RESEARCH ARTICLE

ADVANCED
INTELLIGENT
SYSTEMS
Open Access

www.advintellsyst.com

# Active Learning in Bayesian Neural Networks for Bandgap Predictions of Novel Van der Waals Heterostructures

*Marco Fronzi,\* Olexandr Isayev, David A. Winkler, Joseph G. Shapter, Amanda V. Ellis, Peter C. Sherrell, Nick A. Shepelin, Alexander Corletto, and Michael J. Ford*

The bandgap is one of the most fundamental properties of condensed matter. However, an accurate calculation of its value, which could potentially allow experimentalists to identify materials suitable for device applications, is very computationally expensive. Here, active machine learning algorithms are used to leverage a limited number of accurate density functional theory calculations to robustly predict the bandgap of a very large number of novel 2D heterostructures. Using this approach, a database of ≈2.2 million bandgap values for various novel 2D van der Waals heterostructures is produced.

## 1. Introduction

The electronic bandgap (BG) is one of the fundamental properties of materials. It arises directly from the configuration of electronic structures and corresponds to the minimum energy that an electron requires to be excited into the conduction band.[1–4]

In particular, in 2D materials, the bandgap value determines fundamental physical characteristics such as optical excitation, and electron transport and transfer.[5–9] Therefore, the identification and the control of the bandgap in van der Waals 2D heterostructures (vdWHs) is a viable strategy for design of novel materials for a large variety of electro-optic devices.[10–17]

Depending on the definition used, the bandgap can have more than one meaning. The optical bandgap refers to the minimum energy that an electron must absorb to be excited to the conduction band, where the conduction band minimum and the valence band maximum align, forming an electron–hole pair. If the bands do not align, then the same quantity is referred to as the electrical or transport bandgap, which in general represents the energy threshold for creating an electron–hole pair that is not bound together. In these cases, the optical bandgap is smaller

M. Fronzi
SIT Research Laboratories
Shibaura Institute of Technology
3-7-5, Toyosu, Koto-ku, Tokyo 135-8548, Japan
E-mail: m-fronzi@shibaura-it.ac.jp

M. Fronzi, M. J. Ford
School of Mathematical and Physical Science
University of Technology Sydney
Ultimo, NSW 2007, Australia

O. Isayev
Department of Chemistry
Mellon College of Science
Carnegie Mellon University
Pittsburgh, PA 15219, USA

D. A. Winkler
Monash Institute of Pharmaceutical Sciences
Monash University
381 Royal Parade, Parkville, VIC 3052, Australia

The ORCID identification number(s) for the author(s) of this article can be found under https://doi.org/10.1002/aisy.202100080.

© 2021 The Authors. Advanced Intelligent Systems published by Wiley-VCH GmbH. This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

DOI: 10.1002/aisy.202100080

D. A. Winkler
Department of Biochemistry and Genetics
La Trobe Institute for Molecular Science
La Trobe University
Kingsbury Drive, Bundoora, VIC 3086, Australia

D. A. Winkler
School of Pharmacy
The University of Nottingham
Nottingham NG7 2RD, UK

J. G. Shapter, A. Corletto
Australian Institute for Bioengineering and Nanotechnology
The University of Queensland
St Lucia, Brisbane Qld 4072, Australia

J. G. Shapter
College of Science and Engineering
Flinders University
Bedford Park, Adelaide, SA 5042, Australia

A. V. Ellis, P. C. Sherrell, N. A. Shepelin
Department of Chemical Engineering
University of Melbourne
Parkville, VIC 3010, Australia

N. A. Shepelin
Laboratory for Multiscale Materials Experiments
Paul Scherrer Institut
Villigen CH-5232, Switzerland

than the transport bandgap. These two definitions have physical meaning and correspond to quantities that can be measured experimentally. In contrast, the fundamental bandgap refers to the energy separating unoccupied from occupied one-electron states, and it is meaningful only within a theoretical model. In general, the experimental bandgap does not correspond to the fundamental bandgap. Some approaches that go beyond density functional theory (DFT) have been developed to consider electron excitation for bandgap calculations. However, within the commonly used DFT approximations (e.g., generalized gradient approximation [GGA], meta-GGA, and hybrid functionals used to solve the Kohn–Sham equations for the ground state), optical and fundamental bandgaps are equivalent.[18,19] Furthermore, some studies posit a weak electron–hole binding energy in selected vdWHS and, therefore, to the first approximation they can be treated as the same quantity.[20]

For the simplest approximation of the exchange correlation functional, DFT largely underestimates the value of the bandgap. When calculated within the local density approximation (LDA), the values are typically ≈40% less than experimental values.[21,22] To improve the accuracy of the calculations, functionals considering the gradient of the charge density (GGA) are used. However, significant underestimation of the bandgap values persists and, in many cases, it leads to an incorrect description of the electronic and optical properties of a material.[23–25] To improve the prediction of bandgap calculations, Hartree–Fock exact exchange interactions can be included in the Kohn–Sham equations, leading to the so-called hybrid functional approximations.[26,27] Heyd et al. developed one of the most widely used functionals (HSE06) for predicting electronic structures and properties of solid-state materials, including bandgaps.[28,29] However, the inclusion of exact exchange leads to significantly higher computational costs compared with simpler LDA or GGA functionals, a very important issue when a large number of calculations is required.[30,31]

Due to the large number of theoretically possible 2D monolayers (≈6000 structures), it is possible to generate ≈20 million unique novel bilayer heterostructures by a direct stacking of these 2D monolayers $(N_b = N_m(N_m + 1)/2)$.[https://2dmatpedia.org/),[32]] Calculation of the electronic structures and properties of this many complex materials is currently intractable, even for simple DFT calculations, and the lack of experimentally fabricated vdWHs makes validation challenging.

We have previously shown how machine learning (ML) models can be used to leverage results from DFT calculations of hybrid 2D materials, generating outstanding results.[32] Here we demonstrate how a combination of a relatively small number of computationally expensive DFT-HSE06 calculations and a ML method called active learning (AL) can reliably predict bandgap values for a large set of stable semiconducting bilayers.[32] We show how these ML models can be used to create a database of reliable bandgap values for approximately 2.2 million vdWHs, using a very limited number of first-principles calculations (≈400 bilayers).

## 2. Results

Here, an AL model was built starting from an initial training set of 109 randomly selected bilayers ($X_L$), where each bilayer
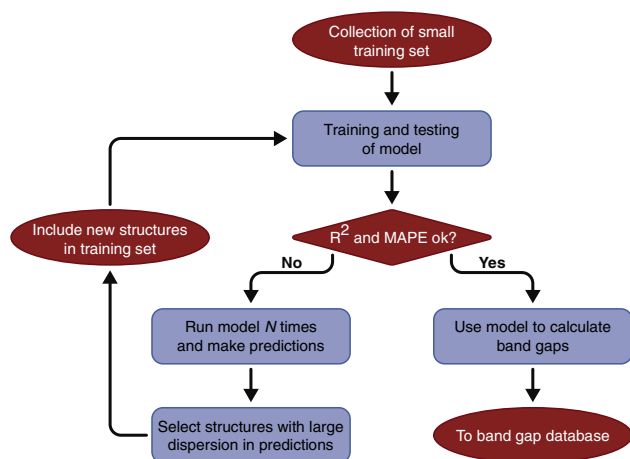
**Table 1.** Bandgap (in eV) comparison between the value calculated here using HSE06 and other HSE06 calculations found in the literature.[58–60] Discrepancies may result from different methodologies used in the approximation of the vdW potential and/or configuration of the supercells used.

| Bilayer | Bandgap [eV] This work | Bandgap [eV] Other works |
| --- | --- | --- |
| $O_2Pt–NiO_2$ | 1.72 | 1.39[58] |
| $OTl_2–GeI_2$ | 1.83 | 1.45[58] |
| $Br_2Mg–Cl_2Zn$ | 5.50 | 5.49[58] |
| $Cl_2Zn–CdCl_2$ | 5.10 | 5.29[58] |
| $OTl_2–O_2Pt$ | 0.62 | 0.86[58] |
| $I_2Yb–Br_2Ge$ | 1.30 | 1.04[58] |
| InSe/AsP | 1.26 | 1.07[59] |
| $HfS_2/MoTe_2$ | 0.59 | 0.35[60] |

consists of two semiconducting monolayers from the 2matpedia database.

Members of this set have an IE, $E \leq -1.0 \, eV \, Å^{-2}$, ensuring that the interactions are vdW, and contain a subset with HSE06 bandgaps relatively close to literature values (see **Table 1**). Possible discrepancies between our results and the literature data may originate from the use of different computational methodologies and/or from configuration of the supercells used in the calculations.

The 109 HSE06 bandgap calculations were used to build the initial Bayesian neural network (BNN) model. Five different versions of the initial model were produced using different selections of the train-test sets, selected by clustering using the k-means algorithm.[32] For each bilayer, bandgap values were calculated using the five versions of the initial model, resulting in five different databases consisting of 2.2M structures ($Y_1, Y_2,... Y_5$). After calculating the mean and standard deviation of each bilayer bandgap, a set of ≈100 bilayers ($X_{AL1}$) that had the largest standard deviations were selected. In other words, the initial BNN models could not find a mapping function of sufficient accuracy for that set of bilayers. The process, shown schematically in **Figure 1**, was repeated four times, each time selecting a new set of ≈100 poorly predicted bilayers ($X_{AL1}...X_{AL3}$), until the quality of the model reached convergence. At that point, the training set contained 473 structures whose bandgaps were predicted with an $R^2$ of 0.81 and mean absolute percentage error (MAPE) of 0.16, and the test set was predicted with an $R^2$ of 0.92 and MAPE of 0.11. An additional set of 52 structures ($X_{AL4}$) was subsequently added to the training set to test for convergence of the parameters (**Table 2**). The $X_{AL}$ subsets used to expand the training set are represented in the uniform manifold approximation and projection (UMAP) in **Figure 2**, showing how the DFT calculations are distributed over the whole set of bilayers. UMAP provides a 2D, visual, and intuitive representation of possible structural–functional similarities of the structures in the hyperspace of descriptors. Points that are close together in the UMAP correspond to structures with similar physicochemical properties. As the training set is approximately uniformly distributed

**Figure 1.** Flowchart of the AL model used to build the BNN model. The structures included in the training set of iteration $M+1$ were selected from the worst performing, after the evaluation of the mean value and the standard deviation calculated for each value from $N$ runs of the model at the iteration $M$.



**Figure 2.** Uniform manifold approximation and projection showing the distribution of DFT (dark blue) calculated bandgaps over the whole 2.2M dataset (light gray).
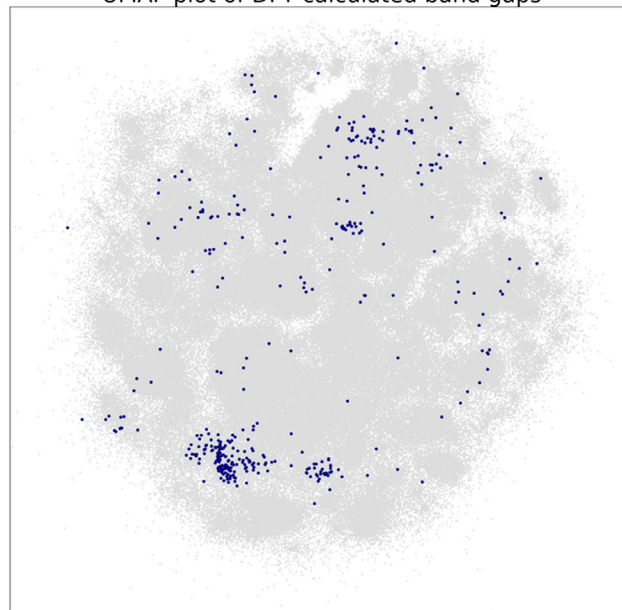
across the UMAC, the 2.2M structures predicted by the ML models lie within the domain of applicability for the model, so should be predicted with reasonable accuracy.

To demonstrate how efficient the AL model was we generated an additional BNN model trained on 425 randomly selected bilayers. This model did not predict the training and test set well, with no $R^2$ values greater than 0.51. This indicated the fundamental role that AL played in reducing the computational cost

**Table 2.** $R^2$, root mean-squared error (RMSE), mean absolute error (MAE), and mean absolute percentage error (MAPE) test and train set predictions for the BNN bandgap models. The results are labeled progressively by four steps where each step adds additional data point sets ($X_{AL1}...X_{AL3}$) to the initial 109 bilayers, selected using an AL algorithm. The fifth run was carried out using additional 52 bilayers ($X_{AL4}$) to test the convergence of the parameters.

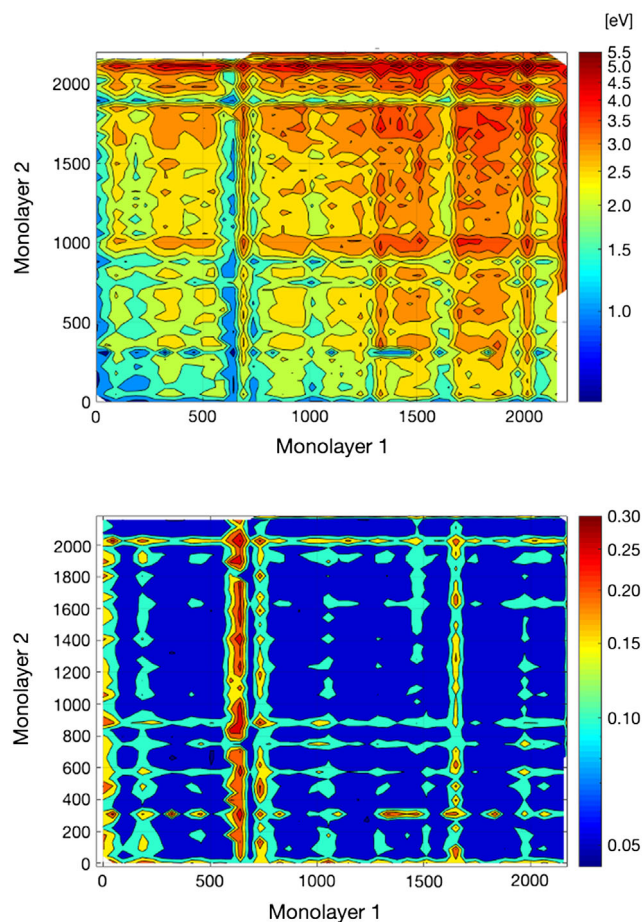| Set | $R^2$ | RMSE [eV] | MAE [eV] | MAPE [%] |
|---|---|---|---|---|
| First run ($X_L$) | | | | |
| BNN-test | 0.37 | 0.92 | 0.66 | 0.6 |
| BNN-train | 0.75 | 0.51 | 0.34 | 0.4 |
| Second run ($X_{AL1}$) | | | | |
| BNN-test | 0.51 | 0.75 | 0.60 | 0.4 |
| BNN-train | 0.77 | 0.51 | 0.35 | 0.3 |
| Third run ($X_{AL2}$) | | | | |
| BNN-test | 0.71 | 0.74 | 0.65 | 0.3 |
| BNN-train | 0.82 | 0.53 | 0.41 | 0.2 |
| Fourth run ($X_{AL3}$) | | | | |
| BNN-test | 0.81 | 0.45 | 0.31 | 0.2 |
| BNN-train | 0.92 | 0.41 | 0.28 | 0.1 |
| Fifth run ($X_{AL4}$) | | | | |
| BNN-test | 0.80 | 0.44 | 0.30 | 0.2 |
| BNN-train | 0.93 | 0.40 | 0.28 | 0.1 |

and increasing the accuracy when leveraging HSE06 calculations of 2D heterostructures.

The distribution of bandgaps and relative errors associated with the predictions, calculated from 600 trial Bayesian networks, as a function of the monolayer building blocks, is shown as a heatmap in **Figure 3**. This shows a logarithmic distribution of bandgaps over the 0.1–8.0 eV range, with ≈9% of the structures having a direct bandgap configuration. Notably, it also indicates a relatively large error (but still below 30%) associated with low bandgap values, apparent around the monolayer $1 = 650–750$. Although the error appears to be associated with specific monolayers, this is due to the sampling used to build the training set that explores the hyperspace of bilayer descriptors in a partly inhomogeneous way. Therefore, any relationships between the BNN error and the nature of chemical bonds or chemical composition of the bilayers are not significant. This is also confirmed by the distribution of bandgaps and largest error structures in the UMAP across the whole set (see **Figure 4**).

The bilayers were grouped by bandgap energy and labeled with their position in the optical spectrum as follows: infrared (IR) $\leq 1.65$ eV $\leq$ red (R) $\leq 1.99$ eV $\leq$ orange (O) $\leq 2.10$ eV $\leq$ yellow (Y) $\leq 2.17$ eV $\leq$ green (G) $\leq 2.50$ eV $\leq$ blue (B) $\leq 2.75$ eV $\leq$ violet (V) $\leq 3.26$ eV $\leq$ ultraviolet (UV). The number of bilayers within each range is shown in **Table 3**.

Within each bandgap group, we counted the frequency of each monolayer in the 2.2M bilayer dataset and shows the most frequent representatives in **Table 4**.
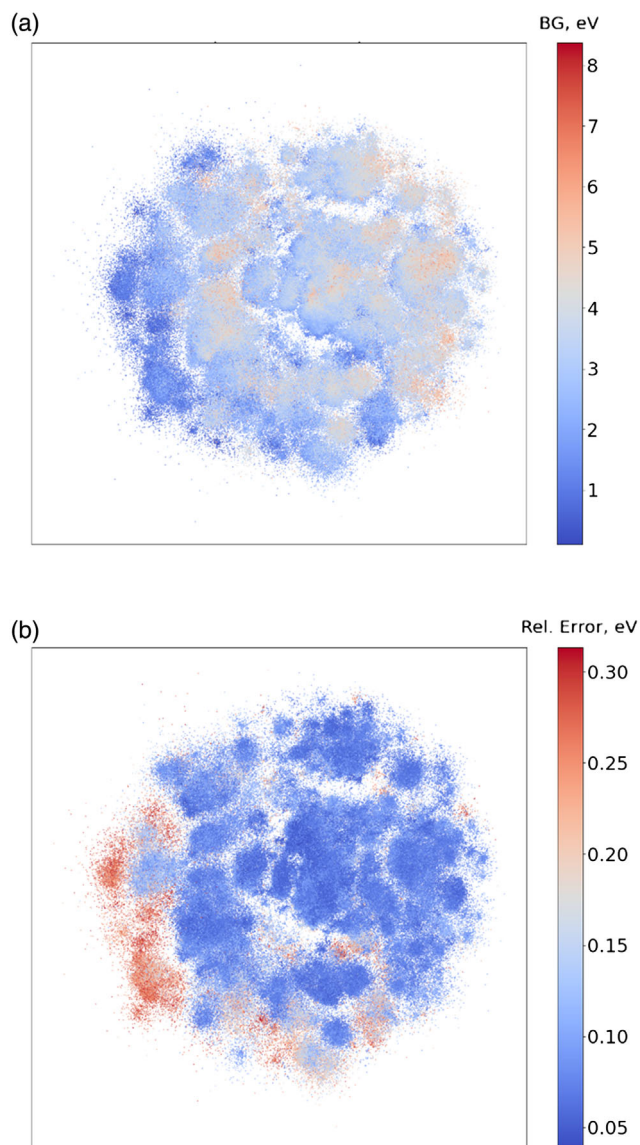
To elucidate the contribution to the bandgap due to inclusion of the exchange term in the DFT functionals, the correlation between the bandgaps calculated within the GGA-PBE and HSE06 approximations was assessed (see **Figure 5**). The

**ADVANCED
SCIENCE NEWS**

www.advancedsciencenews.com

**ADVANCED
INTELLIGENT
SYSTEMS**
Open Access

www.advintellsyst.com

**Figure 3.** Bandgap (top) and relative error (bottom) of the bilayers as a function of the two monolayers building blocks. Absolute errors have been calculated as the standard deviation of the response distribution, using a dropout approach with probability 0.1. Detailed information can be found in Bayesian Neural Networks section. The heatmaps have been generated by interpolating the function $BG = f(x,y)$ and so that the images can provide information by showing potential clustering.



**Figure 4.** Distribution of the calculated a) bandgap values and b) relative error in the UMAP.

Pearson and Spearman correlations between PBE and HSE06 bandgaps were 0.68 and 0.60, respectively, indicating a significant (linear) correlation between the results of the two approaches. The feature importance (FI) of each descriptor and the one of the GGA-PBE bandgaps were calculated, and the results show that the GGA-PBE bandgap FI is only marginally higher than the largest FI of the other descriptors (0.13 vs 0.08). In addition, a least absolute shrinkage and selection operator (LASSO) regression analysis indicated that the GGA-PBE calculated bandgap is a less-effective descriptor in the BNN models compared with the structural descriptors used here, validating the relevance of the descriptors used to represent the bilayers in this study (listed in the Supporting Information).[32–34]

We also assessed the change in the bandgaps as a function of the twist angle. Monolayers with different symmetries along the $x$–$y$ plane (AgI, AlTe, BN, $Br_3Cr$, $ITaTe_4$, and GaSe) were selected, and supercells with twist angles of 0°, 30°, 45°, 60°, and 90° were built. The calculated total energy of each structure

elucidated the resultant effect of charge transfer on the bandgap. The charge transfer suggested a relatively small change at different twist angles because of the weak vdW interactions, resulting in a negligible bandgap change (shown in the Supporting Information).
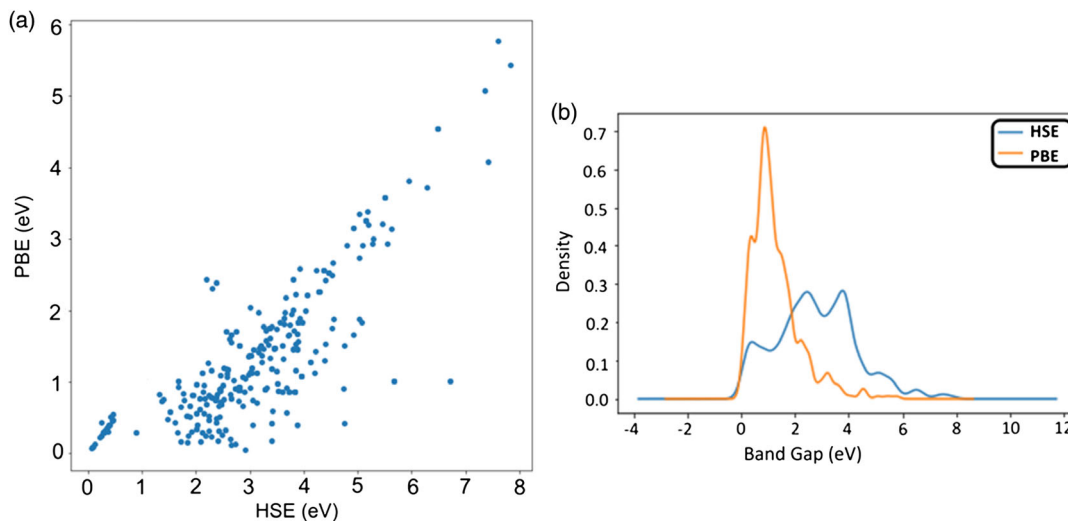
Although the bandgap change with the twist angle was small, care must be taken in interpreting this outcome because of the limited number of twist angles that produce a supercell small enough to make the calculations tractable. The literature suggests that for specific twist angles, the bandgap can vary up to ≈15%, providing a novel way to use geometric parameters of a bilayer to adjust the electronic properties.[35–38] Our work identified the crucial role of crystal symmetries in the band structure configuration. The point symmetries, rotation, reflection, and inversion determine the nature of the bands, supporting the validity of the

**Table 3.** Number of bilayers in the 2.2M set divided in groups depending on the on the position of the vdWHs bandgap in the visible spectrum. Here, the energy ranges as labeled as follows: infrared (IR) $\leq 1.65$ eV $\leq$ red (R) $\leq 1.99$ eV $\leq$ orange (O) $\leq 2.10$ eV $\leq$ yellow (Y) $\leq 2.17$ eV $\leq$ green (G) $\leq 2.50$ eV $\leq$ blue (B) $\leq 2.75$ eV $\leq$ violet (V) $\leq 3.26$ eV $\leq$ ultraviolet (UV).

| Optical spectrum | IR | R | O | Y | G | B | V | UV |
|---|---|---|---|---|---|---|---|---|
| Number of bilayers in band | 307 100 | 217 180 | 86 089 | 59 342 | 320 450 | 267 739 | 474 231 | 538 011 |

**Table 4.** Count of monolayers frequency in the 2.2M bilayers set, grouped depending on the position of the vdWHs bandgap in the visible spectrum. Here, the energy ranges as labeled as follows: infrared (IR) $\leq 1.65$ eV $\leq$ red (R) $\leq 1.99$ eV $\leq$ orange (O) $\leq 2.10$ eV $\leq$ yellow (Y) $\leq 2.17$ eV $\leq$ green (G) $\leq 2.50$ eV $\leq$ blue (B) $\leq 2.75$ eV $\leq$ violet (V) $\leq 3.26$ eV $\leq$ ultraviolet (UV).

| IR | | R | | O | | Y | |
|---|---|---|---|---|---|---|---|
| Monolayer | Count | Monolayer | Count | Monolayer | Count | Monolayer | Count |
| $F_6Li_2O_3Ta_2$ | 2003 | $B_2O_6U$ | 713 | $Cl_4P$ | 343 | $ClP_4$ | 197 |
| $C_6Li_2O_{18}V_3$ | 1863 | $F_3Zr$ | 694 | $F_3Ti$ | 298 | $F_3Ti$ | 179 |
| $LiO_{12}Te_2V_3$ | 1860 | $O_8Pb_3V_2$ | 683 | $H_2Mg_3O_{12}Si_4$ | 227 | $HfI_3$ | 154 |
| $Li_4O_{12}Te_3V$ | 1843 | $Ni_2O_8Te_3$ | 673 | $Nb_2O_3$ | 223 | $Br_3MoTe_6$ | 152 |
| $C_3Li_2O_9V$ | 1805 | $Ge_2Se_5Tl_2$ | 640 | $Br_3MoTe_6$ | 216 | $Cl_8Ge_3$ | 143 |
| **G** | | **B** | | **V** | | **UV** | |
| Monolayer | Count | Monolayer | Count | Monolayer | Count | Monolayer | Count |
| $Cl_3Hf$ | 737 | $Cl_3Zr$ | 596 | $FaI_3Si$ | 997 | $CuNb_2O_8Zn_2$ | 2127 |
| $S_3Y_2$ | 696 | $I_3Ti$ | 531 | $Br_3In$ | 996 | $F_7RbSb_2$ | 2126 |
| $O_2P$ | 674 | $Br_2Mo_2S$ | 516 | $F_3Hf$ | 995 | $B_2CuO_6Pb_2$ | 2120 |
| $HfI_4$ | 655 | $Br_3Sn$ | 508 | $Cl_8CuGa_2$ | 969 | $F_5VZn$ | 2111 |
| $HfI_3$ | 654 | $Br_2Mo$ | 489 | $Cl_2Mo$ | 950 | $Cu_3F_8Li_2$ | 2098 |



**Figure 5.** a) Bilayer GGA-PBE and HSE06 bandgap correlation and b) their relative energy distribution. Their calculated Pearson and Spearman correlation coefficients are 0.68 and 0.60, respectively, indicating a significant (linear) correlation.

present analysis within the high-symmetry/low-energy/bilayer configurations of the vdWHs.[39]

Furthermore, the potential correlations of the HSE06 bandgap with the IE and elastic constant along the z-axis ($C_{33}$), calculated in our previous work were analyzed.[32] These two quantities represent macroscopic features that originate from the charge redistribution upon vdWHs formation, and therefore they may, in principle, have an influence on the band alignment that determines the bandgap. However, the Pearson and Spearman coefficients between the HSE06, and the IE and $C_{33}$ were calculated to be $\approx 0.06$ and 0.03, respectively, indicating a negligible correlation between these quantities, again confirming that the complex nature of the bandgap goes beyond intuitive considerations.

## 3. Conclusions

We have used DFT calculations and AL to generate a database of $\approx 2.2$ million novel vdWHs bandgaps, potentially identifying those with significant potential for technological and scientific utility. No correlations between the BG and other fundamental properties of vdWHs such as IE and $C_{33}$ were found, highlighting the BG calculation as a fundamental problem rather than an emerging property correlated to macroscopic observables. This work did, however, find a relatively strong correlation between the PBE and HSE06 calculated BG. Furthermore, the potential of active ML to very substantially accelerate convergence of ML model prediction accuracies using modest training set sizes was demonstrated, demonstrating the capabilities of the combined ML + DFT approaches used in materials discovery.

## 4. Experimental Section

*DFT Calculations*: To calculate the energy of the structures by DFT, we employed a projector-augmented waves (PAW) approach as implemented in VASP, within both the GGA (Perdew et al. [PBE]) and HSE06.[42–44] In HSE06, the van der Waals correlation correction was applied using the method of Grimme, with Becke–Jonson damping (DFT-D3).[40,41,45] A $(3 \times 3 \times 1)$ point k-mesh, where $x$ and $y$ are in the plane of the 2D layers and $z$ represents the orthogonal stacking axis for the monolayers, and a basis set energy cut-off of 700 eV were used for all geometry optimization calculations. A larger grid of $(9 \times 9 \times 1)$ was used for non-self-consistent bandgap calculations. The energy minimization tolerance was $10^{-6}$ eV, and the force tolerance was $10^{-2}$ eV $\text{Å}^2$. A vacuum of $\approx 15$ Å along the $z$-axis was chosen to avoid interactions with replicas in the periodic boundary conditions.

Structural information was obtained for 6,138 monolayers from an online database (https://2dmatpedia.org/) and 2,132 semiconducting structures were selected to build the 2 270 142 (2.2M) bilayer database. A subset of bilayers was used for the DFT calculations, selected by applying two constraints: the lattice mismatch, $L_m < 2\%$; and the number of atoms in the cell, $N_a < 200$. This ensured both a reliable convergence of DFT calculations and a reasonable computational time.

Only one twist angle between the two monolayers was considered, corresponding to the lowest energy configuration. However, some additional considerations were also required. In both homo- and heterostructures, the weak interlayer forces allowed the angles between the monolayers along the $x$–$y$ plane to adopt different values, with only a small change in the interlayer energy (IE; $\approx 30$ meV). However, as a consequence of the angle-dependent symmetry of the resulting bilayer and the resulting formation of dipoles, the band structure may be significantly affected by the twist angle (changes reported in the literature are between 5% and 15%).[35–37,46] Although this may be useful for engineering the bandgaps of structures, it adds additional complication to the bandgap analysis.[38,47,48] We found a negligible change in the bandgap at high symmetry twist angles for selected bilayers, supporting the accuracy of our calculations within low energy configurations.

*Bayesian Neural Networks*: We used a Bayesian neural network (BNN) ML algorithm to predict the bandgaps of a large number of heterostructures. From the Bayesian point of view, regressions were formulated using probability distributions rather than point estimates.[49] The target property, or response, was not estimated as a single value, but was assumed to be drawn from a probability distribution using a dropout approach.[50–52] Therefore, our BNN also predicted the confidence interval for each value, indicative of the quality of the prediction for each individual heterostructure.[53] Here, a BNN with two hidden layers composed of 32 neurons each was used, where the dropout probability was 0.1. The dropout regularized the network and avoided overfitting by creating a distribution over the calculated response. This was averaged over 600 trial networks giving the mean response value and the associated standard deviation. A detailed description of the BNN used can be found in previous studies.[32]

*Active Learning*: In evaluating the exact exchange energy density, HSE06 calculations have a large computational cost and large memory requirements compared with GGA. The central processing unit (CPU) time ratio of a single-point calculation carried out within HSE06 approximations was $\approx 5$ and 7 times larger than for GGA-PBE calculations, consistent with the literature.[54] Thus, the use of HSE06 calculations to compute the bandgap for selected representative bilayers in the training set was very computationally demanding. Here, an active learning (AL) approach was adopted to restrict the number of HSE06 calculations required to train the BNN.[55,56] The main advantage of AL was that the data used for BNN training can be chosen selectively, leading to a better performance, while requiring substantially less data than traditional static learning methods. It was important to ensure that the training dataset was representative of the true distribution of the data.

The data were divided into a very small, labeled dataset that included up to a few hundred bilayers (the seed $X_L$), and a large unlabeled dataset that included up to $\approx 2.2$ million bilayers ($X_U$). Typically, there was an arbitrary partitioning between the labeled and unlabeled data. After splitting the data, the seed can be used to train the model. Once the model had been trained $N$ times using a different training-test split by k-means clustering, it can be used to predict the response for the unlabeled data ($Y_i = f_i(X_U)$). Each element of the unlabeled data will therefore have $N$ different predicted values. By calculating the mean and standard deviation of the target property (response) for each unlabeled data over the $N$ runs, the worse performing structure was selected, its properties calculated using HSE06, and added to the seed for the next BNN training iteration. With this new labeled dataset, the learner can be retrained, iterating the process until the required accuracy was achieved, evaluated by parameters such as $R^2$, mean square error (MSE), mean absolute error (MAE), and mean absolute percentage error (MAPE). The measures of dispersion were preferred as they are less dependent on the number of parameters in the model and the number of materials in the data set.[57]

## Supporting Information

Supporting Information is available from the Wiley Online Library or from the author.

## Acknowledgements

## Conflict of Interest

The authors declare no conflict of interest.

## Data Availability Statement

The data that support the findings of this study are openly available in Opal at https://opal.latrobe.edu.au. The complete bandgap data are available at https://data.research.uts.edu.au/publication/fae85210bd1e11eba4d3adb3-c726e5fe/. Custom Python codes for data preprocessing and Bayesian

**ADVANCED
SCIENCE NEWS**

www.advancedsciencenews.com

**ADVANCED
INTELLIGENT
SYSTEMS**

www.advintellsyst.com

neural network training and data extrapolation are available at https://github.com/fronzi/projBNN.

[1] E. V. Castro, *Phys. Rev. Lett.* **2007**, *99*, 216802.

[2] J. Feng, X. Qian, C. W. Huang, J. Li, *Nat. Photonics* **2012**, *6*, 866.

[3] A. H. C. Neto, F. Guinea, N. M. R. Peres, K. S. Novoselovv, A. K. Geim, *Rev. Mod. Phys.* **2009**, *81*, 109.

[4] A. M. Smith, S. Nie, *Acc. Chem. Res.* **2010**, *43*, 190.

[5] A. Chaves, J. G. Azadani, H. Alsalman, D. R. da Costa, R. Frisenda, A. J. Chaves, S. H. Song, Y. D. Kim, D. He, J. Zhou, A. Castellanos-Gomez, F. M. Peeters, Z. Liu, C. L. Hinkle, S.-H. Oh, P. D. Ye, S. J. Koester, Y. H. Lee, P. Avouris, X. Wang, T. Low, *npj 2D Mater. Appl.* **2020**, *4*, 29.

[6] A. Kudo, *Int. J. Hydrogen Energy* **2007**, *32*, 2673.

[7] R. Mas-Ballesté, C. Gómez-Navarro, J. Gómez-Herrero, F. Zamora, *Nanoscale* **2011**, *3*, 20.

[8] J. Su, L. Guo, N. Bao, C. A. Grimes, *Nano. Lett.* **2011**, *11*, 1928.

[9] H. Tao, Q. Fan, T. Ma, S. Liu, H. Gysling, J. Texter, F. Guo, Z. Sun, *Prog. Mater. Sci.* **2020**, *111*, 100637.

[10] W. Lei, S. Zhang, G. Heymann, X. Tang, J. Wen, X. Zheng, G. Hu, X. Ming, *J. Mater. Chem. C* **2019**, *7*, 2096.

[11] G. Rao, X. Wang, Y. Wang, P. Wangyang, C. Yan, J. Chu, L. Xue, C. Gong, J. Huang, J. Xiong, Y. Li, *InfoMat* **2019**, *1*, 272.

[12] B. Wang, S. P. Zhong, Z. B. Zhang, Z. Q. Zheng, Y. P. Zhang, H. Zhang, *Appl. Mater. Today* **2019**, *15*, 115.

[13] F. Wang, Y. Zhang, Y. Gao, P. Luo, J. Su, W. Han, K. Liu, H. Li, T. Zhai, *Small* **2019**, *15*, 1901347.

[14] L. Zhang, B. Wang, Y. Zhou, C. Wang, X. Chen, H. Zhang, *Adv. Opt. Mater.* **2020**, *8*, 2000045.

[15] Z. Zhang, Y. Guo, J. Robertson, *Appl. Phys. Lett.* **2020**, *116*, 251602.

[16] Z. Zhang, B. Huang, Q. Qian, Z. Gao, X. Tang, B. Li, *APL Mater.* **2020**, *8*, 041114.

[17] X. Zong, H. Hu, G. Ouyang, J. Wang, R. Shi, L. Zhang, Q. Zeng, C. Zhu, S. Chen, C. Cheng, B. Wang, H. Zhang, Z. Liu, W. Huang, T. Wang, L. Wang, X. Chen, *Light Sci. Appl.* **2020**, *9*, 114.

[18] G. Onida, L. Reining, A. Rubio, *Rev. Mod. Phys.* **2002**, *74*, 601.

[19] A. Dittmer, R. Izsák, F. Neese, D. Maganas, *Inorg. Chem.* **2019**, *58*, 9303.

[20] H. C. Kamban, T. G. Pedersen, *npj Sci. Rep.* **2020**, *10*.

[21] W. Kohn, L. J. Sham, *Phys. Rev.* **1965**, *140*, A1133.

[22] I. N. Yakovkin, P. A. Dowben, *Surf. Rev. Lett.* **2007**, *14*, 481.

[23] X. Liu, L. Li, Q. Li, Y. Li, F. Lu, *Mater. Sci. Semicond. Proc.* **2013**, *16*, 1369.

[24] D. Ma, Y. Chai, V. Wang, E. Li, W. Shi, *Mater. Des.* **2015**, *87*, 877.

[25] J. Yang, Q. Fan, *IOP Conf. Ser.: Mater. Sci. Eng.* **2017**, *167*, 012010.

[26] M. Marsman, J. Paier, A. Stroppa, G. Kresse, *J. Phys.: Condens. Matt.* **2008**, *20*, 064201.

[27] S. Smiga, L. A. Constantin, *J. Phys. Chem. A* **2020**, *124*, 5606.

[28] J. Heyd, G. E. Scuseria, M. Ernzerhof, *J. Chem. Phys.* **2003**, *118*, 8207.

[29] J. Heyd, G. E. Scuseria, M. Ernzerhof, *J. Chem. Phys.* **2006**, *124*, 219906.

[30] F. Furche, J. P. Perdew, *J. Chem. Phys.* **2006**, *124*, 044103.

[31] G.-X. Zhang, A. M. Reilly, A. Tkatchenko, M. Scheffler, *New J. Phys.* **2018**, *20*, 063020.

[32] M. Fronzi, S. A. Tawfik, M. A. Ghazaleh, O. Isayev, D. A. Winkler, J. Shapter, M. J. Ford, *Adv. Theor. Simul.* **2020**, *3*, 2000029.

[33] K. T. Butler, D. W. Davies, H. Cartwright, O. Isayev, A. Walsh, *Nature* **2018**, *559*, 547.

[34] S. A. Tawfik, O. Isayev, C. Stampfl, J. Shapter, D. A. Winkler, M. J. Ford, *Adv. Theor. Simul.* **2019**, *2*, 1800128.

[35] K. P. Nuckolls, M. Oh, D. Wong, B. Lian, K. Watanabe, T. Taniguchi, B. A. Bernevig, A. Yazdani, arXiv:2007.03810 **2020**.

[36] H. Polshyn, M. Yankowitz, S. Chen, Y. Z. K. Watanabe, T. Taniguchi, C. R. Dean, A. Young, *Nat. Phys.* **2019**, *15*, 1011.

[37] S. L. Tomarken, Y. Cao, A. Demir, K. Watanabe, T. Taniguchi, P. Jarillo-Herrero, R. C. Ashoori, *Phys. Rev. Lett.* **2019**, *123*, 046601.

[38] N. Liu, J. Zhang, S. Zhou, J. Zhao, *J. Mater. Chem. C* **2020**, *8*, 6264.

[39] M. Huang *Optoelectronics-Devices and Applications* (Ed: P. Predeep), InTech, Riijeka, Croatia **2011**, Ch. 6.

[40] S. Grimme, J. Antony, S. Ehrlich, S. Krieg, *J. Chem. Phys.* **2010**, *132*, 154104.

[41] S. Grimme, S. Ehrlich, L. Goerigk, *J. Comput. Chem.* **2011**, *32*, 1456.

[42] G. Kresse, J. Hafner, *J. Phys.: Condens. Matt.* **1994**, *6*, 8245.

[43] G. Kresse, D. Joubert, *Phys. Rev. B* **1999**, *59*, 1758.

[44] J. P. Perdew, K. Burke, M. Ernzerhof, *Phys. Rev. Lett.* **1996**, *77*, 3865.

[45] A. M. Ukpong, *Comput. Condens. Matter* **2015**, *2*, 1.

[46] N. Lu, H. Guo, Z. Zhuo, L. Wang, X. Wu, X. C. Zeng, *Nanoscale* **2017**, *9*, 19131.

[47] H. Y. Lee, M. M. Al Ezzi, N. Raghuvanshi, J. Y. Chung, K. Watanabe, T. Taniguchi, S. Garaj, S. Adam, S. Gradečak, *Nano Lett.* **2021**, *21*, 2832.

[48] J.-B. Liu, P.-J. Li, Y.-F. Chen, Z.-G. Wang, F. Qi, J.-R. He, B.-J. Zheng, J.-H. Zhou, W.-L. Zhang, L. Gu, Y.-R. Li, *npj Sci. Rep.* **2015**, *5*.

[49] F. R. Burden, D. A. Winkler, *J. Med. Chem.* **1999**, *42*, 3183.

[50] Y. Gal, Z. Ghahramani, arXiv: 1506.02142v6 **2015**.

[51] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, *J. Mach. Learn. Res.* **2014**, *15*, 1929.

[52] D. Tran, M. W. Dusenberry, M. van der Wilk, D. Hafner, arXiv: 1812.03973 **2018**.

[53] S. Bergmann, S. Stelzer, S. Strassburger, *J. Simul.* **2014**, *8*, 76.

[54] J. Heyd, G. E. Scuseria, *J Chem Phys* **2004**, *121*, 1187.

[55] B. Settles, Computer Sciences Technical Report 1648, University of Wisconsin–Madison, **2010**.

[56] B. Settles, *Synth. Lect. Artif. Intell. Mach. Learn.* **2012**, *6*, 1.

[57] D. L. J. Alexander, A. Tropsha, D. A. Winkler, *J. Chem. Inf. Model.* **2015**, *55*, 1316.

[58] R. Dong, A. Jacob, S. Bourdais, S. Sanvito, *npj 2D Mater. Appl.* **2021**, *5*.

[59] S. Wang, Y. Hu, Y. Wei, W. Li, N. T. Kaner, Y. Jiang, J. Yang, X. Li, *Phys. E* **2021**, *130*, 114674.

[60] X. Yang, X. Qin, J. Luo, N. Abbas, J. Tang, Y. Li, K. Gu, *RSC Adv.* **2020**, *10*, 2615.

Author/s:
Fronzi, M; Isayev, O; Winkler, DA; Shapter, JG; Ellis, A; Sherrell, PC; Shepelin, NA; Corletto, A; Ford, MJ

Title:
Active Learning in Bayesian Neural Networks for Bandgap Predictions of Novel Van der Waals Heterostructures

Date:
2021-08-22

Citation:
Fronzi, M., Isayev, O., Winkler, D. A., Shapter, J. G., Ellis, A., Sherrell, P. C., Shepelin, N. A., Corletto, A. & Ford, M. J. (2021). Active Learning in Bayesian Neural Networks for Bandgap Predictions of Novel Van der Waals Heterostructures. ADVANCED INTELLIGENT SYSTEMS, https://doi.org/10.1002/aisy.202100080.

Persistent Link:
http://hdl.handle.net/11343/283331

File Description:
Published version
License:
CC BY