicates that our data set can statistically represent a sufficient Perovskite space. Fig 2(c) shows the plot of PBE band gaps and HSE decomposition energy. Since we applied HSE calculations for part of the samples, it shows some grouping on low decomposition energy. Compared to PBE plot, more data should be added in high decomposition energy region and suitable band gap region.

**TODO Spectroscopic Limited Maximum Efficiency (SLME) vs PBE Band gap**

Fig 2(d) presents the SLME values related to the PBE band gap. SLME is a proven proxy for photovoltaic performance<sup>1</sup>. SLME measures the absorption efficiency of light for the Perovskite. As Fig 2(d) showing, a peak around 1.5 eV is obvious. The peak indicates that these samples with 1.5 eV PBE band gap will also have best absorption efficiency as photovoltaic materials. As the band gap increases, the SLME value decreases and eventually goes to zero due to the high band gap values.

**TODO Pearson Correlation Results**

No target is adequately explained by a single composition or composition derived axis. But there are helpful relations that aid understanding the significance of key composition derived descriptors.

A Pearson correlation map is produced to check for strong relations. Those that exist, when plotted in detail show some trending, but always with extensive variability. Evidently, an accurate model will have to be formed on a multidimensional domain.

**TODO PCA**

PCA is a method of projecting high dimensional data onto a plane defined by the two linear combinations of axes that

explain as much of the variance as possible.

The method of PCA is the Singular Value Decomposition, a Unitary Transform which generalizes the familiar eigen-decomposition. PCA will "rotate" the N data points in M-D space until their widest 2D cross section is displayed.

At this point it is readily apparent that this dataset is highly topological. The data exists on a mostly bounded domain in high dimensions, so there is some geometry the features constitute.

Our models will prefer to use this this geometric structure in their explanation for why Perovskite properties vary, this can be useful for accuracy, it can also be a bias-inducing hindrance.

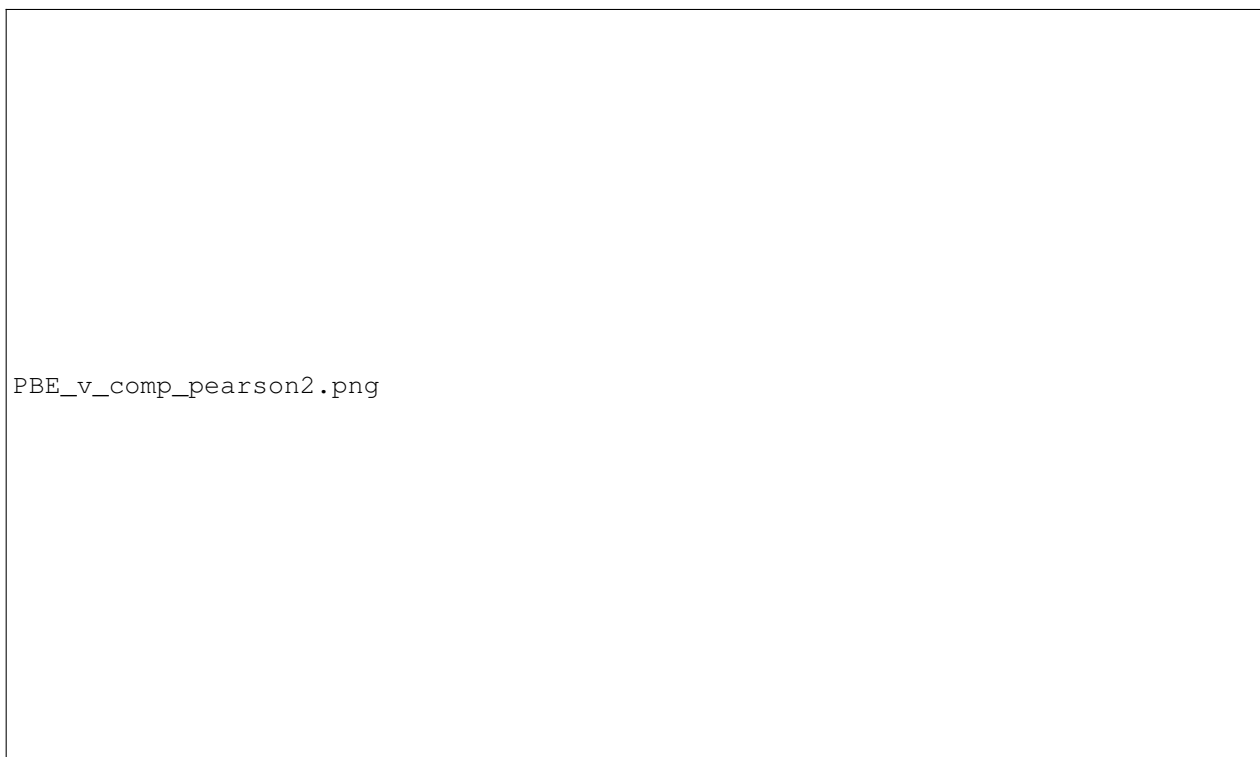
**TODO t-SNE visualization of stability clustering**  
**DONE Screening process of all samples in database**

- State "DONE" from [2022-05-24 Tue 13:31]

The synthetic Perovskite data is collected for sampling variety and without regard for each compound's viability. In order to extract potentially high-performing candidates for synthesis, we screen all samples according to certain constraints, thereby obtaining some Perovskites worth examining with higher Level of Theory DFT calculations and future physical experiments by brute force. This dataset can only confidently claim to sample cubic phase perovskites. Therefore, a promising candidate should have low deformation. Additionally, a lower decomposition energy combined with a maximal SLME value would be ideal for physical testing. Our screening procedure consists of cutting based on a custom Deviation from Cubicity metric, octahedral factor, Goldschmidt tolerance factor, Bartel tolerance factor, decomposition energy, and band gap – acting as less stringent proxy for the SLME spectrum. The following sections will discuss the details of each constraint.

**TODO Deviation from Cubicity**

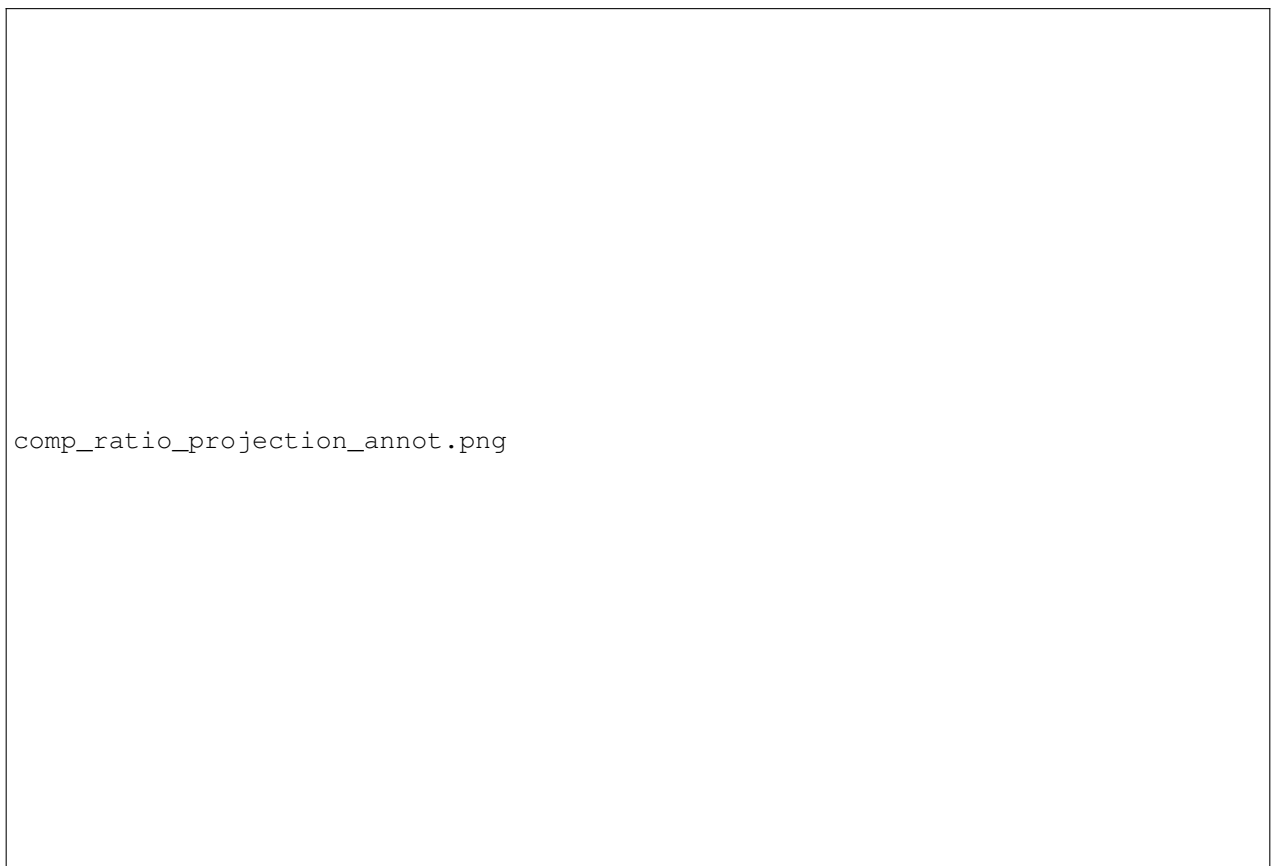
For all the compositions we tested in our data set, some of them will have large strain and deformation because of the combination of elements. For these largely deformed samples, they are no longer remain a cubic perovskite structure. In this section, we want to analyst how the structure is apart from the cubic perovskite structure and rule out these largely deformed samples. Firstly, we need to define deviation of cubicity. For each perovskite sample, we measure the difference between b, c lattice parameter against a lattice, showing in Equation X. If the lattice deviation of cubicity is larger than 10%, we will consider this perovskite is no longer remain a cubic perovskite and exclude it. Similarly, we also need to consider 3 angles,  $\alpha$ ,  $\beta$ , and  $\gamma$ , to make sure the angle also remain 90 degrees. We calculated the



**Fig. 7** Pearson linear correlation coefficients between 14 composition variables, and 4 PBE computed properties

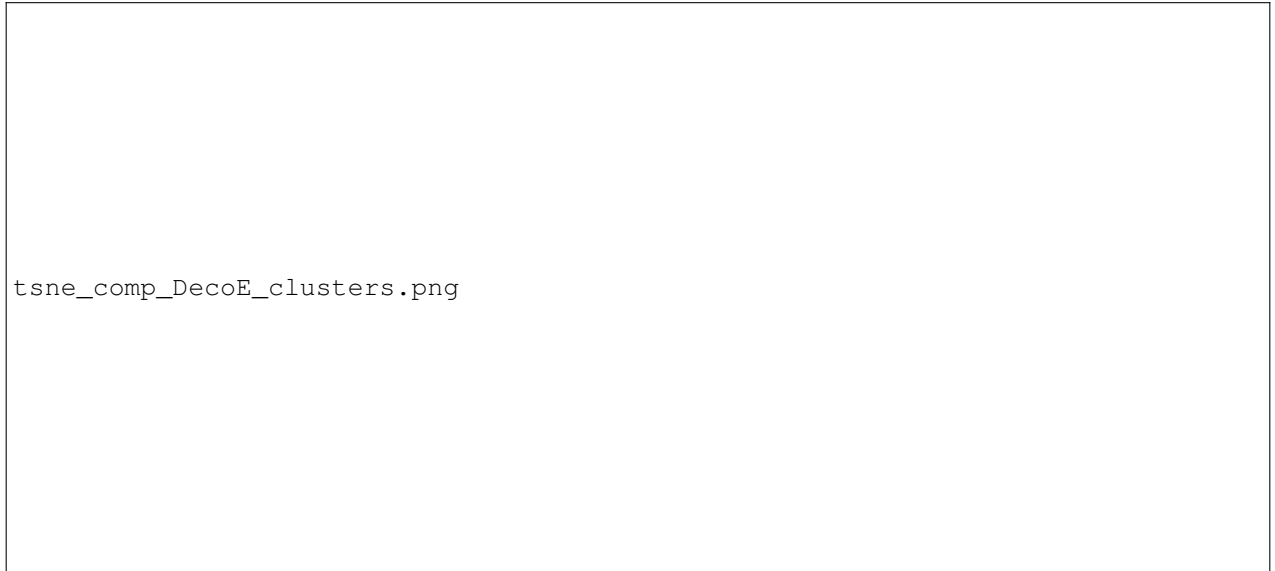


**Fig. 8** Pearson linear correlation coefficients between 36 elemental descriptors and 4 PBE computed properties

A large rectangular box representing a PCA plot. The plot area is mostly empty, with the filename 'comp\_ratio\_projection\_annot.png' centered in the lower-left portion.

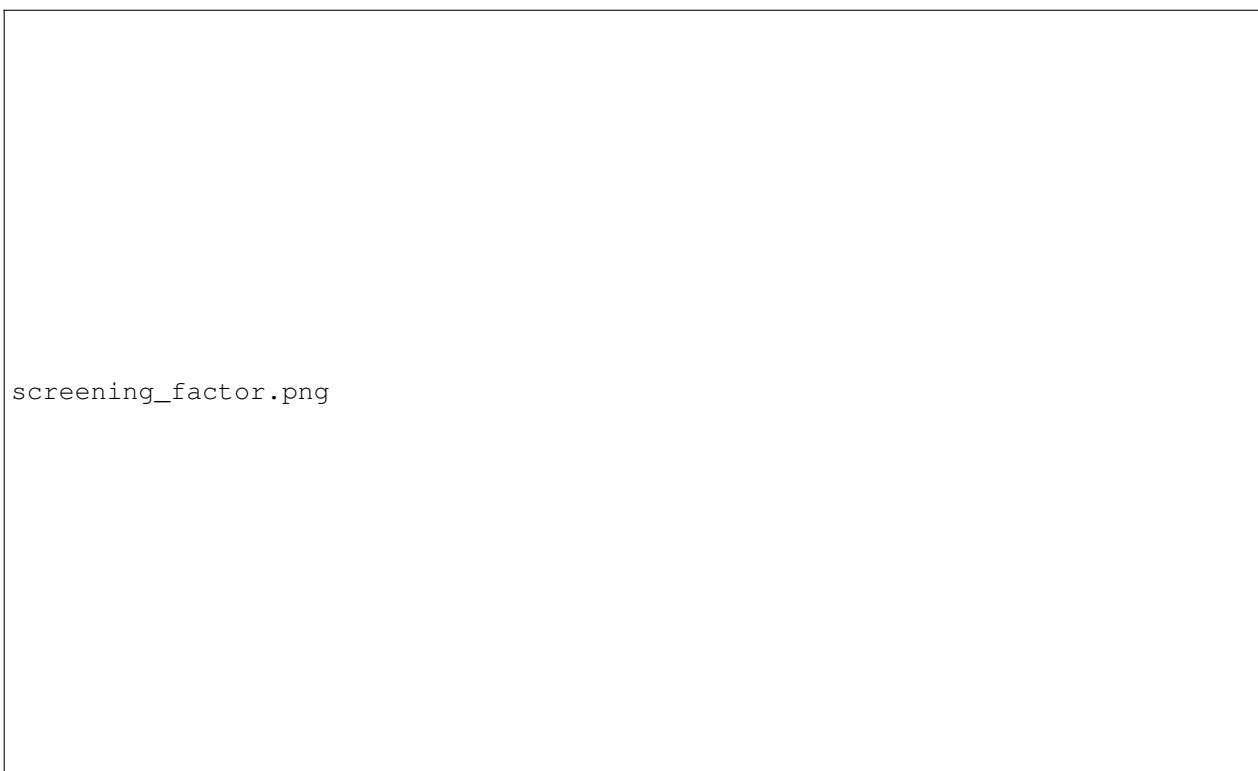
comp\_ratio\_projection\_annot.png

**Fig. 9** PCA

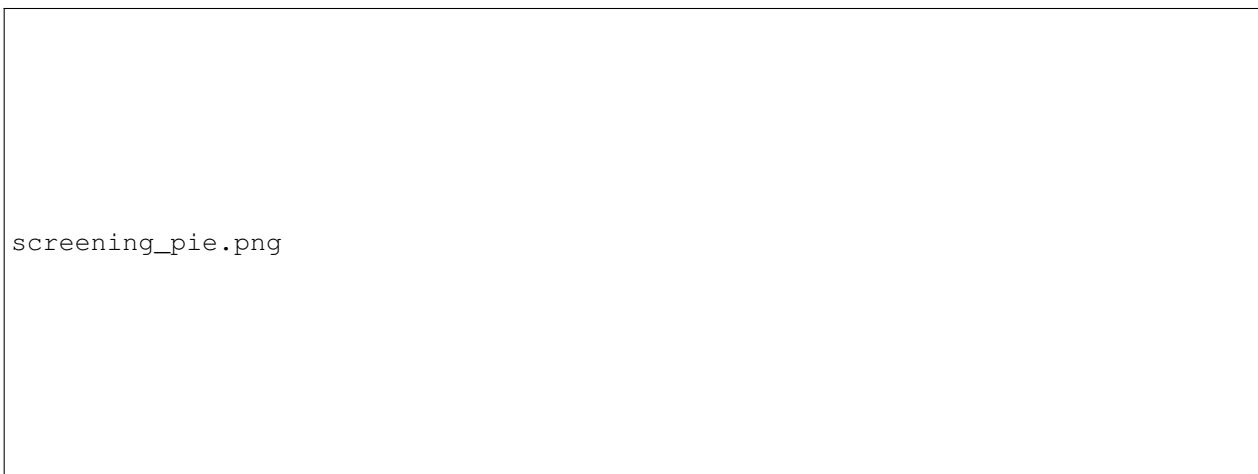
A large rectangular box representing a t-SNE plot. The plot area is mostly empty, with the filename 'tsne\_comp\_DecoE\_clusters.png' centered in the lower-left portion.

tsne\_comp\_DecoE\_clusters.png

**Fig. 10** t-SNE

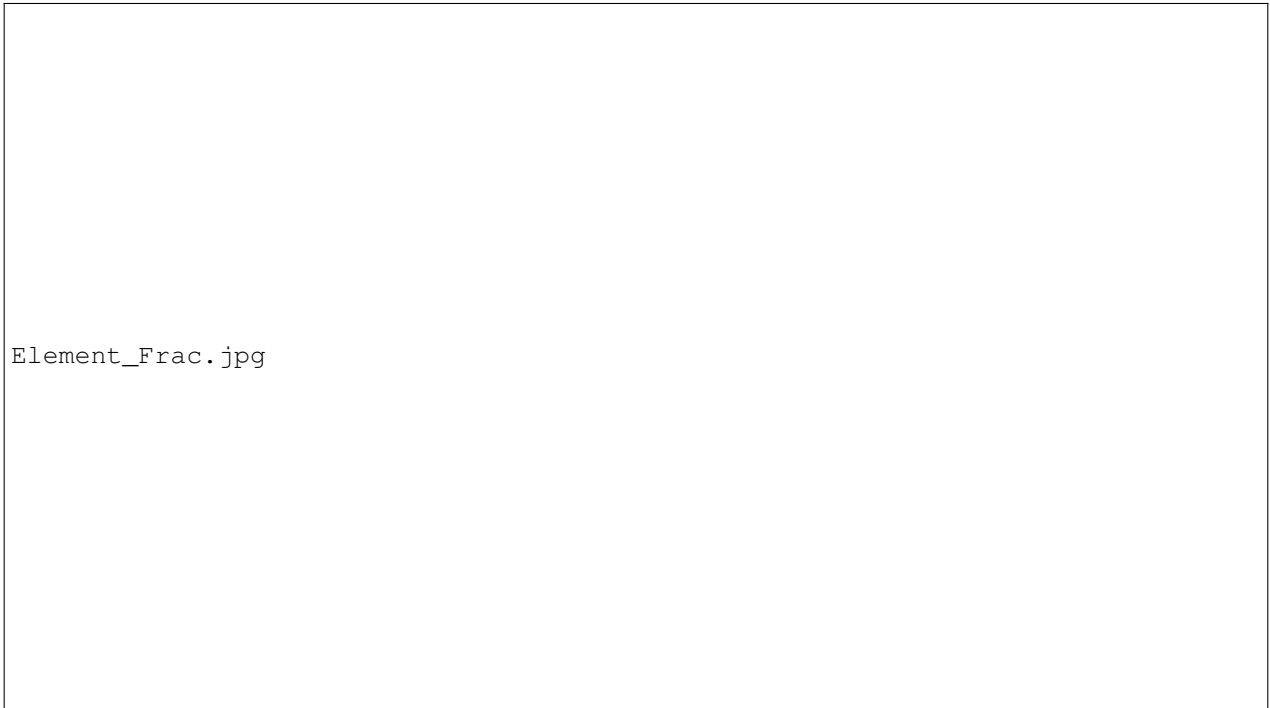


**Fig. 11** Screening results of Decomposition Energy for (a) octahedral factor (b) Tolerance Factor (c) Bartel Tolerance Factor (d) Deviation of Cubicity.



**Fig. 12** 49 Screened perovskites analysis, (a) pie chart of element distribution, (b) pie chart of site mixing





Element\_Frac.jpg

**Fig. 13** Plots of element fraction for A, B, X sites of all screened perovskites.

difference of  $\theta$ , and versus 90 degrees, showing in Equation X, and take this as angular deviation of cubicity. We also consider the samples that have more than 5% of angular deviation of cubicity are not cubic perovskites and exclude them. Fig 7 (d) shows the screening results for angle deviation of cubicity. SI Fig X shows the screening results for b, c lattice and  $\theta$ , angle. Most of the samples with non-cubic perovskites structures are organic-inorganic hybrid perovskites. A large part of the excluded samples are A-mixing hybrid perovskites. It indicates that organic ligands in A site sometimes increase the lattice along some direction and make the perovskite deformed.

#### **TODO Octahedral factor, tolerance factor, and Bartel tolerance factor**

The stability of Perovskite can be predicted by using the atom radius of all components. There are 3 types of factors are usually considered, Octahedral factor, tolerance factor, and Bartel tolerance factor. The formula of these three factors are shown in equation XXX. In our screening process, we set the criteria for Octahedral factor as 0.442 – 0.895. The criteria for tolerance factor is set to be 0.813 – 1.107. The criteria for Bartel tolerance factor is set to be less than 4.18. Fig 7 (a) shows the Octahedral factor versus decomposition energy plot. Fig 7 (b) shows the tolerance factor versus decomposition energy plot. The tolerance factor shows a trend that as the tolerance factor increases, the decomposition energy decreases. Fig 7 (c) shows the Bartel

tolerance factor versus decomposition energy plot. Within the criteria of Bartel tolerance factor, the decomposition energy is rapidly decreased.

#### **TODO Screening Results**

The table XXX shows all 41 promising candidates we screened from our dataset by using the cutoff criteria mentioned above. FigXXX shows the element distribution of all decent samples. It indicates that MA and FA are the most common A site choices for a Perovskite with great stability and suitable band gap and photovoltaic properties. For inorganic A sites, Cs also contributes a large portion. In B site element selection, Pb and Sn are the most selected than other B site elements. For X site halide elements, Br is appears the most, and I is the second. It tells us that Br and I are preferred in a stable Perovskite structure. In figure XXX, the mixing distribution indicates that X-site mixing is more favorable to generate a stable Perovskite with required properties. Tuning in A sites and B sites can also lead to a promising results but more risky. The FigXXX and SI FIGXXX shows the element fraction for each of the screened samples. It indicates if certain element prefer to be mixing or non-mixing. For example, Cl is preferred to be pure in X site than mixing with other elements. Also there is some elements that only appears in pure case, for instance Sr at B site. We can also learn from the result that elements, like K, Cs, Cl and Ba, are more preferred to have small fraction if they are mixing with other elements.

Popular choice, like Pb and Sn, are very flexible in element fraction. These elements can have a variety of possible fraction in promising Perovskite samples. The screening results are showing the trends of how the combinations of elements are selected and what mixing features are preferred in decent Perovskite structures. It will also provide great help on design new Perovskite compounds and tell us what kinds of compounds we should put in to test next.

## TODO Perspective and Future Work

The high-throughput DFT Perovskite dataset is containing huge amount of information of perovskites compounds, structure and properties. It leads us to perform screening, clustering and other process so that lots of trends are revealed. However, there are still many improvements we are eager to accomplish in the future.

Firstly, the dataset are generally constructed based on cubic perovskites structures. But not all perovskites are remaining cubic as their most stable phases. In the future, we are going to expand our dataset that will include high-throughput DFT calculations for multiple phases of perovskites, for example tetragonal, orthorhombic, and hexagonal phases. This will tell us which phase we should choose for each compound and provide more information for analysis.

Also the level of DFT functional should also be included in the future work. Important corrections, like Van der Waals forces, spin orbital coupling, should also be considered for different compound. We are going to expand the dataset with information of multiple level of DFT functional and corrections, which will eventually lead us to decide the best suitable functional for each compounds and get more accurate prediction.

Structural information is another critical part for tuning the properties of Perovskite. Octahedral distortion and rotation is the main part that causes the change in Perovskite structures. We have ongoing work to reveal the relationship between octahedral information and properties. These part of information will also be include in the future dataset.

The design of this dataset is uniquely suited to the exploration of alloying effects on Perovskite properties. The combinatorial space of possible alloys has been sparsely but systematically sampled along four primary alloy schemes. This sample space affords the opportunity for a QM/ML surrogate model to form the basis of an active learning strategy which can begin selecting potentially high performing multi-site alloy candidates based on the current sample set.

For the calculation of future data points obtained by a surrogate optimizer, we will follow a strategy of perform-

ing full structural optimization at a PBE level of theory unless circumstances demand otherwise. This is justified by Figure 4 and the section Comparing synthetic to physical data

Also, we will explore methods for combining the insights of the PBE and HSE datasets in training surrogates, e.g. in section SLME calculations the two SLME spectra could be used together to converge on a physically accurate PCE value.

We anticipate our principal challenge will be in extracting useful predictor variables from the composition information. The basic feature sets examined here are highly correlated, but nonetheless show promise both as a basic screening criterion, and as good classifier features under t-SNE transformation.

Modeling pipelines capable of predicting Perovskite decomposition energy will likely be very achievable using a transduction and invertible equivalent of the t-SNE algorithm, potentially SONG.

For this reason, kernel learning methods appear to be particularly promising for high speed optimization of the space.

## TODO Conclusions

...

## Conflicts of Interest

There are no conflicts to declare.

## Acknowledgements

Extensive discussions with and scientific feedback from UC San Diego researchers David Fenning and Rishi Kumar and Argonne National Lab scientist Maria Chan are acknowledged. This work was performed at Purdue University, under startup account F.10023800.05.002 from the Materials Engineering department. This research used resources of the National Energy Research Scientific Computing Center, the Laboratory Computing Resource Center at Argonne National Laboratory, and the RCAC clusters at Purdue.

## References

- 1 L. Yu and A. Zunger, *Physical Review Letters*, 2012, **108**, year.
- 2 L. Williams, *Sl3me – a Python3 Implementation of the Spectroscopic Limited Maximum Efficiency (SLME) Analysis of Solar Absorbers*, <https://github.com/ldwillia/SL3ME>.
- 3 J. Briones, M. C. Guinto and C. M. Pelicano, *Materials Letters*, 2021, **298**, 130040.
- 4 L. Jiang, J. Guo, H. Liu, M. Zhu, X. Zhou, P. Wu and

- C. Li, *Journal of Physics and Chemistry of Solids*, 2006, **67**, 1531–1536.
- 5 Q. Chen, N. D. Marco, Y. M. Yang, T.-B. Song, C.-C. Chen, H. Zhao, Z. Hong, H. Zhou and Y. Yang, *Nano Today*, 2015, **10**, 355–396.

## References

- 1 L. Yu and A. Zunger, *Physical Review Letters*, 2012, **108**, year.
- 2 L. Williams, *Sl3me – a Python3 Implementation of the Spectroscopic Limited Maximum Efficiency (SLME) Analysis of Solar Absorbers*, <https://github.com/ldwillia/SL3ME>.
- 3 J. Briones, M. C. Guinto and C. M. Pelicano, *Materials Letters*, 2021, **298**, 130040.
- 4 L. Jiang, J. Guo, H. Liu, M. Zhu, X. Zhou, P. Wu and C. Li, *Journal of Physics and Chemistry of Solids*, 2006, **67**, 1531–1536.
- 5 Q. Chen, N. D. Marco, Y. M. Yang, T.-B. Song, C.-C. Chen, H. Zhao, Z. Hong, H. Zhou and Y. Yang, *Nano Today*, 2015, **10**, 355–396.

## Supplemental Material

### Glossary

**PCA** principal component analysis.

**QM/ML** quantum mechanics machine learning.

**SLME** spectroscopic limited maximum efficiency.

**t-SNE** t-distributed stochastic neighbor embedding.

**VASP** Vienna Ab initio Simulation Package.

### Acronyms

**DFT** density functional theory.

**FA** Formamidinium.

**GGA** generalized gradient approximation.

**HSE** Heyd-Scuseria-Ernzerhof.

**MA** Methylammonium.

**PAW** projector augmented wave.

**PBE** Perdew-Burke-Ernzerhof.

**PCE** power conversion efficiency.

**SQS** special quasi-random structures.

HSE\_v\_comp\_pearson2.png

**Fig. 14** Pearson linear correlation coefficients between 14 composition variables, and 4 HSE computed properties


HSE\_v\_site\_prop\_pearson.png

**Fig. 15** Pearson linear correlation coefficients between 36 composition variables, and 4 HSE computed properties



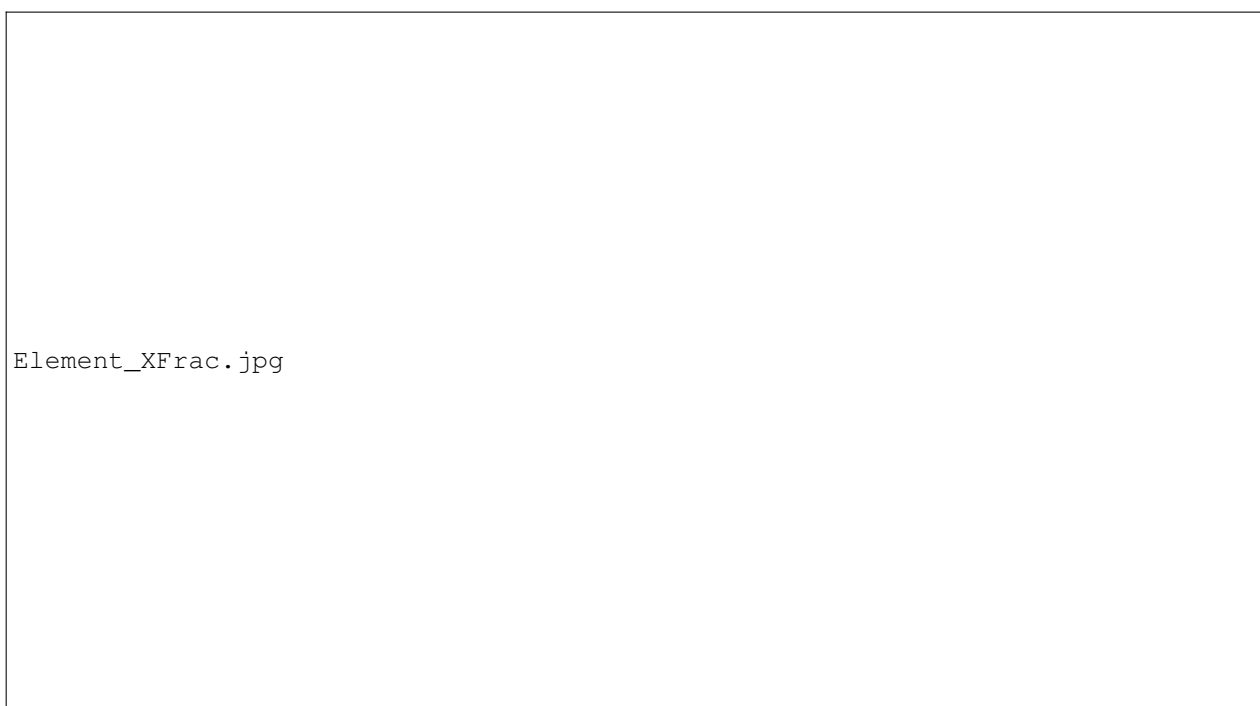
Element\_AFrac.jpg

**Fig. 16** Element Fraction for A site elements among all screened perovskites



Element\_BFrac.jpg

**Fig. 17** Element Fraction for B site elements among all screened perovskites



**Fig. 18** Element Fraction for X site elements among all screened perovskites