

Universal Machine Learning Kohn-Sham Hamiltonian for Materials

Yang Zhong^{1,2}, Hongyu Yu^{1,2}, Jihui Yang^{1,2}, Xingyu Guo^{1,2}, Hongjun Xiang^{1,2,*}, and Xingao Gong^{1,2}

¹*Key Laboratory of Computational Physical Sciences (Ministry of Education), Institute of Computational Physical Sciences, State Key Laboratory of Surface Physics, and Department of Physics, Fudan University, Shanghai, 200433, China*

²*Shanghai Qi Zhi Institute, Shanghai, 200030, China*

Email: hxiang@fudan.edu.cn

Abstract

While density functional theory (DFT) serves as a prevalent computational approach in electronic structure calculations, its computational demands and scalability limitations persist. Recently, leveraging neural networks to parameterize the Kohn-Sham DFT Hamiltonian has emerged as a promising avenue for accelerating electronic structure computations. Despite advancements, challenges such as the necessity for computing extensive DFT training data to explore each new system and the complexity of establishing accurate ML models for multi-elemental materials still exist. Addressing these hurdles, this study introduces a universal electronic Hamiltonian model trained on Hamiltonian matrices obtained from first-principles DFT calculations of nearly all crystal structures on the Materials Project. We demonstrate its generality in predicting electronic structures across the whole periodic table, including complex multi-elemental systems, solid-state electrolytes, Moiré twisted bilayer heterostructure, and metal-organic frameworks (MOFs). Moreover, we utilize the universal model to conduct high-throughput calculations of electronic structures for crystals in GeNOME datasets, identifying 3,940 crystals with direct band gaps and 5,109 crystals with flat bands. By offering a reliable efficient framework for computing electronic properties, this universal Hamiltonian model lays the groundwork for advancements in diverse fields, such as easily providing a huge data set of electronic structures and also making the materials design across the whole periodic table possible.

Introduction

The electronic structure¹⁻⁵ of materials is crucial in understanding and predicting a wide range of physical properties, including electrical conductivity, optical behavior, mechanical strength, chemical reactivity, and magnetic characteristics. Electronic structure calculations provide insights into the electronic band structures, bonding, and reactivity, enabling the design of new materials and the study of chemical reactions.

Among diverse quantum mechanics approaches, density functional theory (DFT)⁶⁻⁹ has become a widely used computational method in electronic structure calculations since DFT has drastically reduced the computational cost by employing the electron density instead of the many-body wave function as the fundamental variable. The price to pay for using the electron density as the fundamental variable is that the Kohn-Sham DFT Hamiltonian is no longer a simple explicit function of the atomic structure, but should be obtained by solving a self-consistent Kohn-Sham equation. However, the computational cost of DFT self-consistent field (SCF) cycles remains expensive for large systems.

In recent years, the use of neural networks to parameterize the DFT Hamiltonian has emerged as an important and effective method to accelerate electronic structure calculations¹⁰⁻²². The ML Hamiltonian models offer a significant advantage by providing a direct mapping from the structure to the self-consistent Hamiltonian matrix, eliminating the need for time-consuming self-consistent iterations typically required in Kohn-Sham DFT. Initially, Hegde and Bowen proposed a kernel ridge regression (KRR) model and successfully applied it to fit the empirical Hamiltonian matrix of the cubic Cu crystal¹⁰. Subsequently, Unke¹² proposed the PhiSNet model, and Schütt¹⁴ introduced the SchNOrb model, both of which demonstrated remarkable performance in accurately fitting the Hamiltonian matrices of various small organic molecules such as water, ethanol, and uracil. In addition, some other researchers have also made important contributions. For instance, Nigam¹⁸ used the Gaussian Process Regression (GPR) model to fit the Hamiltonian matrices of water and benzene molecules, while Zhang²⁰ successfully fitted the Hamiltonian matrix of aluminum using the atomic cluster expansion (ACE) model. Xu and coworkers proposed the graph neural network (GNN)-based DeepH^{17, 21} model, which was used to predict the Hamiltonian matrices of crystals like graphene and MoS₂.

However, despite these advancements, there still exist several challenges in utilizing machine learning for electronic Hamiltonian prediction. Firstly, exploring new systems necessitates retraining a completely new model specifically tailored for that system, which currently lacks automation and can prove time-consuming and computationally expensive, thereby constraining the practical applicability of this approach. Secondly, many practical materials, such as high-entropy alloys²³⁻²⁵ and ceramics^{26, 27}, are composed of a large number of elements, posing challenges in establishing accurate machine-learning models for their electronic Hamiltonians due to the requirement of a substantial amount of DFT training data. Current research endeavors have predominantly concentrated on systems comprising no more than three elements, and surmounting this obstacle to encompass more intricate systems remains a significant challenge in the field. Overcoming these challenges is crucial to enable the broader application of machine learning in predicting electronic Hamiltonians for diverse and multi-elemental materials. Recently, several groups reported developments of universal machine learning interatomic potentials (MLIPs)²⁸⁻³³, which can handle almost all elements across the whole periodic table. Naturally, one might wonder whether it is possible to develop a universal machine-learning model for electronic Hamiltonians. Such a model would enable efficient and accurate calculations of electronic structures

for a wide range of materials and large-scale systems. Recently, we have proposed a transferable Hamiltonian graph neural network (HamGNN)²² that enables the prediction of Hamiltonian matrices for various structures within a chemical space comprising several specific elements, such as SiO₂ isomers and Bi_xSe_y compounds with varying stoichiometric ratios. Compared to the above mentioned models, the excellent transferability of the equivariant GNN-based HamGNN framework makes it the most feasible candidate for constructing a universal Hamiltonian model across the entire periodic table. Unlike the construction of universal MLIPs, achieving a universal Hamiltonian model is not as straightforward and cannot simply rely on training with large-scale datasets due to the high dimensionality and inherent complexity of the Hamiltonian matrix. The training process also plays a crucial role in achieving the universality of the Hamiltonian model.

In this work, we propose a universal Kohn-Sham Hamiltonian model in the sense that this single model is applicable to all elements of the entire periodic table and all structures of a given chemical composition. By employing the HamGNN model, we have developed a universal Hamiltonian model via a 'two-step training procedure' utilizing the Hamiltonian matrices of 5,5000 structures obtained from the Materials Project^{34, 35} as our training dataset. The universality of our Hamiltonian model is demonstrated by its successful prediction of the electronic structures for various bulk or low-dimensional materials with different combinations of chemical elements in the periodic table. The universal model successfully captures not only common systems but also those containing uncommon or rare transition metal elements. Furthermore, it proves high accuracy even in complex multi-element systems comprising more than five elements. By providing a reliable framework for understanding electronic properties across the periodic table, this research paves the way for advancements in material design, catalysis, electronics, and other fields that heavily rely on efficient predictions of electronic structures.

Results

The framework of the universal electronic Hamiltonian model

Constructing an ML model for the electronic Hamiltonian is more complicated than constructing machine learning interatomic potentials (MLIPs)³⁶⁻⁴¹. MLIPs provide a mapping from the two degrees of freedom, atomic types $\{Z_i\}$ and atomic positions $\{\mathbf{r}_i\}$, to the scalar potential energy E . In addition to the two degrees of freedom, the mapping from a crystal structure to the electronic Hamiltonian matrix also necessitates handling the supplementary degree of freedom arising from distinct atomic orbital bases $\{\phi_{ia}\}$ associated with each atom. As the number of elements in the system increases, the relevant degrees of freedom and interactions in the electronic Hamiltonian matrix also increase dramatically, leading to a significant increase in complexity when fitting the electronic Hamiltonian matrix. Therefore, a universal electronic Hamiltonian model across the periodic table requires much more network capacity than the MLIPs to capture all degrees of freedom accurately. For large molecules or crystals, especially in cases where the basis set is large and the system is complex, the dimension of the

electronic Hamiltonian matrix can be extremely large. This increases the difficulty of model training, as it requires handling large-scale matrix and storing a significant amount of parameters. In addition, different sub-blocks of a Hamiltonian matrix are subject to different equivariant constraints under a rotation operation, so the Hamiltonian matrix predicted by the model must also adhere to these constraints.

We have trained a universal Hamiltonian model following the process shown in Figure 1. To develop a universal Hamiltonian model for the whole periodic table, we utilized one of the world's largest open databases for DFT-relaxed crystal structures, namely the Materials Project^{34, 35}. We used OpenMX^{42, 43}, a DFT software package based on norm-conserving pseudopotentials and pseudo-atomic localized basis functions, to calculate the Hamiltonian matrices of approximately 55,000 structures on the Materials Project. Among them, approximately 44,000 structures' Hamiltonian matrices were used as a training set, while validation and test sets were created using the Hamiltonian matrices of around 5,500 structures each. Then, we use these datasets to train a universal HamGNN²² model. HamGNN is a deep learning model based on equivariant graph neural networks, which can automatically learn the features of each element on the entire periodic table without any prior physical or chemical properties of elements. The architecture of the HamGNN model and the training process are shown in Figure 1(b). We will briefly introduce the principle of HamGNN for predicting the Hamiltonian matrix in the following paragraph, and for more network details see Ref. 22.

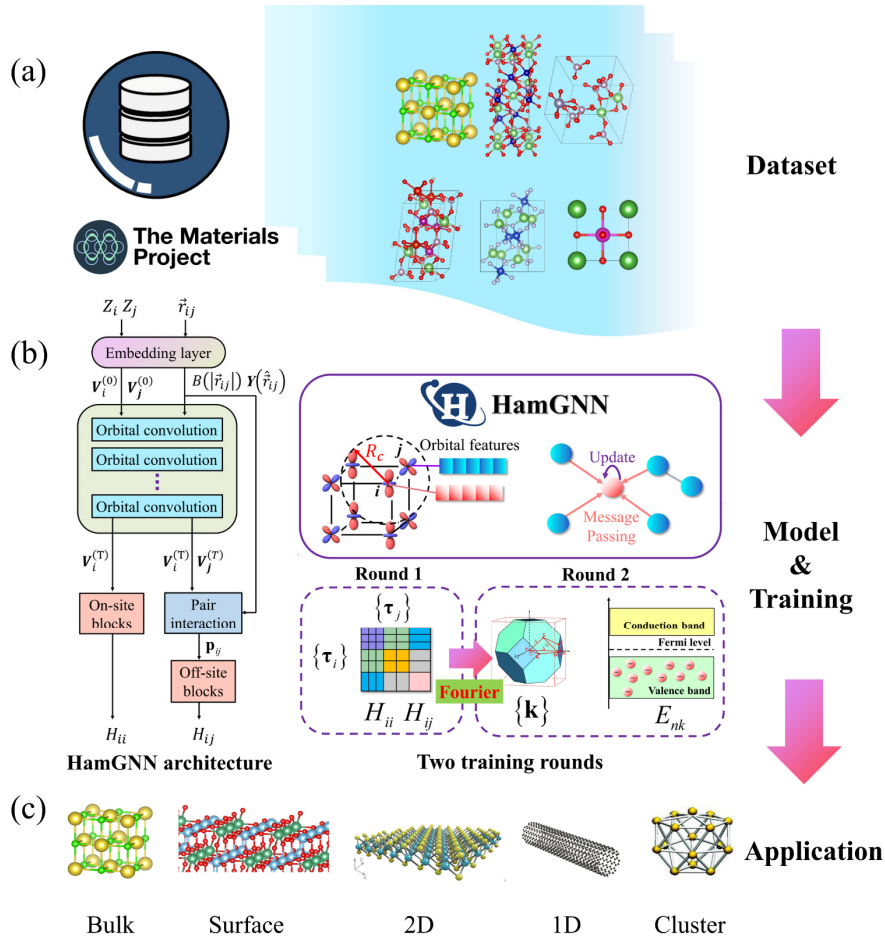


Figure 1. The framework of a universal Hamiltonian model based on HamGNN. (a) Training dataset preparation. The training dataset is generated by calculating the real-space Hamiltonian matrices of crystal structures available on the Materials Project using an ab initio tight-binding software based on numerical atomic orbitals. (b) Model architecture and the training procedure. This dataset is utilized for training the HamGNN model, a deep learning approach that employs equivariant graph neural networks to predict the Hamiltonian matrix. The model can automatically learn the intrinsic features of each element on the periodic table solely based on their atomic numbers, without relying on any prior physical or chemical properties. In HamGNN, the orbital features of the central atom are updated by considering interactions between the neighbor atoms within a cutoff radius of R_c . For atomic pairs beyond this cutoff radius, multi-layer message passing is employed to exchange orbital features. To achieve universality, HamGNN requires two rounds of training. In the first round, the loss function for network training solely considers the error in the real-space Hamiltonian matrix. After this initial round of training, this model can accurately predict the real-space Hamiltonian matrices with high precision. In the second round of training, the real-space Hamiltonian matrices are transformed into the reciprocal Hamiltonian matrices at randomly selected \mathbf{k} points in the Brillouin zone, and the errors of the orbital energies near the Fermi level obtained by diagonalizing the reciprocal Hamiltonian matrices are incorporated into the total loss function. (c) Applications of the universal HamGNN model. After two rounds of training, the universality of the HamGNN model has significantly improved, enabling accurate prediction of electronic structures in crystals with arbitrary periodic boundary conditions and any components.

Since the irreducible representations with rotation orders $l = 0, 1, 2, \dots$ of the $O(3)$ group possess the same rotational equivariance and parity symmetry as the atomic orbitals s, p, d, \dots , HamGNN uses a direct sum of equivariant atomic features \mathbf{V}_l with different rotation orders l to characterize each atom: $\mathbf{V} = \mathbf{V}_0 \oplus \mathbf{V}_1 \oplus \dots \oplus \mathbf{V}_{l_{max}}$. This feature tensor satisfies the rotational equivariance under the rotation operation \hat{R} : $\mathbf{V}_l(\hat{R} \cdot (\vec{r}_1, \dots, \vec{r}_N)) = \mathbf{D}_l(\hat{R})\mathbf{V}_l(\vec{r}_1, \dots, \vec{r}_N)$, where \mathbf{D}_l ($l < l_{max}$) is a Wigner D matrix^{44, 45} of order l . HamGNN updates the equivariant atomic features through an equivariant message-passing function in the orbital convolution layer. After T orbital convolution layers, the atomic features are transformed into on-site Hamiltonian matrices by an “on-site layer”. HamGNN merges the features of atom pairs ij into the edge features \mathbf{P}_{ij} in the “Pair interaction layer”. The edge features \mathbf{P}_{ij} is later transformed into an off-site Hamiltonian matrix through an “off-site layer”. Each subblock of the Hamiltonian matrix can be decomposed into a set of $O(3)$ equivariant irreducible spherical tensors (ISTs) according to the following equation⁴⁴⁻⁴⁷:

$$l_i \otimes l_j = |l_i - l_j| \oplus |l_i - l_j| + 1 \oplus \dots \oplus l_i + l_j \quad (1)$$

The on-site and off-site layers output each sub-block of the Hamiltonian matrix by the following equation

$$H_{l_i m_i, l_j m_j} = \sum_{l=|l_i-l_j|}^{l_i+l_j} \sum_{m=-l}^l C_{m_i, m_j, m}^{l_i, l_j, l} T_m^l \quad (2)$$

where T_m^l is an equivariant IST with rotation order l in $V_i^{(T)}$ or P_{ij} , $C_{m_i, m_j, m}^{l_i, l_j, l}$ is the Clebsch-Gordan coefficient.

To achieve a universal model, HamGNN needs to undergo two rounds of training. During the first round of training, only the error of the real-space Hamiltonian matrix is considered as the loss function for the network's training. After the first round of training, the model is capable of reasonably predicting the real-space Hamiltonian matrix, which lays a good foundation for obtaining accurate band structures in subsequent tasks. To improve the accuracy of the model in predicting the eigenvalues of Bloch states, which constitutes the primary objective of our Hamiltonian model, we further applied fine-tuning techniques for a second round of training. Fine-tuning in neural networks has shown a crucial role in training recently developed large language models to optimize their adaptation to specific tasks or domains^{48, 49}. The results of our tests indicate that the two training rounds are necessary to obtain a truly universal Hamiltonian model. During each training step in the second round, Fourier transformations are performed on the predicted and target real-space Hamiltonian matrices at randomly selected \mathbf{k} points in the Brillouin zone. The orbital energies ε_{nk} are obtained by diagonalizing the reciprocal Hamiltonian matrices, and the error of the orbital energies near the Fermi level is incorporated into the loss function as follows:

$$L = \|\tilde{H} - H\| + \frac{\lambda}{N_{orb} \times N_k} \sum_{k=1}^{N_k} \sum_{n=1}^{N_{orb}} \|\tilde{\varepsilon}_{nk} - \varepsilon_{nk}\| \quad (3)$$

where the variables marked with a tilde refer to the corresponding predictions and λ denotes the loss weight of the orbital energy error. N_{orb} is the number of orbits selected near the Fermi level, N_k is the number of the random \mathbf{k} points generated in each training step.

Figure 2(a) displays the count of each chemical element in the training set for the Hamiltonian. The metallic elements of Group IA and IIA, as well as the non-metallic elements of Group IVA, VA, VIA, and VIIA, have the highest proportion in the training set. Except for some less common transition metal elements that are not included in the training set, this dataset includes all elements supported by OpenMX's PBE pseudopotential library. After the first round of training, the HamGNN model achieved a mean absolute error (MAE) of only 5.4 meV for the real-space Hamiltonian matrix on the test set. The accuracy of the model for the real-space Hamiltonian matrix is shown in Figure 2(b). In the second round of training, we incorporate the error of orbital energies at five randomly selected \mathbf{k} -points in the Brillouin zone into the loss function to restart the training, with a λ value set at 0.01. After the second round of model fine-tuning, the generalization ability of the HamGNN model has significantly improved. The model's high accuracy on the test set is evident from its predictions of energy bands and Fermi surfaces for several crystal structures in the test set (see Supplementary Discussion 1). By addressing potential overfitting issues, this model

demonstrates better adaptability to different datasets and real-world scenarios, showcasing enhanced universality and stability. In the subsequent discussions, we will evaluate HamGNN's prediction accuracy for more complex crystals across the entire periodic table.

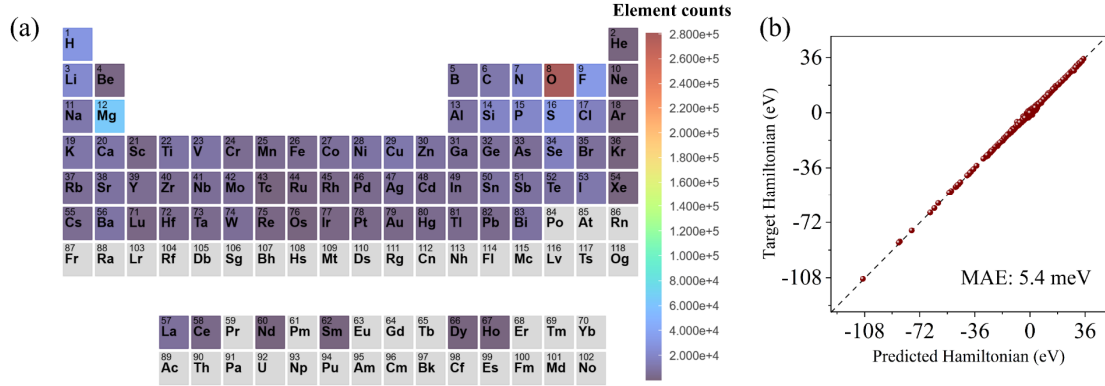


Figure 2. (a) The element distribution of the training dataset. (b) The comparison of the Hamiltonian calculated by HamGNN and OpenMX for the test dataset.

Tests on multi-element materials

We are now applying our universal Hamiltonian model to systems that are not included in the Materials Project to verify the generality and accuracy of our model. Previous ML Hamiltonian models^{10, 15, 20} typically handle crystal structures composed of only 1 to 3 elements, and training an accurate model that can deal with complex crystal structures containing more elements is challenging. The training datasets commonly used by these models are built from structures perturbed using molecular dynamics. However, as the number of atomic species in the crystal increases, the degrees of freedom of the Hamiltonian matrix increase sharply. Consequently, the training samples generated by perturbing a seed structure cannot fully cover the entire configuration space of the crystal. In these cases, the perturbed structures often contain many similar and repetitive patterns, which can cause the Hamiltonian model to be trapped in local minima and fail to accurately fit the Hamiltonian of crystals with multiple elements and complex configurations. However, the universal Hamiltonian model can effectively address such concerns. Through extensive training on a comprehensive and diverse dataset, the universal Hamiltonian model develops a profound understanding of the intricate interactions among atoms in various configurations.

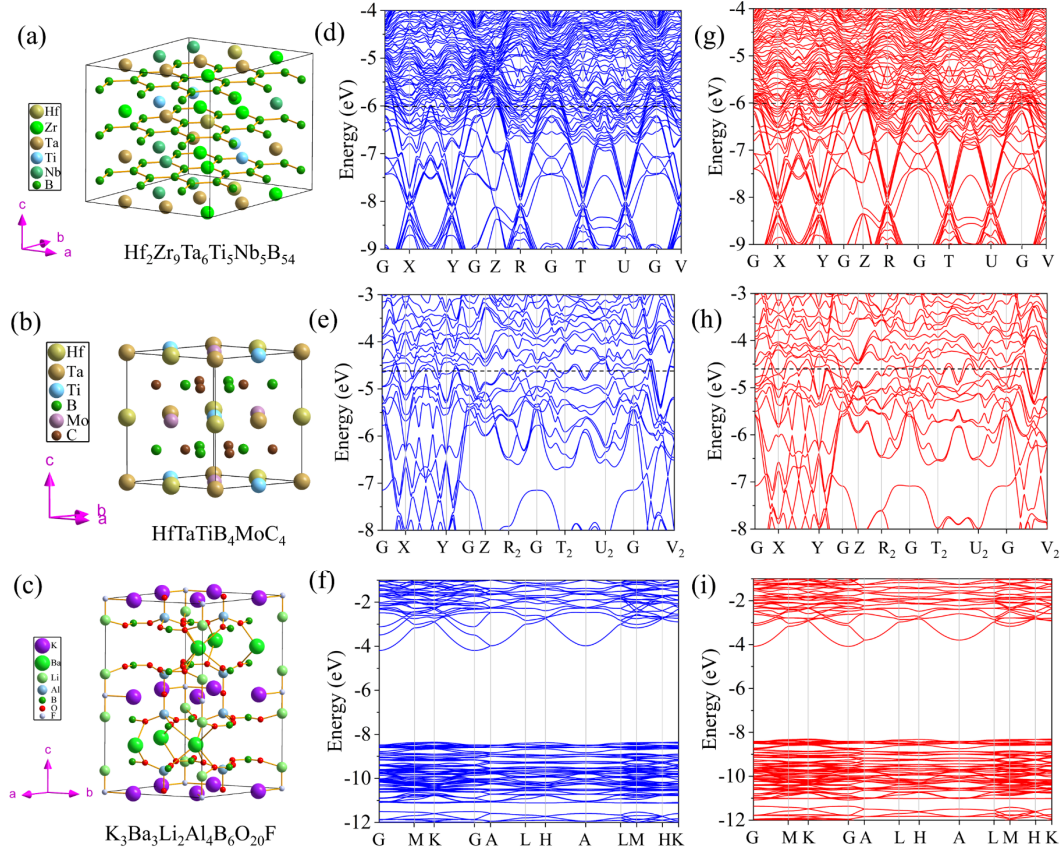


Figure 3. The crystal structures of (a) $\text{Hf}_2\text{Zr}_9\text{Ta}_6\text{Ti}_5\text{Nb}_5\text{B}_{54}$, (b) $\text{HfTaTiB}_4\text{MoC}_4$, and (c) $\text{K}_3\text{Ba}_3\text{Li}_2\text{Al}_4\text{B}_6\text{O}_{20}\text{F}$. The predicted energy bands of (d) $\text{Hf}_2\text{Zr}_9\text{Ta}_6\text{Ti}_5\text{Nb}_5\text{B}_{54}$, (e) $\text{HfTaTiB}_4\text{MoC}_4$, and (f) $\text{K}_3\text{Ba}_3\text{Li}_2\text{Al}_4\text{B}_6\text{O}_{20}\text{F}$. The DFT calculated energy bands of (g) $\text{Hf}_2\text{Zr}_9\text{Ta}_6\text{Ti}_5\text{Nb}_5\text{B}_{54}$, (h) $\text{HfTaTiB}_4\text{MoC}_4$, and (i) $\text{K}_3\text{Ba}_3\text{Li}_2\text{Al}_4\text{B}_6\text{O}_{20}\text{F}$.

The universal HamGNN's generalization performance and accuracy were evaluated by conducting tests on three different crystal structures: $\text{Hf}_2\text{Zr}_9\text{Ta}_6\text{Ti}_5\text{Nb}_5\text{B}_{54}$, $\text{HfTaTiB}_4\text{MoC}_4$, and $\text{K}_3\text{Ba}_3\text{Li}_2\text{Al}_4\text{B}_6\text{O}_{20}\text{F}$. These crystals were specifically chosen to represent a diverse range of compositions and structural complexities⁵⁰⁻⁵². The compound $\text{Hf}_2\text{Zr}_9\text{Ta}_6\text{Ti}_5\text{Nb}_5\text{B}_{54}$ exhibits a hexagonal ω -phase derived structure and belongs to the $P1$ space group of the triclinic crystal system. It consists of five distinct metal ions, which are inserted randomly into the gaps between the hexagonal monolayers of boron. $\text{HfTaTiB}_4\text{MoC}_4$ crystallizes in the triclinic $P1$ space group, with the metal elements inserted into the gaps between the hexagonal monolayers composed of boron and carbon. $\text{K}_3\text{Ba}_3\text{Li}_2\text{Al}_4\text{B}_6\text{O}_{20}\text{F}$ is a deep ultraviolet transparent nonlinear optical crystal, composed of seven elements and possessing a large bandgap⁵³.

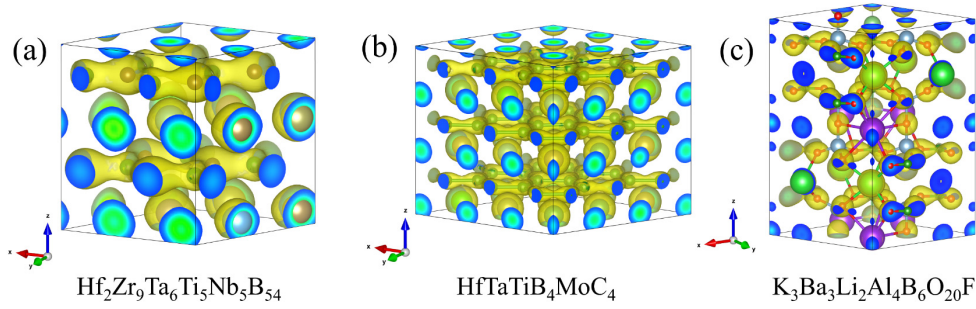


Figure 4. The predicted charge density of (a) $\text{Hf}_2\text{Zr}_9\text{Ta}_6\text{Ti}_5\text{Nb}_5\text{B}_{54}$, (b) $\text{HfTaTiB}_4\text{MoC}_4$, and (c) $\text{K}_3\text{Ba}_3\text{Li}_2\text{Al}_4\text{B}_6\text{O}_{20}\text{F}$.

To assess the generalization performance of HamGNN, a comparison was made between the predicted energy bands generated by the model and the calculated energy bands by OpenMX for each of these crystals, as shown in Figure 3. By examining Figure 3, it becomes evident that HamGNN demonstrates remarkable generalizability in predicting the energy bands for the three crystal structures. The charge densities of these three crystals were also obtained using the predicted Hamiltonian matrix, as shown in Figure 4. The predicted charge densities exhibit excellent agreement with the corresponding DFT-calculated charge densities (shown in Supplementary Figure. S6). This further confirms that HamGNN is capable of accurately capturing not only energy bands but also charge distribution within these materials. Furthermore, to test its versatility and applicability on practical materials, we extended our analysis to a sulfide solid electrolyte with a composition of $\text{Li}_{10}\text{Si}_{1.5}\text{P}_{1.5}\text{S}_{11.5}\text{C}_{10.5}$, a highly disordered system consisting of 7200 atoms. Our universal model shows that this disordered system is insulating with a band gap of 1.4 eV, suggesting that it indeed has excellent electric insulating properties and may serve as a promising solid electrolyte (see Supplementary Discussion 2 for details). These successful evaluations highlight the effectiveness and reliability of HamGNN as a universal model for predicting electronic structures across various complex crystal systems. Its ability to generalize well across different compositions and structural complexities makes it an invaluable tool in high-throughput electronic structure calculations for the periodic table.

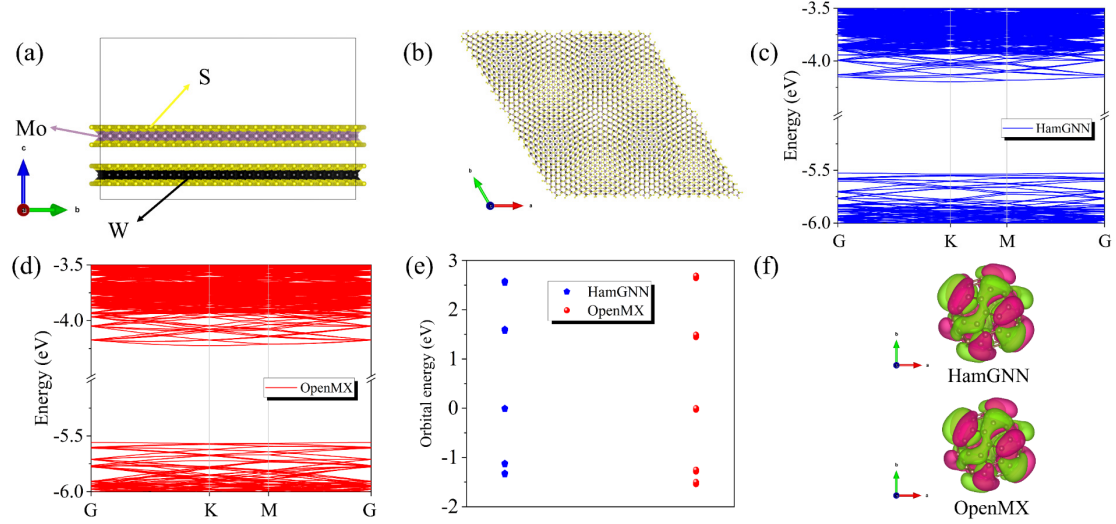


Figure 5. Prediction results of the universal HamGNN model for the bilayer MoS₂/WS₂ heterostructure with a twist angle of 3.5° and the C60 cluster. (a) Side view of the bilayer MoS₂/WS₂ heterostructure. (b) Top view of the bilayer MoS₂/WS₂ heterostructure. (c) Predicted band structure of the bilayer MoS₂/WS₂ heterostructure. (d) Band structure of the bilayer MoS₂/WS₂ heterostructure calculated by OpenMX. (e) The HamGNN-predicted and OpenMX-calculated orbital energies for the C60 cluster. (f) The HamGNN-predicted and OpenMX-calculated wavefunction for the highest occupied molecular orbital (HOMO) of the C60 cluster.

Tests on low-dimensional materials

The above discussions have demonstrated the accuracy of the universal HamGNN model on bulk materials. Now, we will further explore its generalizability and prediction accuracy in the field of low-dimensional materials. We constructed a two-dimensional heterostructure consisting of MoS₂/WS₂ with a twist angle of 3.5°, as shown in Figure 5(a) and Figure 5(b). This structure comprises 1625 atoms and exhibits a higher level of complexity compared to the bulk structures typically included in our training set. The universal HamGNN model effectively captures the interatomic interactions within the bilayer MoS₂/WS₂ heterostructure and demonstrates excellent agreement between the energy bands predicted by the universal HamGNN model (Figure 5(c)) and those calculated by OpenMX (Figure 5(d)). The universal HamGNN model was further tested on the C60 cluster. Figure 5(e) demonstrates a good alignment between the predicted energy level of the C60 cluster near the band gap and the results obtained through DFT calculations, while Figure 5(f) illustrates a close match between the predicted wave function and the wave function calculated by DFT. These tests demonstrate the powerful and wide applicability of the universal HamGNN model in various material systems, ranging from bulk crystals to low-dimensional materials. By utilizing this universal model, researchers can explore a vast array of low-dimensional materials with tailored functionalities and properties.

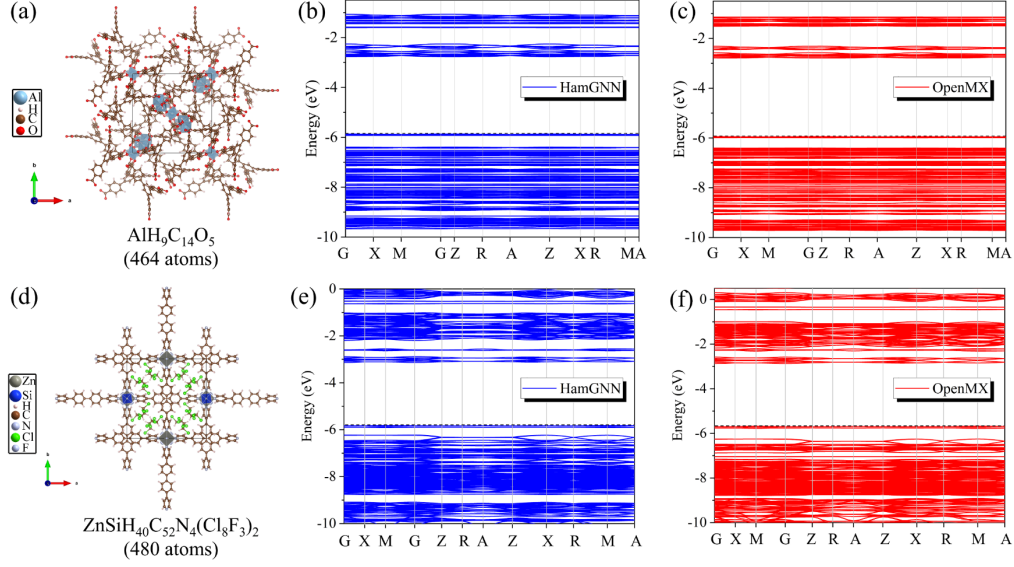


Figure 6. (a) The crystal structure, (b) predicted energy bands, and (c) calculated energy bands of $\text{AlH}_9\text{C}_{14}\text{O}_5$. (d) The crystal structure, (e) predicted energy bands, and (f) calculated energy bands of $\text{ZnSiH}_{40}\text{C}_{52}\text{N}_4(\text{Cl}_8\text{F}_3)_2$.

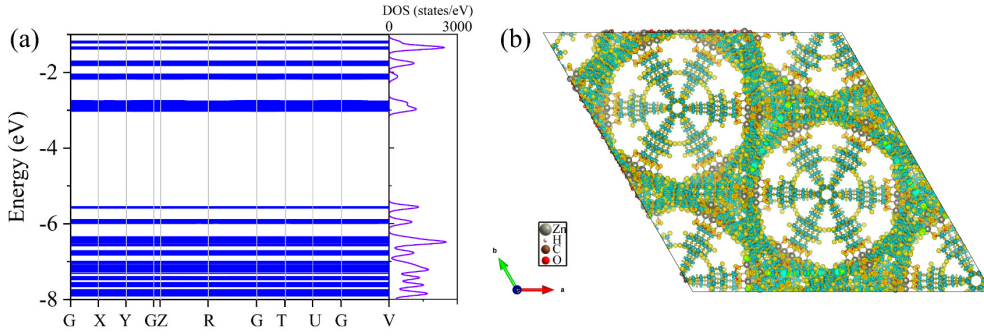


Figure 7. (a) The predicted energy bands and (b) the difference charge density of $\text{Zn}_4\text{H}_{28}\text{C}_{58}\text{O}_{13}$. The yellow and blue colors represent charge accumulation and depletion, respectively.

Application of the universal Hamiltonian model to large-scale hybrid inorganic-organic crystals

In this section, we will utilize the universal Hamiltonian model to predict the electronic structures of metal-organic frameworks (MOFs)⁵⁴⁻⁵⁶ and demonstrate its efficiency and broad applicability in studying hybrid inorganic-organic crystals. MOFs⁵⁴⁻⁵⁶ are a type of porous material composed of metal ions and organic ligands, forming intricate three-dimensional network structures. The electronic structures^{57, 58} of MOF materials, such as band structure, electron density of states, and electron orbital distribution, directly influence the conductivity, optical properties, and potential applications in fields like photocatalysis and photovoltaics. However, due to the complex structure and large size of MOFs, accurately predicting their electronic structures using DFT can be challenging and computationally expensive.

By training on inorganic crystal structures available on the Materials Project, we have successfully developed a universal Hamiltonian neural network model that demonstrates high accuracy across various types of inorganic materials. However, this

model lacks training on organic crystal structures which encompass a significant number of covalent bonds. Consequently, when predicting some complex organic or hybrid inorganic-organic crystal structures, the model may encounter certain challenges. In order to further improve the accuracy of the model in predicting the Hamiltonian matrix of organic crystal structures, we employed incremental training to further train and fine-tune the model. We selected approximately only 1800 small MOF crystal structures from the QMOF⁵⁹ Database, with a maximum number of atoms per unit cell not exceeding 50, as the training set. Subsequently, we further restart the training of the previous universal HamGNN model using this dataset. The mean absolute error of the Hamiltonian matrix predicted by the further trained model on this MOF dataset is about 3.95 meV. The comparison between the Hamiltonian matrix elements predicted by HamGNN and those calculated by OpenMX on the small MOF dataset is shown in Supplementary Figure S8. After undergoing fine-tuning, the Hamiltonian model demonstrates enhanced precision in inferring the interactions among diverse covalent bonds within organic crystalline materials.

As a test, we used the model to predict the electronic structures of two MOF materials, $\text{AlH}_9\text{C}_{14}\text{O}_5$ and $\text{ZnSiH}_{40}\text{C}_{52}\text{N}_4(\text{Cl}_8\text{F}_3)_2$, and compared them with the results obtained from DFT calculations. Both crystal structures have complex topological configurations, and the latter even contains up to seven elements, posing great challenges to the model. Despite these complexities, the energy bands predicted by the model for both structures are in good agreement with those obtained from DFT calculations, as shown in Figure 6. This indicates that our developed model has high accuracy and reliability when dealing with complex organic crystal structures. Furthermore, the difference charge densities predicted for these two materials achieved excellent consistency with those obtained from DFT calculations, as shown in Supplementary Figure S9. Moreover, the speed of DFT calculation can be improved by several orders of magnitude by using machine learning methods. Taking $\text{ZnSiH}_{40}\text{C}_{52}\text{N}_4(\text{Cl}_8\text{F}_3)_2$ as an example, the time cost of using DFT methods is as high as 181 core·hours. However, the HamGNN model only takes 0.33 core·hours to complete the computational task.

With further research, some more complex and larger MOF materials show wide application prospects in catalyst design, photovoltaics, solar cells, light-emitting diodes (LEDs), etc⁶⁰. In this case, the advantage of using machine learning Hamiltonian models for computing large MOF materials becomes apparent. To demonstrate the applicability of our model, we conducted electronic structure calculations on a specific MOF crystal known as $\text{Zn}_4\text{H}_{28}\text{C}_{58}\text{O}_{13}$ (labeled as c6ce00407e_c6ce00407e6_clean) from the CoRE MOF database^{61, 62}. This particular crystal is one of the largest MOFs in the CoRE MOF database and possesses a large unit cell containing 5562 atoms. Using the universal HamGNN model, we obtain the energy bands and density of states for $\text{Zn}_4\text{H}_{28}\text{C}_{58}\text{O}_{13}$, as shown in Figure. 7. It can be seen from the energy bands that the structure has a bandgap value of about 2.6 eV, which makes it a suitable candidate for various applications in the field of energy harvesting and emission, such as photovoltaics and light-emitting diodes. Furthermore, we also predicted the difference charge density of this structure and observed that

electrons primarily transfer from Zn and C atoms to the oxygen and hydrogen atoms. The difference charge density can visually display the distribution of static electric potential and possible catalytic sites.

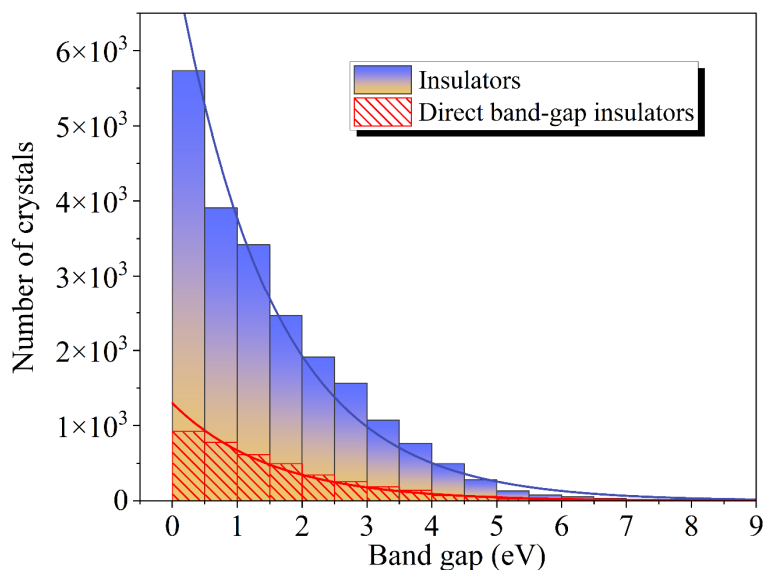


Figure 8. Histogram of the bandgap distribution for insulators and direct bandgap insulators.

High-throughput predictions of electronic structures for materials in the GeNOME dataset

The GeNOME⁶³ database is a remarkable collection of 2.2 million stable crystal structures that have been discovered using large-scale active learning techniques. This dataset has expanded our knowledge of stable materials by almost an order of magnitude, providing an unprecedented opportunity for in-depth investigation of material properties and the development of novel machine learning models. Despite the stability of the GeNOME materials, their electronic structures remain largely unknown. While first-principles calculations could in principle determine the electronic structures, performing such computationally expensive calculations for the entire database would be prohibitively costly. Therefore, we employ the universal HamGNN model, which enables fast, accurate, and high-throughput predictions of electronic structures across the whole dataset, offering valuable insights for future theoretical predictions and experimental synthesis.

We employ the universal HamGNN model to compute the electronic structures of a total of 188,722 structures from the GeNOME database. Our analysis reveals that 21,973 of these structures exhibit insulating properties, with 3,940 being direct bandgap insulators and 18,033 being indirect band gap insulators. Direct bandgap materials have broad applications in photovoltaics and light-emitting devices due to their efficient electron-hole recombination processes. The histogram in Figure 8 illustrates the distribution of bandgaps for both insulators and direct bandgap insulators. The figure shows an exponential decline in the number of insulators as their bandgap values increase. Utilizing this universal Hamiltonian model not only enables the rapid selection of materials with direct band gaps but also facilitates the

identification of crystals exhibiting flat bands near the Fermi level. This unique electronic structure holds great promise as it can give rise to intriguing properties and phenomena, such as fractional quantum Hall effects, Bose-Einstein condensation, unconventional superconductivity, and strong correlation effects⁶⁴⁻⁶⁶. After conducting a search using HamGNN, we have successfully identified 5109 crystals with flat bands. This number is nearly double compared to the 2379 flat-band crystals listed in the Materials Flatband Database. In order to demonstrate the accuracy of this approach for the GeNOME dataset, we will take a gapless material and a semiconductor material with flat bands, Ti_6SeS_7 and $\text{K}_3\text{SrZr}_6\text{BI}_{18}$, from the GeNOME database as examples. Notably, these two crystals are currently not included in the Materials Project.

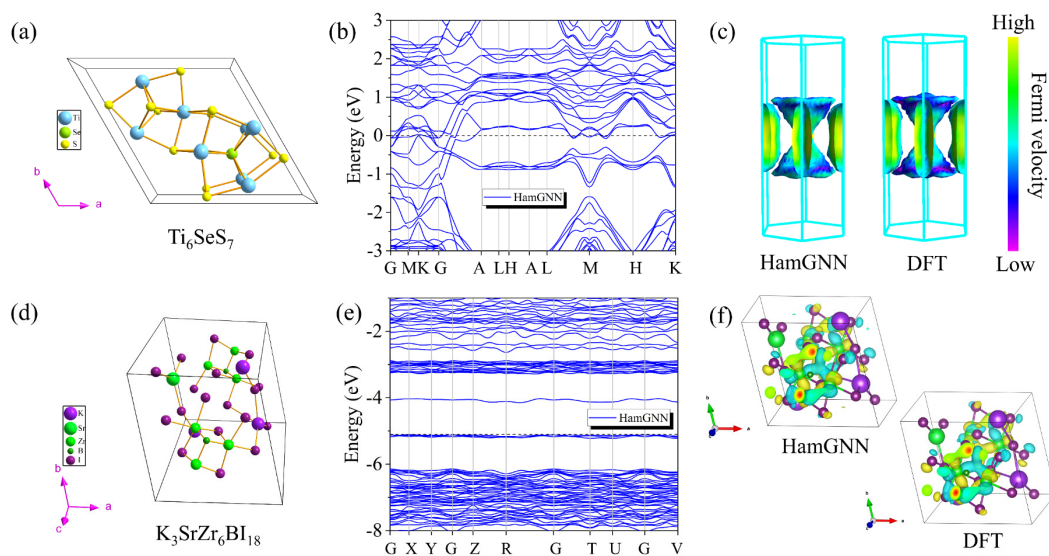


Figure 9. (a) The crystal structure of Ti_6SeS_7 . (b) The predicted energy bands of Ti_6SeS_7 by the universal HamGNN model. (c) Comparison between the predicted Fermi surface of Ti_6SeS_7 and those calculated using DFT. (d) The predicted energy bands of $\text{K}_3\text{SrZr}_6\text{BI}_{18}$ by the universal HamGNN model. (e) The predicted energy bands of $\text{K}_3\text{SrZr}_6\text{BI}_{18}$ by the universal HamGNN model. (f) Comparison of the HamGNN predicted LUMO wave function at the G point and the DFT calculated LUMO wave function at the G point for $\text{K}_3\text{SrZr}_6\text{BI}_{18}$.

The prediction results for the electronic structures of Ti_6SeS_7 and $\text{K}_3\text{SrZr}_6\text{BI}_{18}$ are presented in Figure 9. In Figure 9(b), the predicted energy bands of Ti_6SeS_7 exhibit excellent agreement with the results obtained from DFT calculations, as shown in Supplementary Figure S10(a). In addition, we compare the Fermi surfaces obtained from our predicted Hamiltonian matrices with those computed using DFT in Figure 9(c). The comparison reveals remarkable similarities between the two Fermi surfaces, further validating the accuracy and reliability of our predictions. As shown in Figure 9(e), the presence of flat bands at both HOMO and LUMO levels in the predicted energy bands of $\text{K}_3\text{SrZr}_6\text{BI}_{18}$ is apparent. The predicted flat band electronic structure is validated by the DFT calculation (see Supplementary Figure. S10(b)). Figure 9(f) reveals that the flat band at the LUMO level is primarily formed by atomic orbitals from Zr and I atoms. This observation suggests a significant contribution from these

elements to the unique electronic behavior exhibited by this compound. Furthermore, it is worth noting that our predicted wave function aligns well with actual calculated results, reinforcing the accuracy of our computational approach.

Discussion

This work not only achieves a universal electronic Hamiltonian model for the entire periodic table but also demonstrates its practicality and reliability through successful predictions of electronic structures in various materials. We propose a framework to achieve a universal Hamiltonian model by training HamGNN on the crystal's Hamiltonian matrices obtained from the Materials Project or other large datasets. We have found that incorporating energy eigenvalues in the second training step, in addition to Hamiltonian matrices, is crucial for achieving a truly universal Hamiltonian model.

The universal model can accurately capture not only simple systems but also those containing uncommon or rare transition metal elements. One notable advantage of this universal model is its ability to handle complex multi-element systems comprising more than five elements. This capability opens up new possibilities for studying and understanding the electronic properties of advanced materials that often involve intricate combinations of chemical elements. The reliable framework provided by the universal electronic Hamiltonian model allows for efficient predictions of electronic structures across the periodic table. While the current publicly available GeNOME dataset contains 380,000 stable crystal structures, the complete dataset encompasses an unprecedented 2.2 million stable materials discovered through large-scale active learning.

We anticipate that with more structures being made publicly accessible, the Universal Hamiltonian model can be leveraged to rapidly and accurately compute the electronic band structures of this vast collection through high-throughput calculations with much less computational cost compared to first-principle calculations. This would enable the identification of a significantly larger number of crystals with desirable electronic properties, driving further advancements in materials science. This breakthrough has significant implications for material design, catalysis, electronics, and other fields that heavily rely on accurate knowledge and understanding of electronic properties.

Methods

Network details.

The equivariant node features are $32 \times 0o + 128 \times 0e + 128 \times 1o + 64 \times 1e + 128 \times 2e + 32 \times 2o + 64 \times 3o + 32 \times 3e + 32 \times 4o + 32 \times 4e + 16 \times 5o + 8 \times 5e + 8 \times 6e$, where '32×0o' means that there are 32 channels in this feature part, and the features in each channel are $O(3)$ irreducible representations with $l = 0$ and odd parity. The node features utilized in this universal HamGNN model surpass those employed in our previous work²² to enhance the network capacity for describing the entire periodic table. The universal HamGNN model has five orbital convolution layers. The spherical harmonic basis functions used to expand the interatomic directions are

0e + 1o + 2e + 3o + 4e + 5o + 6e. The interatomic distance between atom i and its neighboring atom j , which falls within the cutoff radius r_c , is expanded utilizing the Bessel basis function:

$$B(|\boldsymbol{\tau}_{ij}|) = \sqrt{\frac{2}{r_c}} \frac{\sin\left(\frac{n\pi|\boldsymbol{\tau}_{ij}|}{r_c}\right)}{|\boldsymbol{\tau}_{ij}|} \quad (4)$$

The atomic neighbors are determined based on the cutoff radius of each atom's orbital basis. The interatomic distance is expanded using a series of Bessel functions with $n = 1, 2, \dots, N_b$, where N_b represents the number of Bessel basis functions. In this study, N_b is set to 64.

DFT details.

All the Hamiltonian matrices in the training set were computed using the PBE functional, with a Monkhorst-pack grid of $6 \times 6 \times 6$, and a convergence criterion of 1.0×10^{-8} Hartree. The energy cutoff employed for discretizing the real space is set at 200 Rydberg.

Reference

1. Marzari, N., Ferretti, A. & Wolverton, C. Electronic-structure methods for materials design. *Nat. Mater.* **20**, 736-749 (2021).
2. McCardle, K. Predicting electronic structure calculation results. *Nat. Comput. Sci.* **3**, 915-915 (2023).
3. Chen, Z. et al. Evolution of the electronic structure in open-shell donor-acceptor organic semiconductors. *Nat. Commun.* **12**, 5889 (2021).
4. Dzade, N. Y. First-principles insights into the electronic structure, optical and band alignment properties of earth-abundant $\text{Cu}_2\text{SrSnS}_4$ solar absorber. *Sci. Rep.* **11**, 4755 (2021).
5. Reidy, K. et al. Direct imaging and electronic structure modulation of moire superlattices at the 2D/3D interface. *Nat. Commun.* **12**, 1290 (2021).
6. Pederson, R., Kalita, B. & Burke, K. Machine learning and density functional theory. *Nat. Rev. Phys.* **4**, 357-358 (2022).
7. Schleder, G. R., Padilha, A. C. M., Acosta, C. M., Costa, M. & Fazzio, A. From DFT to machine learning: recent approaches to materials science—a review. *J. Phys.: Mater.* **2**, 032001 (2019).
8. Makkar, P. & Ghosh, N. N. A review on the use of DFT for the prediction of the properties of nanomaterials. *RSC Adv.* **11**, 27897-27924 (2021).
9. Jones, R. O. Density functional theory: Its origins, rise to prominence, and future. *Rev. Mod. Phys.* **87**, 897-923 (2015).
10. Hegde, G. & Bowen, R. C. Machine-learned approximations to Density Functional Theory Hamiltonians. *Sci. Rep.* **7**, 42669 (2017).
11. Li, H. C., Collins, C., Tanha, M., Gordon, G. J. & Yaron, D. J. A Density Functional Tight Binding Layer for Deep Learning of Chemical Hamiltonians. *J. Chem. Theory. Comput.* **14**, 5764-5776 (2018).

12. Schutt, K. T., Gastegger, M., Tkatchenko, A., Muller, K. R. & Maurer, R. J. Unifying machine learning and quantum chemistry with a deep neural network for molecular wavefunctions. *Nat. Commun.* **10**, 5024 (2019).
13. Gastegger, M., McSloy, A., Luya, M., Schutt, K. T. & Maurer, R. J. A deep neural network for molecular wave functions in quasi-atomic minimal basis representation. *J. Chem. Phys.* **153**, 044123 (2020).
14. Unke, O. T. et al. SE(3)-equivariant prediction of molecular wavefunctions and electronic densities. Preprint at <https://ui.adsabs.harvard.edu/abs/2021arXiv210602347U> (2021).
15. Wang, Z. F. et al. Machine learning method for tight-binding Hamiltonian parameterization from ab-initio band structure. *npj Comput. Mater.* **7**, 11 (2021).
16. Westermayr, J. & Maurer, R. J. Physically inspired deep learning of molecular excitations and photoemission spectra. *Chem. Sci.* **12**, 10755-10764 (2021).
17. Li, H. et al. Deep-learning density functional theory Hamiltonian for efficient ab initio electronic-structure calculation. *Nat. Comput. Sci.* **2**, 367-377 (2022).
18. Nigam, J., Willatt, M. J. & Ceriotti, M. Equivariant representations for molecular Hamiltonians and N-center atomic-scale properties. *J. Chem. Phys.* **156**, 014115 (2022).
19. Schattauer, C., Todorović, M., Ghosh, K., Rinke, P. & Libisch, F. Machine learning sparse tight-binding parameters for defects. *npj Comput. Mater.* **8**, 116 (2022).
20. Zhang, L. et al. Equivariant analytical mapping of first principles Hamiltonians to accurate and transferable materials models. *npj Comput. Mater.* **8**, 158 (2022).
21. Gong, X. et al. General framework for E(3)-equivariant neural network representation of density functional theory Hamiltonian. *Nat Commun* **14**, 2848 (2023).
22. Zhong, Y., Yu, H., Su, M., Gong, X. & Xiang, H. Transferable equivariant graph neural networks for the Hamiltonians of molecules and solids. *npj Comput. Mater.* **9**, 182 (2023).
23. Ye, Y. F., Wang, Q., Lu, J., Liu, C. T. & Yang, Y. High-entropy alloy: challenges and prospects. *Mater. Today* **19**, 349-362 (2016).
24. Feng, R. et al. High-throughput design of high-performance lightweight high-entropy alloys. *Nat Commun* **12**, 4329 (2021).
25. George, E. P., Raabe, D. & Ritchie, R. O. High-entropy alloys. *Nat. Rev. Mater.* **4**, 515-534 (2019).
26. Zhang, R.-Z. & Reece, M. J. Review of high entropy ceramics: design, synthesis, structure and properties. *J. Mater. Chem. A* **7**, 22148-22162 (2019).
27. Oses, C., Toher, C. & Curtarolo, S. High-entropy ceramics. *Nat. Rev. Mater.* **5**, 295-309 (2020).
28. Zhao, P., Xiao, C. & Yao, W. Universal superlattice potential for 2D materials from twisted interface inside h-BN substrate. *npj 2D Mater. Appl.* **5**, 38 (2021).
29. Chen, C. e. a. Expanding materials science with universal many-body graph neural networks. *Nat. Comput. Sci.* **2**, 703-704 (2022).
30. Chen, C. & Ong, S. P. A universal graph deep learning interatomic potential for the periodic table. *Nat. Comput. Sci.* **2**, 718-728 (2022).
31. Takamoto, S. et al. Towards universal neural network potential for material discovery applicable to arbitrary combination of 45 elements. *Nat. Commun.* **13**, 2991 (2022).
32. Deng, B. et al. CHGNet as a pretrained universal neural network potential for charge-informed atomistic modelling. *Nat. Mach. Intell.* **5**, 1031-1041 (2023).

33. Batatia, I. et al. A foundation model for atomistic materials chemistry. Preprint at <https://ui.adsabs.harvard.edu/abs/2024arXiv240100096B> (2023).
34. Jain, A. et al. Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *APL Mater.* **1**, 011002 (2013).
35. Jain, A. et al., The Materials Project: Accelerating Materials Design Through Theory-Driven Data and Tools. In *Handbook of Materials Modeling*, 2018; pp 1-34.
36. Wang, H., Zhang, L. F., Han, J. Q. & E, W. N. DeePMD-kit: A deep learning package for many-body potential energy representation and molecular dynamics. *Comput. Phys. Commun.* **228**, 178-184 (2018).
37. Behler, J. Atom-centered symmetry functions for constructing high-dimensional neural network potentials. *J. Chem. Phys.* **134**, 074106 (2011).
38. Pinheiro, M., Ge, F. C., Ferre, N., Dral, P. O. & Barbatti, M. Choosing the right molecular machine learning potential. *Chem. Sci.* **12**, 14396-14413 (2021).
39. Xie, T. & Grossman, J. C. Crystal Graph Convolutional Neural Networks for an Accurate and Interpretable Prediction of Material Properties. *Phys. Rev. Lett.* **120**, 145301 (2018).
40. Chmiela, S. et al. Machine learning of accurate energy-conserving molecular force fields. *Sci. Adv.* **3**, e1603015 (2017).
41. Glielmo, A., Sollich, P. & De Vita, A. Accurate interatomic force fields via machine learning with covariant kernels. *Phys. Rev. B* **95**, 214302 (2017).
42. Ozaki, T. & Kino, H. Numerical atomic basis orbitals from H to Kr. *Phys. Rev. B* **69**, 195113 (2004).
43. Ozaki, T. Variationally optimized atomic orbitals for large-scale electronic structures. *Phys. Rev. B* **67**, 155108 (2003).
44. Weinert, U. Spherical Tensor Representation. *Arch. Ration. Mech. An.* **74**, 165-196 (1980).
45. Morrison, M. A. & Parker, G. A. A Guide to Rotations in Quantum-Mechanics. *Aust. J. Phys.* **40**, 465-497 (1987).
46. Grisafi, A., Wilkins, D. M., Csanyi, G. & Ceriotti, M. Symmetry-Adapted Machine Learning for Tensorial Properties of Atomistic Systems. *Phys. Rev. Lett.* **120**, 036002 (2018).
47. Thomas, N. et al. Tensor field networks: Rotation- and translation-equivariant neural networks for 3D point clouds. Preprint at <https://arxiv.org/abs/1802.08219v3> (2018).
48. Naveed, H. et al. A Comprehensive Overview of Large Language Models. Preprint at <https://ui.adsabs.harvard.edu/abs/2023arXiv230706435N> (2023).
49. Zhao, W. X. et al. A Survey of Large Language Models 2023. Preprint at <https://ui.adsabs.harvard.edu/abs/2023arXiv230318223Z> (2023).
50. Zeng, L. et al. Superconductivity and non-trivial band topology in high-entropy carbonitride $\text{Ti}_{0.2}\text{Nb}_{0.2}\text{Ta}_{0.2}\text{Mo}_{0.2}\text{W}_{0.2}\text{C}_{1-x}\text{N}_x$. *The Innovation Materials* **1**, 100042 (2023).
51. Zeng, L. et al. Superconductivity in the High-Entropy Ceramics $\text{Ti}_{0.2}\text{Zr}_{0.2}\text{Nb}_{0.2}\text{Mo}_{0.2}\text{Ta}_{0.2}\text{C}_x$ with Possible Nontrivial Band Topology. *Adv. Sci.* **11**, e2305054 (2024).
52. Zeng, L. et al. Discovery of the High-Entropy Carbide Ceramic Topological Superconductor Candidate $\text{Ti}_{0.2}\text{Zr}_{0.2}\text{Nb}_{0.2}\text{Mo}_{0.2}\text{Ta}_{0.2}\text{C}$. *Adv. Funct. Mater.* **33**, 2301929 (2023).
53. Zhao, S. et al. Designing a Beryllium-Free Deep-Ultraviolet Nonlinear Optical Material without a Structural Instability Problem. *J. Am. Chem. Soc.* **138**, 2961-2964 (2016).

54. Zhang, J.-W. et al. Composition, structure and electrochemical performance of LiSiPSCl electrolyte with Li/Li-In anodes in all-solid-state batteries. *Electrochim. Acta* **461**, 142691 (2023).
55. Baumann, A. E., Burns, D. A., Liu, B. & Thoi, V. S. Metal-organic framework functionalization and design strategies for advanced electrochemical energy storage devices. *Commun. Chem.* **2**, 86 (2019).
56. Felix Sahayaraj, A. et al. Metal–Organic Frameworks (MOFs): The Next Generation of Materials for Catalysis, Gas Storage, and Separation. *J. Inorg. Organomet. Polym.* **33**, 1757-1781 (2023).
57. Mancuso, J. L., Mroz, A. M., Le, K. N. & Hendon, C. H. Electronic Structure Modeling of Metal-Organic Frameworks. *Chem Rev* **120**, 8641-8715 (2020).
58. Rosen, A. S. et al. High-throughput predictions of metal–organic framework electronic properties: theoretical challenges, graph neural networks, and data exploration. *npj Comput. Mater.* **8**, 112 (2022).
59. Rosen, A. S. et al. Machine learning the quantum-chemical properties of metal–organic frameworks for accelerated materials discovery. *Matter* **4**, 1578-1597 (2021).
60. Stassen, I. et al. An updated roadmap for the integration of metal-organic frameworks with electronic devices and chemical sensors. *Chem. Soc. Rev.* **46**, 3185-3241 (2017).
61. Chung, Y. G. et al. Computation-Ready, Experimental Metal–Organic Frameworks: A Tool To Enable High-Throughput Screening of Nanoporous Crystals. *Chem. Mater.* **26**, 6185-6192 (2014).
62. Chung, Y. G. et al. Advances, Updates, and Analytics for the Computation-Ready, Experimental Metal–Organic Framework Database: CoRE MOF 2019. *J. Chem. Eng. Data* **64**, 5985-5998 (2019).
63. Merchant, A. et al. Scaling deep learning for materials discovery. *Nature* **624**, 80-85 (2023).
64. He, C. et al. Flat-band based high-temperature ferromagnetic semiconducting state in the graphitic C4N3 monolayer. *Fundam. Res.* (2023).
65. Bhattacharya, A., Timokhin, I., Chatterjee, R., Yang, Q. & Mishchenko, A. Deep learning approach to genome of two-dimensional materials with flat electronic bands. *npj Comput. Mater.* **9**, 101 (2023).
66. Huber, S. D. & Altman, E. Bose condensation in flat bands. *Phys. Rev. B* **82**, 184502 (2010).

Acknowledgment

We acknowledge financial support from the Ministry of Science and Technology of the People's Republic of China (No. 2022YFA1402901), NSFC (grants No. 11825403, 11991061, 12188101), and the Guangdong Major Project of the Basic and Applied Basic Research (Future functional materials under extreme conditions--2021B0301030005).

Data and code availability

The HamGNN code is publicly available from <https://github.com/QuantumLab-ZY/HamGNN>.

The trained network weights for the universal model, the predicted energy bands for the test set, and the identified crystals with flat bands are available on Zenodo (<https://zenodo.org/records/10827117>).

Author contributions

H.J.X. and X.G.G. supervised the project for the universal Hamiltonian framework. Y.Z. wrote the implementation codes of the universal Hamiltonian framework, and conducted network training and testing. X.Y.G. constructed the crystal structure of $\text{Li}_{10}\text{Si}_{1.5}\text{P}_{1.5}\text{S}_{11.5}\text{Cl}_{0.5}$ and discussed the corresponding findings. H.Y.Y. used the universal model to perform the high throughput calculation of electronic structures for the GeNOME dataset. J.H.Y. checked the results in the manuscript. Y.Z., H.Y.Y., H.J.X., and X.G.G. prepared the manuscript. All authors discussed the results and provided comments on the manuscript.

Competing interests

The authors declare no competing interests.