# nature computational science

**Article** 

https://doi.org/10.1038/s43588-024-00661-0

# Virtual node graph neural network for full phonon prediction

Received: 30 November 2023

Accepted: 13 June 2024

Published online: 12 July 2024



Ryotaro Okabe  $\mathbb{D}^{12,12} \boxtimes$ , Abhijatmedhi Chotrattanapituk  $\mathbb{D}^{1,3,12}$ , Artittaya Boonkird<sup>1,4</sup>, Nina Andrejevic<sup>5</sup>, Xiang Fu $\mathbb{D}^3$ , Tommi S. Jaakkola<sup>3</sup>, Qichen Song  $\mathbb{D}^6$ , Thanh Nguyen<sup>1,4</sup>, Nathan Drucker  $\mathbb{D}^{1,7}$ , Sai Mu<sup>8</sup>, Yao Wang  $\mathbb{D}^9$ , Bolin Liao  $\mathbb{D}^{10}$ , Yongqiang Cheng  $\mathbb{D}^{11} \boxtimes \mathbb{A}$  Mingda Li $\mathbb{D}^{1,4} \boxtimes$ 

Understanding the structure-property relationship is crucial for designing materials with desired properties. The past few years have witnessed remarkable progress in machine-learning methods for this connection. However, substantial challenges remain, including the generalizability of models and prediction of properties with materials-dependent output dimensions. Here we present the virtual node graph neural network to address the challenges. By developing three virtual node approaches, we achieve  $\Gamma$ -phonon spectra and full phonon dispersion prediction from atomic coordinates. We show that, compared with the machine-learning interatomic potentials, our approach achieves orders-of-magnitude-higher efficiency with comparable to better accuracy. This allows us to generate databases for Γ-phonon containing over 146,000 materials and phonon band structures of zeolites. Our work provides an avenue for rapid and high-quality prediction of phonon band structures enabling materials design with desired phonon properties. The virtual node method also provides a generic method for machine-learning design with a high level of flexibility.

The structure–property relationship defines one of the most fundamental questions in materials science<sup>1,2</sup>. The ubiquitous presence of structure–property relationships profoundly influences almost all branches of materials sciences, including structural materials<sup>3</sup>, energy harvesting, conversion and storage materials<sup>4–6</sup>, catalysts<sup>7</sup> and polymers<sup>8</sup>, and quantum materials<sup>9</sup>. However, building an informative structure–property relationship can be nontrivial despite its central importance to materials design. On the one hand, the number of stable

structures grows exponentially with unit-cell size<sup>10</sup>, and the structure design efforts have been largely limited to crystalline solids with relatively small unit cells. On the other hand, certain material properties are challenging to acquire due to experimental or computational complexities. In the past few years, data-driven and machine-learning methods have played an increasingly important role in materials science and substantially boosted the research on building structure–property relationships<sup>11–13</sup>. Complex structures such as porous materials<sup>14,15</sup>,

<sup>1</sup>Quantum Measurement Group, Massachusetts Institute of Technology, Cambridge, MA, USA. <sup>2</sup>Department of Chemistry, Massachusetts Institute of Technology, Cambridge, MA, USA. <sup>3</sup>Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, USA. <sup>4</sup>Department of Nuclear Science and Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA. <sup>5</sup>Argonne National Laboratory, Lemont, IL, USA. <sup>6</sup>Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA, USA. <sup>7</sup>Department of Applied Physics, School of Engineering and Applied Sciences, Harvard University, Cambridge, MA, USA. <sup>8</sup>SmartState Center for Experimental Nanoscale Physics, Department of Physics and Astronomy, University of South Carolina, Columbia, SC, USA. <sup>9</sup>Department of Chemistry, Emory University, Atlanta, GA, USA. <sup>10</sup>Department of Materials, University of California, Santa Barbara, Santa Barbara, CA, USA. <sup>11</sup>Chemical Spectroscopy Group, Spectroscopy Section, Neutron Scattering Division, Oak Ridge National Laboratory, Oak Ridge, TN, USA. <sup>12</sup>These authors contributed equally: Ryotaro Okabe, Abhijatmedhi Chotrattanapituk.

Se-mail: rokabe@mit.edu; chengy@ornl.gov; mingda@mit.edu

nanoalloys 16,17 and grain boundaries 18 are becoming more feasible to handle, and properties ranging from mechanical strength to quantum ordering can be learned with increased confidence 19,20. One particularly powerful approach is through graph neural networks (GNNs)<sup>21</sup>. By representing atoms as graph nodes and interatomic bonds as graph edges, GNNs provide a natural representation of molecules and materials. For crystalline solids, crystallographic symmetry offers a further boost on the GNN performance, with a few symmetry-augmented GNNs being proposed<sup>22-24</sup>. A few fundamental challenges still exist. For one, many materials properties are not naturally represented as a weighted aggregation of each atom in real space, such as reciprocal and energy space properties. For another, the output property length is usually fixed, like the heat capacity<sup>25</sup>, as a single scalar. In contrast, many materials' properties have unique degrees of dimensions, such as the number of electronic and phononic bands<sup>26</sup>, frequency ranges with optical responses, and features of magnetic structure such as propagation vectors.

In this Article, we propose the virtual node graph neural network (VGNN) as a generically applicable approach to augment GNN. In contrast to symmetry-augmented GNN, which focuses on reducing the input data volume, VGNN focuses on handling the output properties with variable or even arbitrary dimensions. To demonstrate the application of VGNN, we study materials' phonon spectra and dispersion relations, since phonon band structures are challenging to compute or measure with high computational costs and limited experimental resources. By using the phonon spectra as examples, we present three versions of VGNN: the vector virtual nodes (VVN), the matrix virtual nodes (MVN) and the momentum-dependent matrix virtual nodes (k-MVN). All three VGNN models take atomic structures as input without prior knowledge of interatomic forces. The VVN is the simplest VGNN that takes in crystal structure with m atoms and outputs 3m branches of Γ-phonon energies. The MVN is a more involved VGNN that shows higher accuracy for complex materials with slightly higher computational costs. Finally, the k-MVN is a VGNN that can predict full phonon band structure at arbitrary k points in the Brillouin zone. To achieve so, the crystal graphs contain 'virtual-dynamical matrices' (VDMs), which are matrix structures that resemble phonon dynamical matrices<sup>27</sup>. Instead of performing direct ab initio calculations on each material, all matrix elements are learned from the neural network optimization process using training data composed of all other materials. We prove that the proposed VGNN approach could reduce the computational cost and run time without sacrificing accuracy, compared with the more common machine-learning interatomic potential (MLIP) approach. Our work offers an efficient technique that can compute zone-center phonon energies and full phonon band structures directly from atomic structures in complex materials and enables phonon property optimization within a larger structure design space. The prediction methods have enabled us to acquire relevant information on materials such as group velocities, heat capacities and density of states (DoS) as byproducts. Meanwhile, the virtual node structures also shed light on future flexible GNN design, that is, to use intermediate crucial quantities (for example, dynamical matrix) as key learning parameters without having to put target properties (for example, phonon band structures) as output.

#### Results

# Virtual node augmentation for GNNs

Figure 1 gives an overview of the VGNN method as a generic approach to augment GNN. For a crystal with m atoms per unit cell (Fig. 1a), a typical GNN model converts the crystal into a crystal graph, where each graph node represents an atom and each graph edge represents the interatomic bonding as shown in Fig. 1b. The node features associated with each atomic node (Fig. 1b, gray arrays) are updated by neighborhood nodes and edges connecting the nodes (Fig. 1b, gray arrows). After iterative layers of graph convolutions, m final-layer node features are

obtained that represent the atomic features (local features) from each of the *m* atoms. The final graph output (global feature) can be obtained by aggregating the final-layer node features into one fixed-sized output.

Figure 1c,d describes the general idea of VGNN that endows a GNN with greater flexibility for prediction. On top of the conventional. real-node GNN, virtual atoms are added into crystal (Fig. 1c, yellow nodes), which become the virtual nodes in the corresponding GNN (Fig. 1d, yellow nodes). Each model places the virtual atoms by different rules (see 'VVN method', 'MVN method' and 'k-MVN' sections in Methods for more details). As Fig. 1d illustrates, just like the bidirectional message passing between real atomic nodes (double-arrow gray lines), the message passing (double-arrow yellow lines) between virtual nodes is also bidirectional. On the other hand, to preserve the structure of the conventional GNN, the messages from real nodes to virtual (single-arrow gray-to-yellow gradient lines) are unidirectional. Given the flexibility of the choice of the virtual nodes, a VGNN gains flexibility to predict materials-dependent outputs with arbitrary lengths, solving the conventional mismatch between the number of node features m and that of the target predicted values n illustrated in Fig. 1d. We will introduce three VGNN methods for phonon prediction with increased levels of predictive power and complexity.

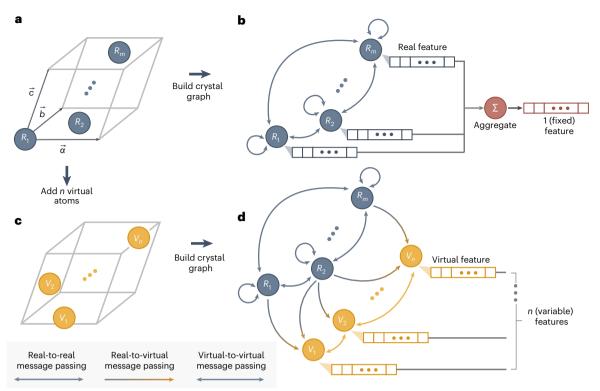
#### VVNs for Γ-phonon prediction

As illustrated in Fig. 1, VGNN makes it possible to adjust output dimension on the basis of input information with flexibility. We first introduce the VVN method, which is the simplest approach to acquire 3m phonon branches when inputting a crystal with m atoms per unit cell (see 'VVN method' section in Methods for more detail). Figure 2 shows the VVN approach to predict Γ-phonon spectra. Since the virtual nodes do not pass information to real nodes, there is additional flexibility in choosing the position of the virtual node. Without loss of generality, we assign the position of the virtual nodes evenly spaced along the diagonal line of the unit cell. The crystal graph is constructed with virtual and real nodes (Fig. 2a). After updating node features in each convolution layer, the feature vectors pass a linear layer so that virtual node features  $V_i$ ,  $i \in [1, 3m]$  are converted to 3m scalars, which represent the predicted Γ-phonon energies. Throughout this work, the GNN part is implemented through the Euclidean neural networks<sup>23</sup> that are aware of the crystallographic symmetry. Data preparation, neural network architectures and optimizations are described in Supplementary Information sections 1-3.

The main results using the VVN for Γ-phonon prediction are shown in Fig. 2b. The three-row spectral comparison plots are randomly selected samples from the test set within each error tertile (top-to-bottom rows are top-to-bottom performance tertiles, respectively). The first four columns are taken from the same database as the training set from high-quality density-functional perturbation theory (DFPT) calculations<sup>28</sup>, and the fifth column contains additional test examples with much larger unit cells from a frozen-phonon database<sup>29</sup>. Our results show that the prediction loss becomes larger and distributed broader as the input materials are more complicated (Fig. 2c). From the correlation plot of predicted and ground-truth phonon frequencies (Fig. 2d), most data points are along the diagonal line, indicating good prediction between VNN prediction and ground truth from DFPT calculations with the number of atoms per unit cell  $m \le 24$  (blue dots). For complex materials, the correlation performance is degraded (orange dots). More test results are shown in Supplementary Information sections 4.1-4.3.

#### MVNs for Γ-phonon with enhanced performance

In this section, we introduce another type of virtual nodes approach, the MVNs. The MVN approach performs better  $\Gamma$ -phonon prediction than VVN, especially for complex materials, with a slightly higher computational cost. Moreover, the structure of MVN lays the groundwork for the full phonon band structures to be discussed in the next section.



**Fig. 1**| **Overview of VGNN. a**, The atomic structure of a crystalline material with m atoms per unit cell, where atom i is represented by a dark-blue  $R_i$  real node.  $\vec{a}$ ,  $\vec{b}$  and  $\vec{c}$  are unit-cell lattice parameters. **b**, A GNN converts the atomic structures into a crystal graph. The GNN handles the nodes as the output gates. Therefore, the output dimension is restricted to the same number of nodes. However, each node output mostly contains local information dictated by the message passing scheme. Hence, after layers of graph convolutions (omitted for simplicity), the

final node features are aggregated into a single fixed-sized output feature.  $\mathbf{c}$ , A flexible number of n virtual atoms are added into the crystal structure, where virtual atom, j is represented by a yellow  $V_j$  virtual node.  $\mathbf{d}$ , After forming the crystal graph with both real and virtual nodes, the flexibility of virtual nodes enables the choices of output not necessarily from real-node aggregation but can have variable lengths in different spaces.

In MVN, m copies of virtual crystals are generated for material with m atoms per unit cell, and each copy contains m virtual nodes that share the same crystal structure as the real crystal (Fig. 3a). This results in a total of  $m^2$  virtual nodes  $V_{ij}$ ,  $i,j \in [1,m]$  with more involved node connectivity (see 'MVN method' section in Methods for more detail).

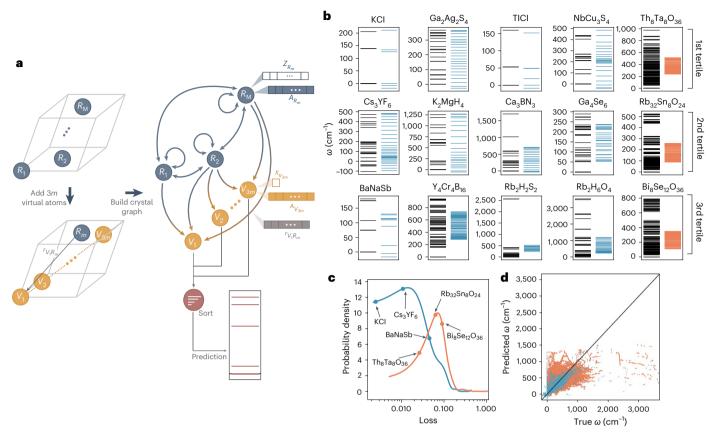
With this graph construction scheme, after the neural network training, the virtual nodes  $V_{ij}$  would capture the essence of the connection between  $R_i$ , and  $R_j$ . Hence, after the message passes in each convolutional layer, each virtual node feature is further converted into a three-by-three matrix. Each of  $V_{ij}$  is assembled to form a (i,j) block of a supermatrix  $\tilde{D}$  of shape (3m,3m). Given the structural similarity between this matrix and the dynamical matrix expressed in equation (2) with k=0, we predict  $\Gamma$ -point phonon energies squared by solving for m eigenvalues of the matrix m. It is still worthwhile mentioning that, although the matrix shares a similar feature with the dynamical matrix, the matrix elements are learned from neural network training and are not necessarily the matrix elements from the real dynamical matrix. An intuitive comparison is that the edge of GNN does not necessarily reflect true chemical bonding, but is more like an atomic neighbor connection.

The predicted phonons using MVN are summarized in Fig. 3b, which shares the same structure with Fig. 2b as error tertile plots from the high-quality DFPT database (blue) and database for complex materials (orange). MVN shows comparable performance with VVN for simple materials (Figs. 2c and 3c, blue curves) but shows substantial performance improvement for complex materials. The prediction loss distribution of MVN shows a heavier distribution toward a lower loss regime compared with VVN (Figs. 2c and 3c, orange curves), and the phonon frequencies in the correlation plot align better toward ground

truth (Figs. 2d and 3d, orange dots). More results of the MVN method are shown in Supplementary Information sections 4.1-4.3.

#### MVNs for phonon band structure prediction

The structure of MVN inspires us to take one step further and construct full momentum-dependent VDMs by taking into account the unit-cell translation, termed k-MVN. We construct VDMs following equation (2). In contrast to the MVN, which focuses on Γ-point phonons by taking  $\vec{k} = 0$ , here, in k-MVN, we include the phase factor  $e^{i\vec{k}\cdot\vec{T}}$  when defining the VDMs, where  $\vec{T}$  is the relative unit-cell translation of a neighboring unit-cell origin relative to the chosen reference unit cell  $\overrightarrow{T}_0$  (Fig. 4a). If a total number of t neighboring unit cells are included, each with translation  $T_h, h \in [0, t-1]$  (reference cell included), then a total t copies of MVN-type virtual nodes matrices will be generated, with a total number of  $tm^2$  virtual nodes  $V_{ij}^h, h \in [0, t-1], i, j \in [1, m]$  in k-MVN. To obtain the phonon band structure, each set of virtual nodes at a given  $\overrightarrow{T_h}$  needs to multiply by the phase factor  $e^{i\vec{k}\cdot\vec{T_h}}$ , and all virtual nodes at each  $\overrightarrow{T_h}$  are summed in equation (3) in Methods. Thanks to the graph connectivity within the cutoff radius (see 'k-MVN' section in Methods), only a small number of t is needed as long as crystal graph connectivity can be maintained. In practice, t is material-dependent, and t = 27 (nearest neighbor unit cells) is sufficient for many materials and does not need to go beyond t = 125 (next-nearest neighbor unit cells) in all cases. Intuitively, such a supercell approach resembles the ab initio band structure calculations with frozen phonons. To facilitate the training, phonons from selected high-symmetry points are included in the training data without using full phonon energies in the entire Brillouin zone. This substantially facilitates the training process while maintaining accuracy. More details are discussed in 'k-MVN' section in Methods.



**Fig. 2**| **The VVN method to predict Γ-point phonons. a**, Schematic of VVN model construction and prediction. For material with m atoms per unit cell, 3m virtual nodes are augmented along the diagonal vector  $\vec{v} = \vec{a} + \vec{b} + \vec{c}$  of the unit cell. We embedded the components of the crystal when building the GNN model. For instance, atomic numbers of the mth real atom ( $A_{R_m}$ ) and that of the 3mth virtual atom ( $A_{V_{3m}}$ ) are embedded as the attributes of each nodes. The atomic mass of the mth real atom ( $Z_{R_m}$ ) is set as the initial feature of that node. The relative position of the node  $V_1$  with respect to  $R_m$  is  $\vec{r}_{V_1R_m}$ , which is used to embed the edge attribute between the two nodes. The model predicts Γ-phonon spectra by sorting the scalar output features from virtual nodes. **b**, Spectral prediction samples in the test set within each error tertile compared with ground truth (black), test from the same database as the training set (blue) and a different database containing complex materials (orange). Phonon is displayed with a unit

of cm<sup>-1</sup>. **c**, Evaluation of the test accuracy through loss distribution on test datasets with the same color scheme as in **b**. The curves represent probability density of the model prediction loss when a random material is sampled from test datasets. Three points on each curve indicate loss values of example materials from different error tertile in **b**. Note that the loss axis is in log scale, which makes it possible to visualize both distributions in the same plot but exaggerates the width of the blue curve. **d**, Evaluation of the test accuracy through correlation plot between ground-truth and predicted phonon frequencies with the graph y = x as reference. The concentration of distribution at low loss regime of the distribution plot and the agreement along the diagonal line of the correlation plot for the test set (blue) indicates a good phonon prediction at least for relatively simple materials with the number of atoms per unit cell  $m \le 24$ . The loss becomes higher with reduced performance for more complex materials (orange).

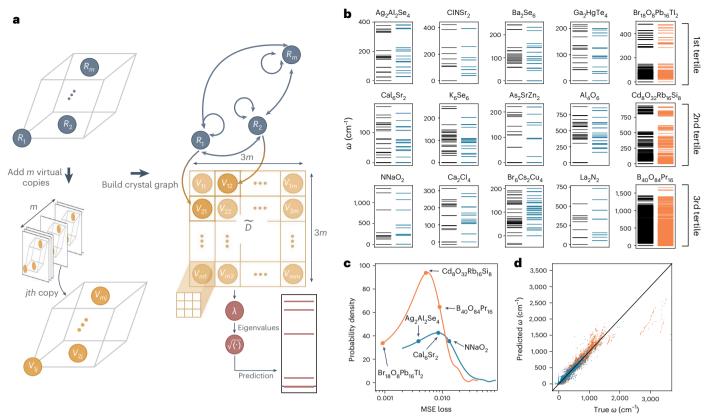
Figure 4b shows the prediction results of phonon band structures. Here, 12 materials are selected from the same dataset for training (blue color) and additional dataset for complex materials (orange color). Despite the complexity of a generic phonon band structure, the k-MVN model could predict the positions and the shapes of the phonon bands, such as gaps between different optical branches. The dispersion relations of the acoustic phonons are also well generated around the  $\Gamma$ -points on the left three columns, even though we do not enforce that acoustic  $\Gamma$ -phonons have to be gapless with zero energy known as the acoustic sum rule  $^{30}$ . This may enable the prediction of crystal stability for future works. While there are risks that prediction performance could be degraded for the phonon bands of higher frequencies, most of the predicted phonons follow the references, including the complex materials with more than 40 atoms per unit cell. More test results are shown in Supplementary Information sections 5.1–5.3.

## Model validation and benchmarking

We demonstrate the prediction of phonons directly from the materials' atomic coordinates, using three different types of virtual node augmentation approach. The comparison between them is summarized in Table 1. VVN directly acquires the phonon spectra directly from the

virtual nodes. The assignment of 3m virtual nodes ensures that the output phonon band number is always 3m for a crystal containing m atoms per primitive unit cell. The bottleneck step of this scheme is the message passing, where there can be up to  $O(m^2)$  connections. In MVN, instead of computing phonon energies directly, a VDM is constructed first, from which the phonon energies are solved as an eigenvalue problem. This step is crucial to gain robustness for complex materials prediction since intermediate quantities like force constants and dynamical matrices are considered more 'fundamental' than final phonon energies to reflect the interatomic interactions. The k-MVN goes one step further, using the unit-cell translations to generate the momentum dependence that could be used to obtain the full phonon band structure. This requires multiple MVN connections each representing different unit-cell translations. For k-MVN that considers t unit-cell translations, the memory requirement is stacked up t times of the MVN case, but the bottleneck step still is the eigenvalue solving.

To validate the necessity of virtual node, through the comparison of  $\Gamma$ -phonon prediction, we show that VGNN can substantially outperform the conventional GNN techniques without virtual nodes. One popular approach for flexible output dimension is to adjust the output dimensions by supplementing arbitrary values such as zeros to a



**Fig. 3** | **The MVN method to predict Γ-point phonons. a**, Augment  $m^2$  virtual nodes as m sets of virtual crystals (left) and the message passing scheme and postprocessing of virtual node features (right). The legends are the same as in Fig. 2. In contrast to VVN, where each node  $V_j$  is a scalar, here, each node  $V_{ij}$  is a 3 × 3 matrix. The phonon spectra in MVN are obtained by solving the eigenproblems instead of direct output, as done in VVN. **b**, Selected test examples within each error tertile. Tests from the same dataset as the training

set and additional tests containing complex materials are predicted in blue and orange, respectively. **c**, Comparison of prediction loss distribution with three examples of materials from each error tertile. **d**, The correlation plots of phonon frequencies with the graph y = x as reference. Better performance for MVN is achieved than VVN for complex materials (orange color), which can be seen from both the loss distribution and the average phonon frequencies.

fixed-length vector, termed as zero padding, which shows limited accuracy even in a well-trained model (Extended Data Fig. 1a). VGNN solves this issue by rigorously handling the output dimensions for each graph system (Extended Data Fig. 1b). Extended Data Fig. 1c reveals the difference in these performances, where the zero-padding model without virtual nodes (NOVGNN) predicted only the bands of low frequencies (red), while MVN is able to handle the entire frequencies (green) and decently matching with the ground-truth values (black). More detail and results of the tests are in Supplementary Information section 4.4.

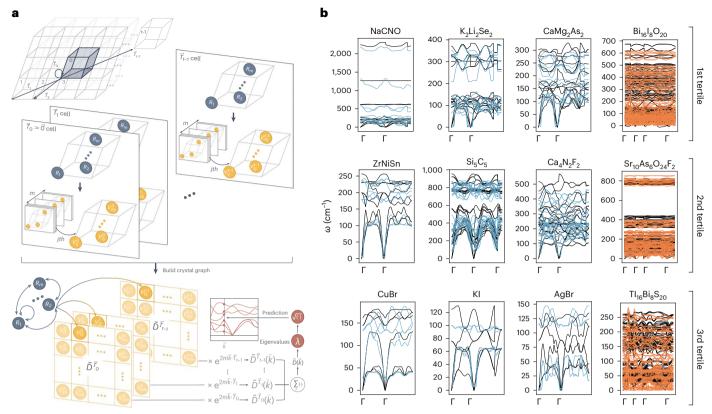
Given that MLIPs are leading approaches for machine-learningdriven phonon prediction, we benchmark VGNN against MLIP in terms of both accuracy and efficiency. M3GNet<sup>31</sup> is one of the state-of-the-art MLIPs predicting interatomic potentials. Though there exist various other MLIP models, such as MACE<sup>32</sup>, ALIGNN<sup>33</sup>, CHGNet<sup>34</sup> and so on, with varying degree of performance<sup>35</sup>, we decided to use M3GNet as the representative MLIP for benchmarking. When comparing a pretrained VGNN model with MLIP, VGNN demonstrates a systematic reduction in run time by two to three orders of magnitude, which can be seen either as a function of number of atoms per unit cell (Fig. 5a) or as a function of chemical element types (Fig. 5b). This improvement can be attributed to VGNN's unique approach of directly inferring dynamical matrix elements, bypassing the traditional MLIP method of calculating forces, second derivatives of interatomic potentials and subsequent Fourier transform. This capability is especially advantageous for materials characterized by large unit cells.

Moreover, it is worthwhile mentioning that VGNN reaches high efficiency without loss of accuracy, comparable to or even slightly

outperforming MLIP and other machine-learning methods. We examine *k*-MVN by comparing prediction accuracy of full phonon band structures with M3GNet, and prediction of phonon DoS with Mat2Spec<sup>36</sup> and E3NN<sup>37</sup>, other leading methods predicting only phonon DoS. We use the phonon band structure, phonon DoS and heat capacity as key metrics to compare the models. For phonon band prediction, the *k*-MVN model outperforms M3GNet with nearly fivefold reduction of error. For phonon DoS prediction, the Mat2Spec and E3NN yielded the smaller prediction error but still comparable to our *k*-MVN method. Finally, all models demonstrate robust performance in predicting heat capacity. Remarkably, our *k*-MVN model gave the smallest error. More detail and results of the tests are further discussed in Supplementary Information section 6.5.

# Application of VGNN method in generating large-scale databases

Today, ab initio calculations such as frozen-phonon and DFPT remain the most accurate methods for phonon calculations. Even so, since the VGNN-based phonon calculation skips the direct calculation of the material-by-material dynamical matrix, it shows substantially faster computation speed while maintaining reasonable accuracy. We can take advantage of this speed in many computationally expensive tasks, such as fast and cheap band structure verification with inelastic scattering experiment (Supplementary Information section 6.1) and high-entropy alloy band structure calculation (Supplementary Information section 6.2). Given the critical role of zeolites in ion exchange, catalysis and gas separations, we use k-MVN to build a zeolite phonon



**Fig. 4** | **The** k-**MVN to predict full phonon band structures. a**, Top: augment  $m^2$  virtual nodes for each translation vector  $\vec{T}$ , and with a total of t neighboring unit cells, a total of  $tm^2$  virtual nodes are generated. Bottom: by multiplying a phase factor by each translated unit cell, a full VDM can be constructed. b, Selected

phonon band structure prediction examples for the high-quality DFPT database (blue) and additional complex materials test (orange) from each error tertile compared with their ground truth (black).  $\Gamma$ -point positions are labeled for each spectrum.

Table 1 | Comparison of how the virtual nodes contribute to phonon prediction in terms of physics and computational costs

	VVN	MVN	k-MVN
Force constants	-	-	Reflected in VDM
Dynamical matrices	-	VDM	VDM
Phonon data	Virtual nodes	Eigenvalues	Eigenvalues
Run time (per phonon wave vector)	O(m²)	O(m <sup>2.37</sup> )	O(m <sup>2.37</sup> )
Storage	O(m)	O(m²)	O(t×m²)
Generalization to larger systems	False	True	True

Here, m and t indicate the number of atoms per unit cell and the number of the unit-cell counts, respectively.

band structure database containing 177 zeolite materials, where the ground truth can be challenging to obtain with ab initio methods. Such a database could support the understanding of phonon-assisted adsorption, catalysis and reaction kinetic processes (Supplementary Information section 6.3). We also conduct further prediction of thermal properties derivable from the phonon band structures, including phonon DoS and heat capacity. With additional ab initio anharmonic force constants calculations in 180 solid-state materials, we demonstrate the capability of directly predicting the temperature-dependent thermal conductivity originating from phonon anharmonicity (Supplementary Information sections 6.4 and 6.7). Finally, by using MVN, we build a database containing the  $\Gamma$ -phonon spectra for over 146,000

 $materials\ listed\ in\ the\ Materials\ Project\ (Supplementary\ Information\ section\ 6.6).$ 

#### **Discussion**

Our proposed VGNN approach offers a versatile framework for predicting material properties with variable dimensions, a capability we have demonstrated through phonon predictions in this study. In general. by taking advantage of the flexibility endowed by virtual nodes, other properties that are challenging to predict for a conventional GNN can be predicted similarly such as electronic band structures, tight-binding and  $k \cdot p$  effective Hamiltonian with a variable number of bands, optical properties such as flexible optical absorption peaks as in the Lorentz oscillator model, and magnetic properties such as the number of propagation vectors. Furthermore, the efficiency of VGNN would enable a different paradigm of the material design. Our VGNN phonon database took a system composed of eight graphics processing units less than 5 h to obtain over 146,000 results, including materials of over 400 atoms per unit cell. The high efficiency of VGNN may further enable materials search and optimization in a broader systems, including alloys, interfaces and even amorphous solids, with superior engineered phonon properties toward thermal storage, energy conversion and harvesting, and superconductivity applications.

Our investigation also reveals some limitations of our current models, specifically for k-MVN in its band structure prediction. We use materials containing light atoms, which generally give high phonon frequencies, and materials with some negative eigenvalues of their dynamical matrices. We also present our additional efforts to solve the issues by fulfilling material data of high phonon frequencies and imposing the restraints of symmetry on the modeled dynamical matrix, respectively, which show improvements in prediction quality in both

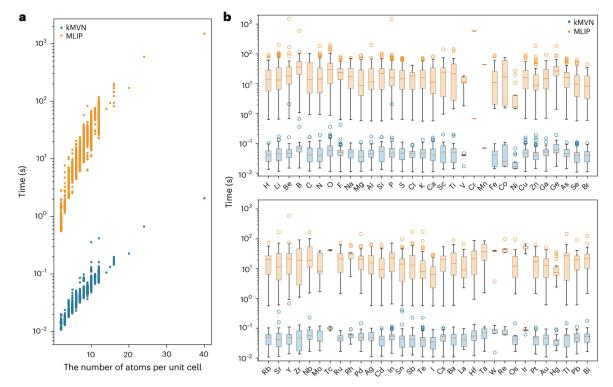


Fig. 5 | The comparison computation time run time between VGNN and MLIP. a, The time spent for each input material (seconds) is plotted in a logarithmic scale for all the 1,500 materials in a DFPT database  $^{28}$ . We present the results of  $\emph{k}$ -MVN (blue) and MLIP (orange). b, The element-wise box plot shows the computation time per material containing each element. The boxes indicate from the first quartile (Q1) to the third quartile (Q3). A horizontal red line goes

through the box at the median. The lower whisker extends from Q1 to the smallest data point within 1.5 times the interquartile range (IQR) from Q1. The upper whisker extends from Q3 to the largest data point within 1.5 times IQR from Q3. The outliers (yellow and blue circles) are the data points that fall outside the range defined by the whiskers.

cases (more detail on these limitations is described in Supplementary Information sections 5.4–5.5). Furthermore, since k-MVN is inspired by the dynamical matrix concept, the model cannot capture long-range interaction contribution to the band structure (Supplementary Information section 5.3). This can potentially be improved by assigning additional virtual nodes tailored to learn long-range corrections.

We can see that these limitations can be addressed either with a proper training data procurement or a better design of virtual node adding scheme emphasizing the flexibility in virtual node method. Moreover, with a potential for incorporating physical information into the machine-learning design like in MVN and k-MVN, virtual node method can be utilized for physically interpretable machine-learning models. Together, the VGNN opens the possibility for GNNs to be used more broadly for predicting more diverse material properties and to design and optimize material in a high-throughput manner.

### Methods

## Phonon data preparation

We trained all of our models against an ab initio DFPT computational database for phonon dispersion in harmonic model<sup>28</sup>. The dataset contains material structures (the same as the primitive structure obtained from the Materials Project<sup>38</sup>), second-order derivatives of energies with respect to atomic perturbations for regular points inside the irreducible zone, and phonon dispersion along highly symmetric paths of 1,521 crystalline inorganic materials. These materials have 2–40 atoms per unit cell, with an average of 7.38. For this work, we only used the highly symmetric path phonon dispersion as our training data. The dispersion is between wave vectors  $\vec{k}$  in the fractional reciprocal unit and response spectra in cm<sup>-1</sup>. All models randomly split the data into 90% training (1,365 materials) and 10% testing (156 materials) sets. Furthermore, we trained our models with a five-fold cross-validation scheme.

We also obtained phonon dispersion of complex (more number of atoms per unit cell) materials from Atsushi Togo's phonon database<sup>29</sup>. We used seekpath<sup>39,40</sup> module to get the highly symmetric path of each material. Then, we fed it alongside POSCAR, FORCE\_SET and phonopy. config files from the database to Phonopy's<sup>41</sup> Python command to calculate the phonon dispersion along such path. To quality control the data, we selected materials whose lowest  $\Gamma$ -phonon band is higher than  $-0.07~\rm cm^{-1}$ . We also only selected materials with more than  $40~\rm atoms$  per unit cell. Finally, we randomly selected 156 (the same as the number of data in the testing set for ease of comparison) out of 505 filtered materials. We used them as our complex material dataset. Further profiles of the dataset are described in Supplementary Information section 1.

#### **Computation environments**

We implemented the models in Python 3.9.13 and trained them on our graphics processing unit cluster with CUDA version 10.2. To facilitate the model implementation and training, we used some important Python modules: Pymatgen<sup>42</sup> and ase<sup>43</sup> for handling material structure files (.CIF), PyTorch<sup>44</sup> for managing the model training framework and e3nn<sup>23</sup> for implementing our neural network models in a form that is equivariant for the Euclidean group.

#### **VGNN**

We have developed a scheme for a GNN for it to be able to have variable output dimensions depending on the input size. For ease of understanding, we will explain the method with our work on phonon prediction.

Considering a material with m atoms per unit cell, we add n additional virtual atoms. We can adjust the number n depending on the model architecture. Using both real and virtual atoms, we convert the crystal structures into periodic graphs with m real nodes for the actual

atoms and n virtual nodes for the added virtual atoms. Then, we connect nodes with edges indicating the message-passing process. To preserve the structural information of the materials and limit the computational cost, we apply the following rules for connections. First, if the distance between the two real nodes is within a specified cutoff radius  $r_{\text{max}}$ , the real nodes are connected through bidirected edges. We also set up an edge between a real node and a virtual node according to the model description, but this edge is directed from real to virtual nodes. Lastly, we embed the information of radial distance vector, for example,  $\vec{r}_{ab}$  from atom b to a, in the form of radial basis functions and spherical harmonics on the corresponding edge as edge attributes, which represent the distance and the direction of  $\vec{r}_{ab}$ , respectively.

Since each node represents an atom in the unit cell, we embedded the atomic numbers A information as node attributes  $\mathcal A$  by passing one-hot representation vectors of length 118 through an embedding layer. As for the model's input, we embedded the atomic masses Z information as input node features  $\mathcal Z$  by passing the product of atomic mass and one-hot representation of atomic number through an embedding layer. We describe the background embedding atomic number and atomic mass into the model in Supplementary Information section 2. Furthermore, we discuss the evidence of using atomic mass as the node features among all other physical descriptors in Supplementary Information section 5.2.

The constructed graph is then passed through the model message passing that operates on the features with multiple convolutions (Supplementary Information section 2) and gated activation layers  $^{45}$ . After the final layer, which consists of only a convolution (no gated activation), each of the n virtual node features is collected and passed through the postprocessing block, which outputs the 3m predicted phonon branches. The postprocessing block is different and will be explained in detail in the subsequent section of each model.

In this work, we utilized e3nn<sup>23</sup> as the framework for implementing GNN architecture. e3nn is designed specifically for three-dimensional Euclidean data. It utilizes the symmetries of the Euclidean group in dimension 3, which includes rotations, translations and mirroring. By leveraging irreducible representations (irreps), e3nn enables GNNs to efficiently process and learn from complex three-dimensional structures while capturing the underlying symmetries. This enhances the expressiveness and interpretability of GNN models for analyzing and modeling three-dimensional Euclidean data.

The model is optimized by minimizing the prediction error. Note that we used the mean squared error (MSE) as the loss function. Phonon energies are the quantities fed into the MSE loss function. We normalize phonon energy values of the training dataset and the predicted ones by the maximum phonon frequency of each material before applying MSE. The full network structure is provided in Supplementary Information section 2. We explore the setup of training and hyperparameter tuning in Supplementary Information section 3.

#### **VVN** method

VVN is a VGNN we designed for learning to predict  $\Gamma$ -phonon spectra from material structures. Since for a material with m atoms per unit cell, there are 3m phonon bands, one sensible choice of adding virtual atoms is to add 3m virtual atoms, each of which outputs the prediction of one of the bands. Hence, when there are m atoms in the unit cell of crystalline material, we assign the position  $\vec{r}_{V_i}$  of the virtual nodes  $V_i$ ,  $i \in [1, 3m]$  following equation (1). We can set the atomic species of the virtual node as anything, and we use Fe after optimization.

$$\vec{r}_{V_i} = \frac{i-1}{3m} \left( \vec{a} + \vec{b} + \vec{c} \right) \tag{1}$$

Here,  $\vec{a}$ ,  $\vec{b}$  and  $\vec{c}$  indicate the unit-cell vectors of the material. In other words, 3m virtual atoms are placed along the diagonal line from (0,0,0) to  $\vec{a} + \vec{b} + \vec{c}$  with equal spacing. By keeping the distances between the virtual nodes in the real space, it is possible to give position

dependencies to the feature updating process. In that sense, equation (1) can consistently keep virtual nodes away from each other and enables us to use the virtual 3m virtual nodes as the output nodes of the network. To get information from the whole structure, each of the 3m virtual nodes is connected to all of the real nodes via directed edges from real to virtual nodes. After each convolution layer, the virtual node features are passed to a linear layer, converted to a scalar output and sorted on the basis of their magnitudes. The outputted 3m scalars represent the predicted  $\Gamma$ -phonon.

#### MVN method

MVN is a VGNN we designed with the influence of the dynamic matrix representation of a periodic harmonic system for learning to predict  $\Gamma$ -phonon spectra from material structures. Given the momentum vector  $\vec{k}$ , the dynamical matrix element  $D_{ij}(\vec{k})$ , which is a three-by-three matrix representing 3D harmonic interaction between atom  $R_i$  and  $R_j$ , can be written as the Fourier transform of the force constant matrix  $\Phi^{\alpha\beta}_{ij}$  following equation (2). Here,  $Z_{R_i}$  is  $R_i$  atom's atomic mass, and  $\vec{r}_{\alpha}$  is the  $\alpha$ th unit-cell position. Note that, for each k vector, the system has 3m degrees of freedom and frequencies where m is the number of atoms per unit cell. We can get the phonon dispersion relations  $\omega(\vec{k})$  by solving eigenvalues  $\omega^2(\vec{k})$  of  $\vec{D}(\vec{k})$ , which is a matrix with shape (3m,3m) that is composed of  $m^2$  blocks of  $\vec{D}_{ij}(\vec{k})$  for  $i,j \in [1,m]$ ,

$$\tilde{D}_{ij}(\vec{k}) = \sum_{\alpha,\beta} \frac{\Phi_{ij}^{\alpha\beta}}{\sqrt{Z_{R_i} Z_{R_j}}} e^{i\vec{k} \cdot (\vec{T}_{\alpha} - \vec{T}_{\beta})}.$$
 (2)

In the MVN method, we generate a matrix that could work like a dynamical matrix as is written in equation (2). Here, we focus on the prediction of Γ-phonon, that is  $\vec{k} = \vec{0}$ . So, the contributions of the same atom pair, for example,  $R_i$ , and  $R_i$  from every unit-cell separation  $\vec{T}_{\alpha} - \vec{T}_{\beta}$ are summed without the  $\vec{k}$ -dependent exponential phase factor. Hence, the model needs to predict a matrix with shape (3m, 3m) representing such summation. To do that, while preserving the relation of each matrix element, we generate m virtual crystals  $C_i$ ,  $j \in [1, m]$  each of which has m virtual nodes  $V_{ii}$ ,  $i \in [1, m]$  of the same atomic species and at the same positions as the real atoms  $R_i$ ,  $i \in [1, m]$ . Here, a virtual node  $V_i$  represents the interaction term from a real node  $R_i$  to another real node  $R_i$  by adding a directed edge from  $R_i$  to  $V_{ii}$  whenever there is an edge connecting  $R_i$  to  $R_i$ . After each convolution layer, the virtual node features are passed to a linear layer and converted to complex-valued output vectors with length 9. For each output feature, we reshape the output features into three-by-three matrices and arrange them such that  $V_{ii}$ 's matrix is the (i, j) block of  $\tilde{D}$  supermatrix with shape (3m, 3m). The method to convert the irreps as a vectors with length 9 into three-by-three Cartesian matrices for constructing  $\tilde{D}$  supermatrix is described in Supplementary Information section 2. Finally, we solve  $\tilde{D}$  for its 3m eigenvalues, which work as the square of  $\Gamma$ -phonon prediction.

#### k-MVN

k-MVN is a generalization of MVN model with nonzero  $\vec{k}$ . Unlike the MVN case, the k-MVN model needs to predict matrices representing interactions between atoms from a unit cell, for example,  $\vec{T}_{\beta}$ , to the different unit cell, for example,  $\vec{T}_{\alpha}$ . Since the phase factor depends only on the difference in unit-cell positions, we can redefine  $\vec{T}$  to be such a difference and simplify equation (2) into

$$\tilde{D}_{ij}(\vec{k}) = \sum_{\vec{T}} \frac{\Phi^{\vec{T}}_{ij}}{\sqrt{Z_{R_i} Z_{R_i}}} e^{i\vec{k}\cdot\vec{T}} := \sum_{\vec{T}} D^{\vec{T}}_{ij} e^{i\vec{k}\cdot\vec{T}}.$$
 (3)

With this simplification, for each  $\vec{T}$ , we generate m virtual crystal  $\vec{C_{j\in[1,m]}}$  the same way as in MVN. However, in this case,  $\vec{V_{ij}}$  represents the interaction term from a real node  $R_j$  to another real node  $R_i$  that is in

the unit cell with unit-cell position  $\vec{T}$  with respect to  $R_j$ 's. In other words, we add a directed edge from  $R_j$  to  $V_{ij}^{\vec{T}}$  whenever there is an edge connecting  $R_j$  to  $R_i$  and that edge represents  $\vec{r}_i - \vec{r}_j = \vec{r}_i' + \vec{T} - \vec{r}_j'$ . Here,  $\vec{r}'$  is the atomic position relative to its unit cell. Since GNN only considers edges with interatomic distance less than  $r_{\text{max}}$ , the model can generate, with this scheme, a nonzero matrix for only a finite number of  $\vec{T}$  that satisfy

$$\min_{i,j \in [1,m]} |\vec{r}_i' + \vec{T} - \vec{r}_j'| \le r_{\text{max}}. \tag{4}$$

Hence, before the virtual crystal generations, the model also iterates through atom pairs to find all viable  $\vec{\tau}$ .

Similar to the MVN model, we convert virtual node features into three-by-three matrices and merge them into a matrix with shape (3m, 3m) representing  $\vec{D}^{\vec{l}}$  for each  $\vec{T}$ . Finally, we weight sum these matrices with their phase factor to get  $\vec{D}$  and solve for its 3m eigenvalues as a phonon spectrum at wave vector  $\vec{k}$ .

# **Data availability**

Source data are provided with this paper. The Γ-phonon database generated with the MVN method, the zeolite phonon band structure database generated with the *k*-MVN method, and the anharmonic phonon calculation data are available at Open Science Framework (OSF) at https://doi.org/10.17605/OSF.IO/K5UTB (ref. 46). We also put the training dataset and source data files generated from models' training and testing, which we used for all analyses, at the same location.

#### **Code availability**

The source code is available at Zenodo (https://doi.org/10.5281/zenodo.8028365)<sup>47</sup>. The GitHub repository presents the instructions for reproducing the results of our simulations and machine learning (https://github.com/RyotaroOKabe/phonon\_prediction).

#### References

- Oganov, A. R. & Glass, C. W. Crystal structure prediction using ab initio evolutionary techniques: principles and applications. J. Chem. Phys. 124, 244704 (2006).
- Le, T., Epa, V. C., Burden, F. R. & Winkler, D. A. Quantitative structure-property relationship modeling of diverse materials properties. Chem. Rev. 112, 2889–2919 (2012).
- Cheng, Y. & Ma, E. Atomic-level structure and structure– property relationship in metallic glasses. *Progress Mater. Sci.* 56, 379–473 (2011).
- Mishra, A., Fischer, M. K. & Bäuerle, P. Metal-free organic dyes for dye-sensitized solar cells: from structure:property relationships to design rules. *Angew. Chem. Int. Ed.* 48, 2474–2499 (2009).
- Dresselhaus, M. S. et al. New directions for low-dimensional thermoelectric materials. Adv. Mater. 19, 1043–1053 (2007).
- Liu, Z. et al. Antiferroelectrics for energy storage applications: a review. Adv. Mater. Technol. 3, 1800111 (2018).
- Zheng, W. & Lee, L. Y. S. Metal-organic frameworks for electrocatalysis: catalyst or precatalyst? ACS Energy Lett. 6, 2838–2843 (2021).
- Jancar, J. et al. Current issues in research on structure-property relationships in polymer nanocomposites. *Polymer* 51, 3321–3343 (2010).
- Kumar, N., Guin, S. N., Manna, K., Shekhar, C. & Felser, C. Topological quantum materials from the viewpoint of chemistry. Chem. Rev. 121, 2780–2815 (2020).
- Oganov, A. R., Pickard, C. J., Zhu, Q. & Needs, R. J. Structure prediction drives materials discovery. *Nat. Rev. Mater.* 4, 331–348 (2019).
- 11. Zhu, T. et al. Charting lattice thermal conductivity for inorganic crystals and discovering rare earth chalcogenides for thermoelectrics. *Energy Environ*. Sci. **14**, 3559–3566 (2021).

- Dunn, A., Wang, Q., Ganose, A., Dopp, D. & Jain, A. Benchmarking materials property prediction methods: the matbench test set and automatminer reference algorithm. npj Comput. Mater. 6, 1–10 (2020).
- Peng, J. et al. Human and machine-centred designs of molecules and materials for sustainability and decarbonization. *Nat. Rev. Mater.* https://doi.org/10.1038/s41578-022-00466-5 (2022).
- 14. Altintas, C., Altundal, O. F., Keskin, S. & Yildirim, R. Machine learning meets with metal organic frameworks for gas storage and separation. *J. Chem. Inform. Model.* **61**, 2131–2146 (2021).
- 15. Schwalbe-Koda, D. et al. A priori control of zeolite phase competition and intergrowth with high-throughput simulations. Science **374**, 308–315 (2021).
- Yao, Y. et al. High-entropy nanoparticles: synthesis-structureproperty relationships and data-driven discovery. Science 376, eabn3103 (2022).
- 17. Hart, G. L., Mueller, T., Toher, C. & Curtarolo, S. Machine learning for alloys. *Nat. Rev. Mater.* **6**, 730–755 (2021).
- Wagih, M., Larsen, P. M. & Schuh, C. A. Learning grain boundary segregation energy spectra in polycrystals. *Nat. Commun.* 11, 1–9 (2020).
- Guo, K., Yang, Z., Yu, C.-H. & Buehler, M. J. Artificial intelligence and machine learning in design of mechanical materials. *Mater. Horizons* 8, 1153–1172 (2021).
- Stanev, V., Choudhary, K., Kusne, A. G., Paglione, J. & Takeuchi, I. Artificial intelligence for search and discovery of quantum materials. Commun. Mater. 2, 1–11 (2021).
- Fung, V., Zhang, J., Juarez, E. & Sumpter, B. G. Benchmarking graph neural networks for materials chemistry. npj Comput. Mater. 7, 1–8 (2021).
- 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).
- 23. Geiger, M. and Smidt, T. e3nn: Euclidean neural networks. Preprint at https://arxiv.org/abs/2207.09453v1 (2022).
- 24. 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).
- 25. Delaire, O. et al. Phonon density of states and heat capacity of La<sub>2-v</sub>Te<sub>4</sub>. Phys. Rev. B **80**, 184302 (2009).
- Baroni, S., Gironcoli, S. & Corso, A. Phonons and related crystal properties from density-functional perturbation theory. Rev. Mod. Phys. 73, 515 (2001).
- Kong, L. T. Phonon dispersion measured directly from molecular dynamics simulations. Comput. Phys. Commun. 182, 2201–2207 (2011).
- 28. Petretto, G. et al. High-throughput density-functional perturbation theory phonons for inorganic materials. *Sci. Data* **5**, 1–12 (2018).
- Togo, A. Phonon database at Kyoto University. Kyoto University https://github.com/atztogo/phonondb/tree/main (2015).
- 30. Sham, L. Electronic contribution to lattice dynamics in insulating crystals. *Phys. Rev.* **188**, 1431 (1969).
- Chen, C. & Ong, S. P. A universal graph deep learning interatomic potential for the periodic table. Nat. Comput. Sci. 2, 718–728 (2022).
- Batatia, I., Kovacs, D. P., Simm, G., Ortner, C. & Csanyi, G. MACE: higher order equivariant message passing neural networks for fast and accurate force fields. In 36th Conference on Neural Information Processing Systems (eds Koyejo, S. et al.) 11423–11436 (Curran Associates, 2022).
- 33. Choudhary, K. & DeCost, B. Atomistic Line Graph Neural Network for improved materials property predictions. *npj Comput. Mater.* **7**, 185 (2021).
- Deng, B. et al. CHGNet as a pretrained universal neural network potential for charge-informed atomistic modelling. *Nat. Mach. Intell.* 5, 1031–1041 (2023).

- Yu, H., Giantomassi, M., Materzanini, G., Wang, J. & Rignanese, G. M. Systematic assessment of various universal machine-learning interatomic potentials. Preprint at https://arxiv.org/ abs/2403.05729v2 (2024).
- 36. Kong, S. et al. Density of states prediction for materials discovery via contrastive learning from probabilistic embeddings. *Nat. Commun.* **13**, 949 (2022).
- 37. Chen, Z. et al. Direct prediction of phonon density of states with Euclidean neural networks. *Adv. Sci.* **8**, 2004214 (2021).
- 38. Jain, A. et al. Commentary: The Materials Project: a materials genome approach to accelerating materials innovation. *APL Mater.* **1**, 011002 (2013).
- 39. Hinuma, Y., Pizzi, G., Kumagai, Y., Oba, F. & Tanaka, I. Band structure diagram paths based on crystallography. *Comput. Mater. Sci.* **128**, 140–184 (2017).
- 40. Togo, A. & Tanaka, I. spglib: a software library for crystal symmetry search. Preprint at https://arxiv.org/abs/1808.01590v2 (2018).
- 41. Togo, A. & Tanaka, I. First principles phonon calculations in materials science. *Scripta Mater.* **108**, 1–5 (2015).
- 42. Ong, S. P. et al. Python materials genomics (pymatgen): a robust, open-source Python library for materials analysis. *Comput. Mater. Sci.* **68**, 314–319 (2013).
- 43. Mortensen, J. et al. The atomic simulation environment—a Python library for working with atoms. *J. Phys. Condens. Matter* **29**, 273002–273002 (2017).
- Paszke, A. et al. PyTorch: an imperative style, high-performance deep learning library. Adv. Neural Inf. Process. Syst. 32, 8024–8035 (2019).
- Miller, B. K., Geiger, M., Smidt, T. E. & Noé, F. Relevance of rotationally equivariant convolutions for predicting molecular properties. Preprint at arXiv https://arxiv.org/abs/2008.08461v4 (2020).
- Okabe, R. & Chotrattanapituk, A. Virtual node graph neural network for full phonon prediction. OSF https://doi.org/10.17605/ OSF.IO/K5UTB (2024).
- Okabe, R. & Chotrattanapituk, A. Virtual node graph neural network for full phonon prediction. *Zenodo* https://doi. org/10.5281/zenodo.8028365 (2024).

#### **Acknowledgements**

R.O., A.C., A.B. and M.L. thank M. Geiger, S. Fang, T. Smidt, K. Persson and S. Yip for helpful discussions and acknowledge the support from the US Department of Energy (DOE), Office of Science (SC), Basic Energy Sciences (BES), award no. DE-SC0021940, and National Science Foundation (NSF) Designing Materials to Revolutionize and Engineer our Future (DMREF) Program with award no. DMR-2118448. A.B. is partially supported by NSF ITE-2345084. R.O. acknowledges support from Heiwa Nakajima Foundation. B.L. acknowledges the support of NSF DMREF with award no. DMR-2118523. T.N., N.D. and M.L. are partially supported by DOE BES award no. DE-SC0020148. T.N. acknowledges the support from Mathworks Fellowship and Sow-Hsin Chen Fellowship. Q.S. acknowledges the support from the Harvard Quantum Initiative. Y.C. is partially supported by the Artificial Intelligence Initiative as part of the Laboratory Directed Research

and Development (LDRD) program of Oak Ridge National Laboratory (ORNL), managed by UT-Battelle, LLC, for the US Department of Energy under contract DE-ACO5-00OR22725. Computing resources for a portion of the work were made available through the VirtuES project, funded by the LDRD Program and Compute and Data Environment for Science (CADES) at ORNL. Another portion of simulation results were obtained using the Frontera computing system at the Texas Advanced Computing Center. M.L. acknowledges the support from the Class of 1947 Career Development Chair and the support from R. Wachnik.

#### **Author contributions**

R.O. and A.C. contributed equally to this work. R.O., A.C. and M.L. conceived the work. R.O. and A.C. set up the framework models, with support from A.B., N.A. and M.L. R.O., A.C. and A.B. performed the machine-learning tests with support from X.F., T.S.J., Q.S., T.N., N.D. and M.L. R.O., A.C. and Y.C. performed high-throughput computations with support from S.M., Y.W. and B.L. R.O. built the databases in the work with support from Y.C. M.L. supervised the project. R.O. wrote the manuscript with revisions from M.L., A.C., Y.C. and input from all authors.

# **Competing interests**

The authors declare no competing interests.

#### **Additional information**

**Extended data** is available for this paper at https://doi.org/10.1038/s43588-024-00661-0.

**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s43588-024-00661-0.

**Correspondence and requests for materials** should be addressed to Ryotaro Okabe, Yongqiang Cheng or Mingda Li.

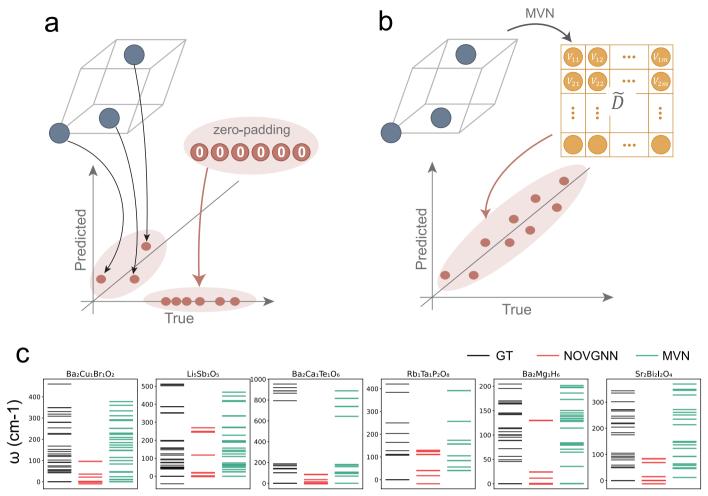
**Peer review information** *Nature Computational Science* thanks Gian-Marco Rignanese for their contribution to the peer review of this work. Primary Handling Editor: Jie Pan, in collaboration with the *Nature Computational Science* team. Peer reviewer reports are available.

**Reprints and permissions information** is available at www.nature.com/reprints.

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.

@ The Author(s), under exclusive licence to Springer Nature America, Inc. 2024



**Extended Data Fig. 1**| **The role of virtual nodes for phonon prediction with flexibility in dimensions. a.** NOVGNN: GNN without virtual nodes. It needs zero-filling to adjust the output dimension to 3m. **b.** Phonon prediction with MVN method. **c.** Comparative plot of  $\Gamma$ -phonon prediction with NOVGNN (red) and MVN (green).