

ARTICLE OPEN



Benchmarking graph neural networks for materials chemistry

Victor Fung¹ , Jiaxin Zhang², Eric Juarez¹  and Bobby G. Sumpter¹ 

Graph neural networks (GNNs) have received intense interest as a rapidly expanding class of machine learning models remarkably well-suited for materials applications. To date, a number of successful GNNs have been proposed and demonstrated for systems ranging from crystal stability to electronic property prediction and to surface chemistry and heterogeneous catalysis. However, a consistent benchmark of these models remains lacking, hindering the development and consistent evaluation of new models in the materials field. Here, we present a workflow and testing platform, MatDeepLearn, for quickly and reproducibly assessing and comparing GNNs and other machine learning models. We use this platform to optimize and evaluate a selection of top performing GNNs on several representative datasets in computational materials chemistry. From our investigations we note the importance of hyperparameter selection and find roughly similar performances for the top models once optimized. We identify several strengths in GNNs over conventional models in cases with compositionally diverse datasets and in its overall flexibility with respect to inputs, due to learned rather than defined representations. Meanwhile several weaknesses of GNNs are also observed including high data requirements, and suggestions for further improvement for applications in materials chemistry are discussed.

npj Computational Materials (2021)7:84; <https://doi.org/10.1038/s41524-021-00554-0>

INTRODUCTION

In the search for materials with various functional applications ranging from catalysis to energy storage to electronics, machine learning (ML) has quickly gained traction as a powerful and flexible approach, especially where a broad exploration of the materials space is needed^{1–6}. The adoption of ML for materials discovery is expected to expand even further with the ongoing growth in the availability of high-throughput density functional theory (DFT) datasets and continued advancements in ML algorithms^{7–14}. Conventionally, ML models in materials chemistry are descriptor-based, where the key descriptors representing the system must first be defined prior to fitting a suitable ML model for prediction. General examples of these descriptors include stoichiometry, the elemental properties such as group, period, electronegativity and radius, and electronic properties such as partial charges and s, p, d-band positions. A number of structural descriptors have also been proposed satisfying translation and rotational invariance, including but not limited to the Coulomb matrix¹⁵, atom-centered symmetry functions (ACSFs)¹⁶, and smooth overlap of atomic positions (SOAP)¹⁷. However, finding effective descriptors can prove challenging for problems with a large amount of compositionally and structurally diverse materials.

In recent years, graph neural networks (GNNs)^{18–20} have received increasing attention as a method that could potentially overcome the limitations of static descriptors by learning the representations on flexible, graph-based inputs. Within this overarching class of ML method, a number of GNN models have been proposed for chemistry-related problems, with the earliest adopters focusing on molecular systems^{21–24}. Subsequently, GNNs have also been used in materials prediction, with a number of studies tackling systems such as periodic crystals^{25–30} and surfaces^{11,31–34}. These systems are generally described by their atomic structures, where the atoms can be represented by nodes and the neighbors encoded by the edges. Information regarding the atoms and bonds such as element type and bond distances, respectively, can be further encoded in the node and edge

attributes. GNNs operate on these atom-based graphs to create node-level embeddings through convolutions with neighboring nodes and edges. This differs from static or pre-defined descriptors which are obtained via a dictionary lookup or a fixed function with no trainable parameters, whereas in a GNN embeddings are obtained from a trainable neural network³⁵.

Given the rapid advances in GNNs for computational materials chemistry currently, a critical evaluation of the current state-of-the-art (SOTA) is warranted. To accomplish this, several criteria should be met: (1) the same datasets should be used across the evaluated models, (2) the datasets used should represent diverse problems in materials chemistry, (3) the same input information and representation should be used, (4) the hyperparameters of the models should be optimized to the same extent, and (5) these should be performed in a reproducible manner. In this work, we attempt to address these criteria and provide an open-source framework, MatDeepLearn, which can be used for further studies in this area (Fig. 1)³⁶. Provided with atomic structures and target properties from any dataset of choice, MatDeepLearn handles the processing of structures to graphs, offers a library of GNN models, and provides hyperparameter optimization for the models. Within this framework, improvements and additions to the input representations and model architectures can be easily made, which can greatly reduce the development time needed for new ML models. GNNs can also be critically evaluated, and actionable information can be quickly obtained when applied to specific problems in chemistry and the materials sciences. We then use this framework to benchmark several SOTA models and provide a timely snapshot of current progress and suggestions for future progress.

RESULTS

Evaluation of performance across datasets

We evaluated a total of seven ML models for regression tasks on the five datasets, summarized in Table 1, with the specific

¹Center for Nanophase Materials Sciences, Oak Ridge National Laboratory, Oak Ridge, TN, United States. ²Computer Science and Mathematics Division, Oak Ridge National Laboratory, Oak Ridge, TN, United States. ✉email: fungv@ornl.gov

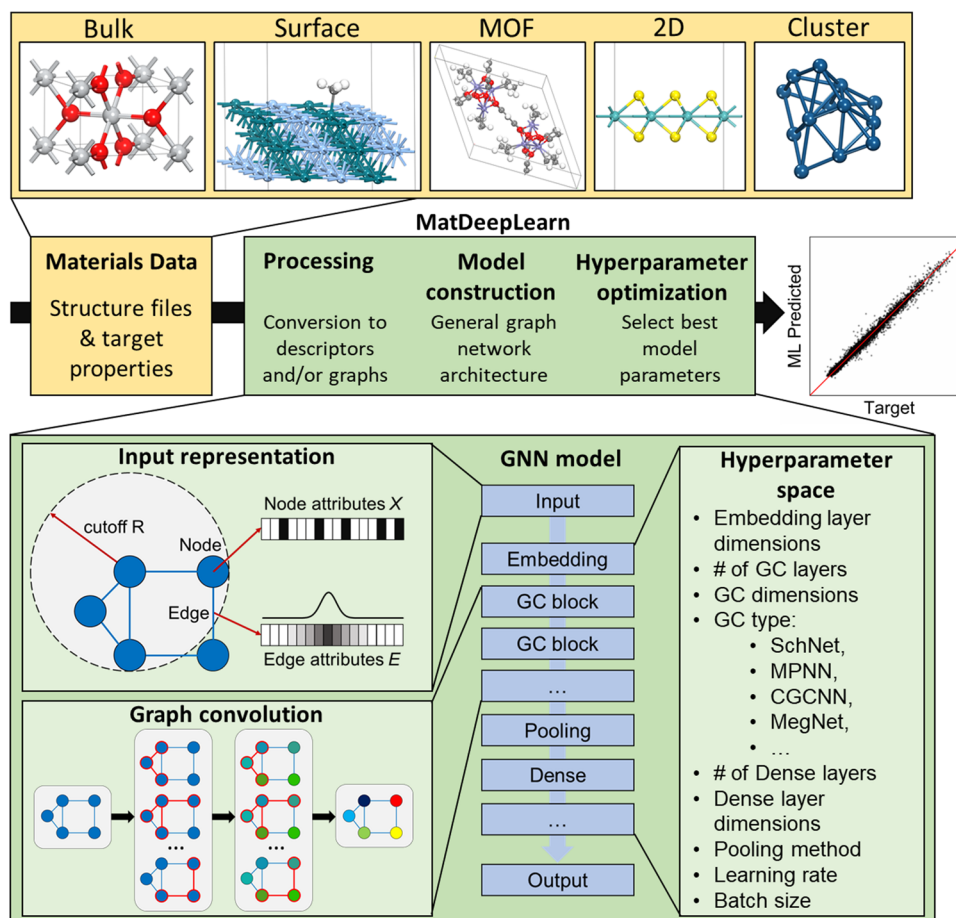


Fig. 1 Scheme of the GNN testing framework and workflow. A general outline of the approach is presented, starting with materials data (with several examples shown) in the form of structure files and target properties, followed by data processing, model construction, and hyperparameter optimization as part of the MatDeepLearn framework. A general graph neural network architecture is constructed, taking in graphs containing nodes, edges, node attributes, and edge attributes, inputted into an embedding layer, GC blocks, pooling, and dense layers. This allows models to be similarly compared within a shared hyperparameter space.

Table 1. Benchmarking results—models.

ML models	Mean Absolute Error (MAE)				
	Datasets				
	Bulk crystals (eV/Atom)	Alloy surfaces (eV)	MOFs (eV)	2D materials (eV)	Pt clusters (eV)
SchNet	0.050	0.063	0.228	0.214	0.151
MPNN	0.046	0.058	0.245	0.204	0.182
CGCNN	0.049	0.060	0.233	0.208	0.205
MEGNet	0.048	0.069	0.253	0.224	0.180
GCN	0.067	0.175	0.355	0.304	0.577
SOAP	0.047	0.118	0.318	0.203	0.143
SM	0.394	0.621	0.608	0.607	0.460
Baseline	0.978	1.480	0.984	0.773	4.984

hyperparameters for each top model listed in Supplementary Table 1, with GNN entries shaded in gray. The MAE is obtained through five-fold cross-validation over the entire dataset. We find all four SOTA models (SchNet, MPNN, CGCNN, MEGNet) performed very well for all tested datasets, once again demonstrating the capability of GNNs for accurate predictions when ample training data is available. Unexpectedly, we find the SOTA models all performed equally well in most cases, which shows

hyperparameter optimization can be just as important as the choice of graph convolutional operator and overall machine learning architecture. The SOTA models also offer advantages over simpler GNNs that were not designed for materials chemistry, such as the GCN model, which performed markedly worse for nearly all datasets except the bulk crystals. The particularly poor GCN performance, especially for Pt clusters, suggests convolutions involving edge attributes containing interatomic distances in the

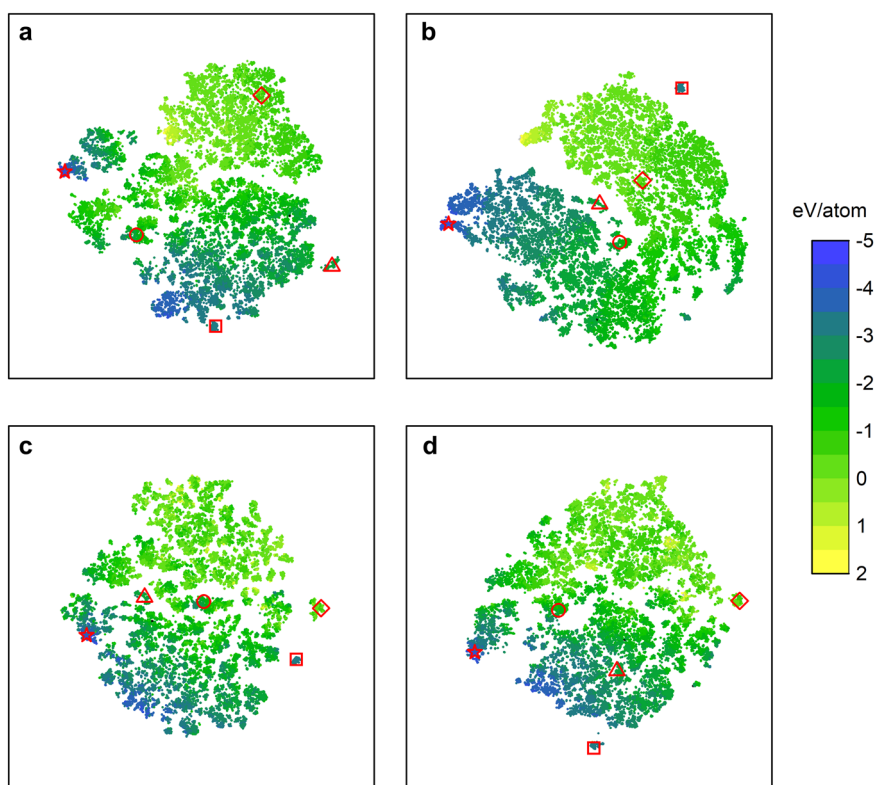


Fig. 2 Visualization of graph-wide feature space. t-distributed stochastic neighbor embedding (t-SNE) plot of the graph-level embedding from the readout/pooling layer for (a) CGCNN, (b) MPNN, (c) SchNet, (d) MEGNet trained on the bulk dataset, with each point representing an individual crystal. Colors for each point are mapped to formation energies. Selected structures are marked by red shapes, with the same shape corresponding to the same crystal in each plot.

Gaussian basis, which the SOTA models share, are much more effective for capturing the spatial information of atomic structures.

Meanwhile, for non-GNN models, SM performed poorly across the datasets, but SOAP performed surprisingly well for bulk crystals and 2D materials and had a performance similar to or better than the SOTA models for Pt clusters. On one hand, this suggests the SOAP descriptors, like other similar descriptors such as ACSFs, remain excellent choices for capturing spatial information in atomic structures, particularly for use in machine-learned potentials. On the other hand, this nevertheless highlights another strength in GNN-based models in that a similar level of performance can still be achieved without any pre-defined descriptors or existing knowledge.

Relative to the Baseline, the SOTA GNNs had the highest relative errors in predicting work functions for the 2D dataset, and this is likely due to the much smaller size of the dataset compared to the other datasets. This is followed by the MOFs dataset, which performed second worst. Next, the GNNs performed similarly well for bulk and surfaces, with both containing ample amounts of data for training. Finally, the GNNs performed best for the clusters relative to the Baseline; however, these errors, at approximately 0.015 eV/atom, remain relatively high for use as machine-learned potentials. A more thorough evaluation would be needed in the future for GNN performance in ML potentials, which may require more significant changes in the model architecture and training data and the inclusion of forces.

Visualization of GNN features

We then compare the ability of the different GNNs in learning the structure and composition representations by obtaining the graph-level embedding from the output of the readout/pooling layer and visualizing with t-distributed stochastic neighbor

embedding (t-SNE) in Fig. 2, using the Bulk dataset as an example. The plots in Fig. 2 represent a combined structure-composition latent space for the trained materials where points within a grouping can be expected to share similarities in both their atomic structures and elemental compositions. Each of the four SOTA GNNs is able to generate viable representations leading to groupings of crystals delineated by similar formation energies. In some cases, each GNN obtained a similar grouping in the latent space, such as for structures labeled with the star or square symbols. In other cases, the GNNs obtained dissimilar grouping, for example with the ones denoted by the diamond, and to a lesser extent the triangle and circle symbols. Ultimately, for purposes of regression each of the plotted latent spaces is equally valid and provides similar prediction errors; though additional investigation may be needed in the future to reveal additional differences for applications such as generative machine learning.

Training size dependence

Next, we investigated the training size dependence of the models, another avenue of comparison between models and datasets, and estimate the maximum performance with respect to data size in Fig. 3, with specific values in Supplementary Table 2. The model performance is obtained from five parallel runs for each training size using a different train/test split and averaging the errors. We find, in almost all situations, the training size dependence to be very similar between the GNN models for a particular dataset, approximately irrespective of the number of parameters in the model. Unsurprisingly, SOAP has a better training size scaling for the cluster dataset, likely thanks to the effectiveness of the pre-defined descriptors. However, when applied to compositionally diverse systems, SOAP loses this advantage over GNNs and has a similar training size scaling with the other models. We also

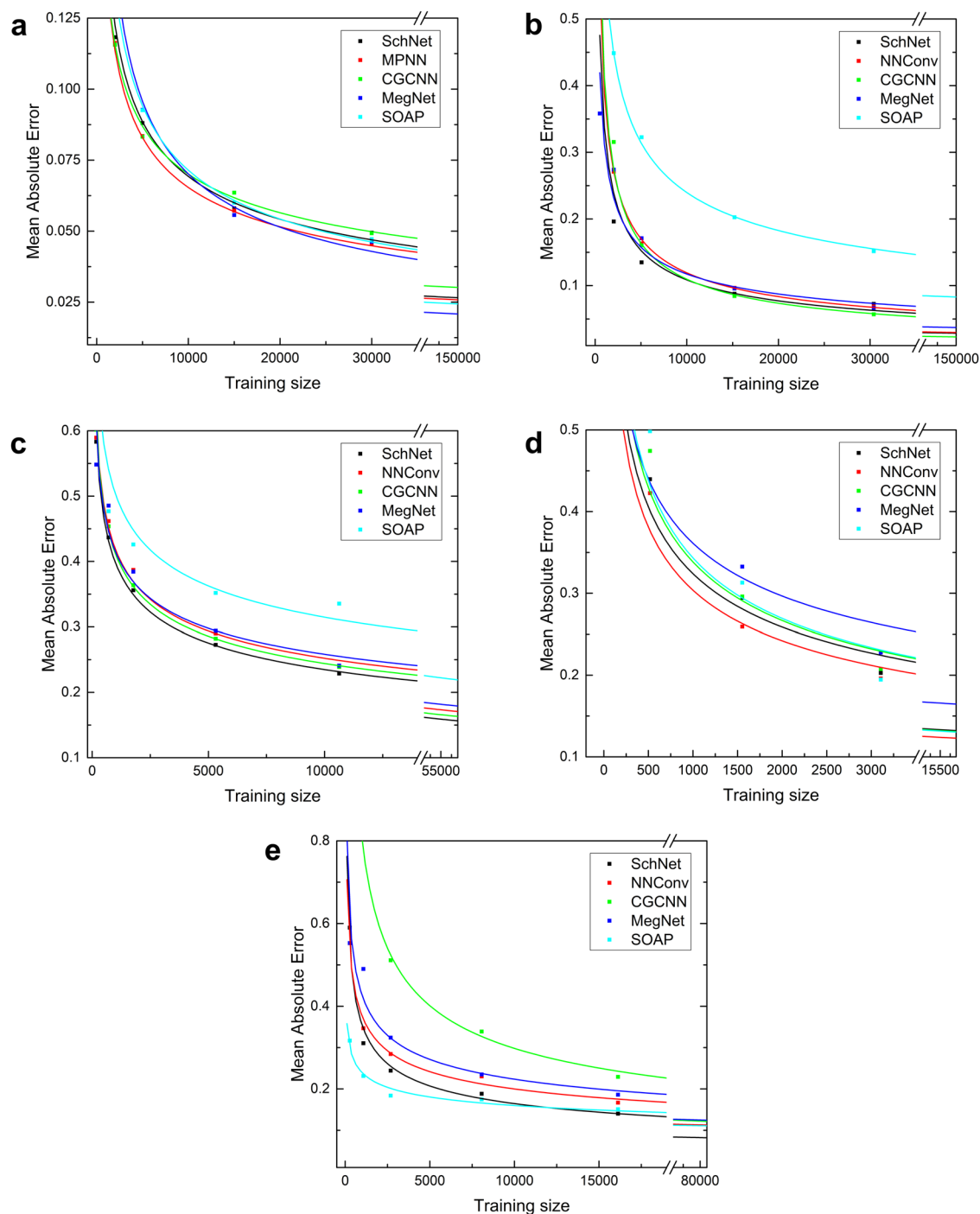


Fig. 3 Training size dependence and extrapolation. The training size dependence for (a) bulk, (b) surface, (c) MOF, (d) 2D, (e), and cluster datasets are plotted. Each point represents the average of five separate runs using different train/test splits. Solid lines are obtained from fitting the power law.

estimate performance on data sizes by extrapolating training size dependence curves using a power-law function, $\epsilon(m) = am^b$ where ϵ is the error as a function of samples, m ³⁷. We extrapolate to five times the current data set size for each case. For the bulk datasets, fitting the power-law function gives exponents of ~ -0.3 for the GNNs, and extrapolating predicts MAEs of ~ 0.02 – 0.03 for 5x data at $\sim 150,000$ training data. For surfaces, power-law fitting suggests a better scaling than bulk with exponents of ~ -0.5 and provides an MAE estimate of ~ 0.03 – 0.04 eV. Meanwhile, a recent study reached an opposite conclusion, finding worse scaling for

surface adsorption compared to the bulk¹¹. This is likely due to unrelaxed structures being used as inputs, which significantly increases the dimensionality and difficulty of the task.

Evaluation of representation sensitivity

So far, we have focused primarily on model performance and maintained a consistent representation for all models and datasets. We examine the soundness of this methodology by testing several different representations and observing their impact on the

Table 2. Benchmarking results—representations.

Input representation	MAE				
	Datasets				
	Bulk crystals (eV/Atom)	Alloy surfaces (eV)	MOFs (eV)	2D materials (eV)	Pt clusters (eV)
Default	0.049	0.060	0.233	0.208	0.205
Neighbors: 4	0.054	0.082	0.238	0.195	0.410
Neighbors: full	0.047	0.060	0.261	0.198	0.203
Node features: element One-hot	0.067	0.058	0.232	0.217	0.184
Node features: blank	0.319	0.570	0.410	0.562	0.187
Edge length: 10	0.050	0.059	0.257	0.218	1.328

performance (Table 2). Here we use the optimal hyperparameters in Supplementary Table 1, and test using the CGCNN model. Possible avenues for modifying the representations include changing the edge cutoffs and selection criteria for generating the graphs and changing the properties of the node and edge attributes. First, we find reducing the number of considered neighbors from 12 to 4 did not significantly change prediction accuracy for most cases besides the cluster dataset where the performance decreases significantly. Meanwhile, a fully connected graph did not improve prediction accuracy beyond the default, while greatly increasing the computational cost. Moving to node attributes, we found using node features based simply on element identity with one-hot encoding was also sufficient, and including additional elemental properties such as electronegativity and radius were generally unnecessary for sufficiently large data sets. We also tested using blank input node attributes (with all zeroes), thereby removing any information regarding the elemental composition; this significantly reduced the performance for all cases except the cluster dataset (which only has one element in the system). Finally, a test was performed where the length of edge attributes was reduced from 50 to 10, thereby greatly reducing the structural resolution encoded by the Gaussian basis. This had little to no impact on 4 out of 5 datasets, but greatly reduced the prediction ability for clusters as the information needed to distinguish small changes in Pt-Pt bond distances is lost. Thus, we find the default representation used in this work to be satisfied within the scope of including only structure and atomic information as inputs. Variations in the representation, up to an extent, did not appreciably affect performance likely due to the property of GNNs in producing learned representations via training.

DISCUSSION

With the rapid expansion of readily accessible high-throughput DFT datasets and new GNN models which can train on them, there is a strong demand for a general benchmarking tool to assess which ML models are best for a given target application in materials chemistry. We have developed such a framework, MatDeepLearn, and used it to gauge the performance of several GNNs on different materials data sets. We use the Pytorch and Pytorch-Geometric libraries which allow for fast ML calculations that are optimized for GPU-based computing resources, and the Ray library which provides distributed hyperparameter optimization on multiple nodes. In this work, training time ranges from ~10 min to ~1–2 h on an NVIDIA Tesla V100 GPU for the smallest and largest datasets, respectively. Hyperparameter training with 160 trials can be finished in roughly 2 days on a single GPU node containing 8 V100 GPUs. In this work, five GNN models were tested, but this list can be rapidly expanded with the provided message-passing network class in Pytorch-Geometric, and with approximately 40 existing methods already implemented for use. Consequently, the development time needed for GNNs can be shortened significantly, along with the rapid testing and

benchmarking of these developed methods. In general, we find MatDeepLearn can provide a highly competitive baseline performance with little to no human interaction or effort required; only a suitable dataset containing atomic structures is needed as the inputs.

Our current study also provides some general observations regarding GNNs for materials applications. For the prediction tasks in this work, a GNN which contains nonlinear update functions for nodes (usually neural networks), and a sufficiently descriptive representation of bond distances generally performs quite well, and differences between the current tested graph convolutional operators become small through hyperparameter optimization. Within the scope of our study, including edge updates did not appear to improve the performance perceptibly. In moving forward, a more exhaustive screening of the GNN design space may be fruitful, such as those proposed in a recent study by You et al.³⁸ Structure-agnostic models have also been developed in recent years which rely only on composition and their related properties^{39,40}, with some approaches also using the GNN architecture⁴¹. For certain applications such as bulk crystal properties, these have shown to reach comparable performance with the methods in this work.

Training the GNNs can be data-intensive, and we find a minimum of 10^3 – 10^4 data points are needed to achieve adequate accuracy in the models. This general sentiment is echoed in a similar study finding descriptor-based models performing better than GNNs with small data sizes, but with GNNs performing better when ample data is available²⁵. For applications with very small datasets, pre-defined descriptors still work best, but will quickly fail when moving outside of its domain of applicability. Methods that can improve the quality of the initial guess structures could help reduce the training costs significantly.

Approaches to effectively incorporate domain knowledge with GNN models which complements their existing flexibility and ability for learnable representations would also likely improve performance. For example, additional information about the system can be incorporated into the node or edge attributes, such as the bond types, aromaticity, and chirality in the case of molecular systems, which is not inherently known from just the atomic structure. Along these lines, recent approaches incorporating some form of features relating to atomic orbitals and their interactions have also yielded promising results^{30,42}. Alternatively, the inclusion of physical constraints in the loss functions or through other means which have been used for ML, in general, could also be considered for the GNNs here.

Additionally, effectively incorporating knowledge from pre-existing datasets through transfer learning is also another promising avenue for improvement by leveraging current high-throughput computational databases. Besides the datasets used in this study, many other computational databases with 10^4 or more data entries are also available which would be well-suited for training with GNNs^{7–11}. GNN methods that can train on multi-

Table 3. Dataset information.

Target property	Datasets				
	Bulk crystals	Alloy surfaces	MOFs	2D materials	Pt clusters
Formation energy (eV/atom)		Adsorption energy (eV)	Band gap (eV)	Work function (eV)	Total energy (eV)
Number of data	Approx. 37000	Approx. 37000	Approx. 13000	Approx. 4000	Approx. 20000
Material size range (atoms)	1 to 200	13 to 16	17 to 150	2 to 12	10 to 13
Material composition range (elements)	87	42	78	60	1
Calculation method	DFT: PBE+U	DFT: BEEF-vdw	DFT: PBE	DFT: PBE+U	DFT: PBE
Source	Materials project	CatHub	QMOF	C2DB	Literature

fidelity data have been demonstrated recently which would improve integration with multiple datasets⁴³.

METHODS

Datasets

The datasets we used were each chosen to reasonably represent a variety of different classes of inorganic materials ranging from 3D to 0D, each with different target properties, summarized in Table 3. 3D materials are periodic in three dimensions and include bulk crystals, while 2D materials refer to solids with the finite thickness (single or few-layer), and 0D materials refer to nanoclusters and nanoparticles which are non-periodic. All five datasets were generated through DFT calculations. Four of the datasets were obtained from curated, open computational databases.

For bulk materials, we compiled approximately 37,000 structures from the Materials Project¹², which contains the widest elemental diversity of the datasets in this work, as well as the greatest diversity in structure size from the number of atoms. Bulk crystals have been used as a de facto benchmark system in the computational materials literature, with many successful examples of regression and classification using both GNN models^{25–30}. The property selected for this dataset is formation energy in units of eV/atoms. We note this database is constantly expanding, and here we used a snapshot of a subset of the available data.

For surfaces, we employed a dataset of another roughly 37,000 structures from the Catalysis-Hub database, containing eleven adsorbates on approximately 2000 unique metal alloy surfaces⁴⁴. For this dataset, we used the relaxed surfaces containing the adsorbate as the input structures and adsorption energy in eV as the target property. Unrelaxed surfaces are not used in this study due to the ambiguities involved in the placement of the adsorbate which will not be explored in detail here. With this dataset, we intend to evaluate the general ability of the model for predicting surface chemistry properties, which are relevant for catalysis.

For porous materials, we used a metal-organic framework (MOF) dataset from the QMOF database containing roughly 13,000 processed structures at the time of access, from experimentally synthesized MOFs and containing band gaps in eV as the target property⁴⁵.

For 2D materials, we used a dataset from C2DB, containing around 4000 structures, with the target property being work function in eV, another important electronic property⁴⁶. This dataset is smaller than the other examples in this work while still being compositionally diverse with 60 elements included, and may prove challenging for GNNs with their generally high data requirements.

The last dataset is compositionally narrow with only one element, Pt, but is structurally diverse with around 20,000 different nanoclusters ranging from 10 to 14 atoms, and with total energies in eV as the target. These clusters were obtained from basin-hopping global optimization in a previous study⁴⁷. This dataset was included to evaluate the ability of GNNs to capture structural sensitivity for potential use in machine-learned potentials, a task for which descriptor-based methods are currently most commonly used.

Structure/graph representations

The data are organized as a set of structures and a set of target properties associated with the structures. The structures are described by a set of atomic positions in Cartesian coordinates and an accompanying lattice vector describing the dimensions of the periodic cell. Non-periodic

structures can also be represented by simply placing the molecule/cluster in a sufficiently large empty cell and ensuring the distance between images are larger than distance cutoffs for edges. Each atom is represented by a node in the graph, and the edges are usually determined by interatomic distances within a certain radius cutoff. The node attributes in this work contain elemental properties of the atoms which are one-hot encoded as described by Xie et al.²⁷. Edge attributes encode the interatomic distances, described later in the section.

Machine learning models and training

The general architecture of overall GNN models used so far in materials chemistry contains several shared characteristics, which we unify into a general architecture here, illustrated in Fig. 1. First, an embedding or preprocessing layer is present which transforms the node attributes from the input to a specified dimension. This is followed by N number of graph convolutional blocks, which perform the convolution and aggregation of the nodes. We used the same graph convolutional operators used in the original GNN models. This is followed by a graph-wide readout/pooling layer which provides an overall graph representation by aggregating the node attributes; in this work, we choose from max, average, sum, and set2set pooling. This is followed by M dense layers and the scalar output. The dimensions of the embedding, graph convolutional, and dense blocks, the type of pooling layer used, layer counts N and M, batch size, and learning rate are hyperparameters that are selected through optimization for each model and each dataset (search ranges in the SI).

We tested five different graph convolutional operators here:

The SchNet²² convolutional operator:

$$x'_i = \sum_{j \in N(i)} x_j \odot h_{\Theta} \left(\exp(-\gamma(d_{ij} - \mu))^2 \right) \quad (1)$$

Here, d_{ij} is defined as the interatomic distance between atoms i and j , and h_{Θ} is a neural network containing dense layers which generate filters from interatomic distances. Before being fed to the neural network, the distances are expanded by a Gaussian basis function, which provides a continuous, non-sparse representation. This was also later successfully applied for use in other GNNs such as CGCNN and MEGNet. In the rest of this work, we define this term as

$$e_{ij} = \exp(-\gamma(d_{ij} - \mu)^2) \quad (2)$$

and use these as the edge attributes.

Additionally, update functions are applied to the node attributes in the form of dense layers. In the original SchNet model, atom features are transformed to scalar atom-wise values before being summed over the whole graph. This is common practice for machine-learned potentials which first calculate the atom-wise energies before sum pooling to obtain the total energy^{48,49}. Instead, in this work we pool first to obtain the graph-level features before predicting the scalar property using dense layers; this is consistent with other approaches for materials property predictions such as MEGNet and CGCNN.

The Message Passing Neural Network (MPNN)²¹ convolutional operator:

$$x'_i = \Theta x_i + \sum_{j \in N(i)} x_j \odot h_{\Theta}(e_{ij}) \quad (3)$$

h_{Θ} is a neural network containing dense layers, and update functions are also applied to the nodes, this time in the form of a gated recurrent unit.

The Crystal Graph Convolutional Neural Network (CGCNN)²⁷ convolutional operator:

$$x'_i = x_i + \sum_{j \in N(i)} \sigma(z_{ij}W_f + b_f) \odot g(z_{ij}W_s + b_s) \quad (4)$$

Here, $z_{ij} = x_i \oplus x_j \oplus e_{ij}$, and σ and g are sigmoid and softplus functions respectively.

The MatErials Graph Network (MEGNet)²⁶ convolutional operator:

$$e'_{ij} = h_{\text{oe}}(x_i \oplus x_j \oplus e_{ij}) \quad (5)$$

$$x'_i = h_{\text{ov}} \left(\left(\frac{1}{N(i)} \sum_{j \in N(i)} e_{ij} \right) \oplus x_i \right) \quad (6)$$

Two dense layers are added at the beginning of each MEGNet graph convolutional block to preprocess the inputs. Here h_{oe} and h_{ov} are edge and node update functions, which are also dense layers. The updates follow the order of edges, nodes, and global attributes. A skip connection adds the unprocessed input attributes of each block with the output attributes. In the original work, the global attributes were left blank for the inorganic crystals dataset and are similarly unused here.

The Graph Convolutional Network (GCN)⁵⁰ convolutional operator:

$$x'_i = \Theta \sum_j \frac{1}{\sqrt{\hat{d}_i \hat{d}_j}} x_j \quad (7)$$

We include this as a baseline graph convolutional with much simpler construction and not specifically developed for materials chemistry applications. In contrast to the earlier models, a purely linear update function Θ is used here for the node attributes, and edge attributes are not included. Instead, the edge weights are used, containing the inverse normalized atomic distances.

For each model and dataset, hyperparameter optimization was performed for 160 trials using the HyperOpt optimizer with an 80–20 split for training and validation, and the model with the lowest validation error was selected⁵¹. The hyperparameter search space can be found in the Supplementary Information. To confirm whether 160 trials were sufficient, two tests were performed for the bulk dataset with 640 trials and found no significant improvement in performance. The reported performance is then obtained from five-fold cross-validation on the selected model. We found the performance may differ from the values in the literature; for example, the error for CGCNN in this study is 0.049 while it is 0.039 in the original paper²⁷. This may be the result of many potential factors, such as data selection, processing, and model optimization. Here, we limited the training to within 200 epochs to make hyperparameter optimization computationally affordable, while a more thoroughly optimized model may have a better performance. Meanwhile, we also note large variations in performance for the same model in the literature^{25,29,30}. This furthermore reinforces the need for a consistent framework for fairly and reproducibly comparing models.

In addition, we tested two non-GNN models, using the Sine matrix⁵² (SM) and the Smooth Overlap of Atomic Positions¹⁷ (SOAP) descriptors, for comparison. For the SM descriptor, the size of the matrix is padded with zeros to the maximum atomic size of the dataset, and only the eigenvalues are used, sorted in descending order of their absolute values. For the SOAP descriptor, Gaussian type orbital basis functions are used and an “inner” average is obtained for each element present in the dataset, whereby the average is taken over the sites before summing up the magnetic quantum numbers. The SOAP parameters, the distance cutoff, the number of radial basis functions, the degree of spherical harmonics, and the standard deviation of the Gaussians in the basis functions are all considered as hyperparameters to be optimized. For both models, the descriptors are fed into dense layers with variable layer count and size as determined from hyperparameter optimization. Finally, an overall baseline is provided for comparison in the form of a dummy regressor which only returns the mean of the training dataset.

The code was written in Python 3.7 and uses PyTorch v1.6 and PyTorch-Geometric⁵³ v1.6 libraries for the ML models³⁶. The Dscribe library was used to obtain SM and SOAP descriptors⁵⁴. We use the Ray library which provides distributed hyperparameter optimization on multiple nodes⁵⁵.

DATA AVAILABILITY

The datasets used for testing are also provided in full or with instructions included at <https://github.com/vxfung/MatDeepLearn>.

CODE AVAILABILITY

The code used in this work, MatDeepLearn, is open-source and available at <https://github.com/vxfung/MatDeepLearn>.

Received: 20 January 2021; Accepted: 14 May 2021;

Published online: 03 June 2021

REFERENCES

- Schmidt, J., Marques, M. R. G., Botti, S. & Marques, M. A. L. Recent advances and applications of machine learning in solid-state materials science. *Npj Comput. Mater.* **5**, 83 (2019).
- Alberi, K. et al. The 2019 materials by design roadmap. *J. Phys. D* **52**, 013001 (2018).
- Chen, C., Zuo, Y., Ye, W., Li, X., Deng, Z. & Ong, S. P. A critical review of machine learning of energy materials. *Adv. Energy Mater.* **10**, 1903242 (2020).
- Batra, R., Song, L. & Ramprasad, R. Emerging materials intelligence ecosystems propelled by machine learning. *Nat. Rev. Mater.* 1–24 (2020).
- 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).
- Schlexer Lamoureux, P. et al. Machine learning for computational heterogeneous catalysis. *ChemCatChem* **11**, 3581–3601 (2019).
- Saal, J. E., Kirklin, S., Aykol, M., Meredig, B. & Wolverton, C. Materials design and discovery with high-throughput density functional theory: the Open Quantum Materials Database (OQMD). *JOM* **65**, 1501–1509 (2013).
- Choudhary, K. et al. The joint automated repository for various integrated simulations (JARVIS) for data-driven materials design. *Npj Comput. Mater.* **6**, 173 (2020).
- Curtarolo, S. et al. AFLOWLIB.ORG: a distributed materials properties repository from high-throughput ab initio calculations. *Comput. Mater. Sci.* **58**, 227–235 (2012).
- Draxl, C. & Scheffler, M. The NOMAD laboratory: from data sharing to artificial intelligence. *J. Phys. Mater.* **2**, 036001 (2019).
- Chanussot, L. et al. The Open Catalyst 2020 (OC20) Dataset and Community Challenges. *ACS Catal.* **11**, 6059–6072 (2021).
- Jain, A. et al. Commentary: the materials project: a materials genome approach to accelerating materials innovation. *APL Mater.* **1**, 011002 (2013).
- Clement, C. L., Kauwe, S. K. & Sparks, T. D. Benchmark AFLOW data sets for machine learning. *Integr. Mater. Manuf. Innov.* **9**, 153–156 (2020).
- Wang, A. Y.-T. et al. Machine learning for materials scientists: an introductory guide toward best practices. *Chem. Mater.* **32**, 4954–4965 (2020).
- Rupp, M., Tkatchenko, A., Müller, K.-R. & von Lilienfeld, O. A. Fast and accurate modeling of molecular atomization energies with machine learning. *Phys. Rev. Lett.* **108**, 058301 (2012).
- Behler, J. Atom-centered symmetry functions for constructing high-dimensional neural network potentials. *J. Chem. Phys.* **134**, 074106 (2011).
- De, S., Bartók, A. P., Csányi, G. & Ceriotti, M. Comparing molecules and solids across structural and alchemical space. *Phys. Chem. Chem. Phys.* **18**, 13754–13769 (2016).
- Battaglia, P. W. et al. Relational inductive biases, deep learning, and graph networks. Preprint at <https://arxiv.org/abs/1806.01261> (2018).
- Zhou, J. et al. Graph neural networks: a review of methods and applications. *AI Open* **1**, 57–81 (2020).
- Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C. & Yu, P. S. A comprehensive survey on graph neural networks. *IEEE Trans. Neural Netw. Learn. Syst.* **32**, 4–24 (2021).
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O. & Dahl, G. E. Neural message passing for quantum chemistry. *ICML* 1263–1272 (2017).
- Schütt, K., Kindermans, P.-J., Sauceda Felix, H. E., Chmiela, S., Tkatchenko, A. & Müller, K.-R. SchNet: a continuous-filter convolutional neural network for modeling quantum interactions. *NeurIPS* **30**, 991–1001 (2017).
- Duvenaud, D. K. et al. Convolutional networks on graphs for learning molecular fingerprints. *NeurIPS* 2224–2232 (2015).
- Wu, Z. et al. MoleculeNet: a benchmark for molecular machine learning. *Chem. Sci.* **9**, 513–530 (2018).
- 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**, 138 (2020).
- Chen, C., Ye, W., Zuo, Y., Zheng, C. & Ong, S. P. Graph networks as a universal machine learning framework for molecules and crystals. *Chem. Mater.* **31**, 3564–3572 (2019).
- 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).

28. Louis, S.-Y. et al. Graph convolutional neural networks with global attention for improved materials property prediction. *Phys. Chem. Chem. Phys.* **22**, 18141–18148 (2020).
29. Park, C. W. & Wolverton, C. Developing an improved crystal graph convolutional neural network framework for accelerated materials discovery. *Phys. Rev. Mater.* **4**, 063801 (2020).
30. Karamad, M., Magar, R., Shi, Y., Siahrostami, S., Gates, I. D., Barati & Farimani, A. Orbital graph convolutional neural network for material property prediction. *Phys. Rev. Mater.* **4**, 093801 (2020).
31. Back, S., Yoon, J., Tian, N., Zhong, W., Tran, K. & Ulissi, Z. W. Convolutional neural network of atomic surface structures to predict binding energies for high-throughput screening of catalysts. *J. Phys. Chem. Lett.* **10**, 4401–4408 (2019).
32. Palizhati, A., Zhong, W., Tran, K., Back, S. & Ulissi, Z. W. Toward predicting intermetallics surface properties with high-throughput DFT and convolutional neural networks. *J. Chem. Inf. Model.* **59**, 4742–4749 (2019).
33. Gu, G. H., Noh, J., Kim, S., Back, S., Ulissi, Z. & Jung, Y. Practical deep-learning representation for fast heterogeneous catalyst screening. *J. Phys. Chem. Lett.* **11**, 3185–3191 (2020).
34. Palizhati, A., Zhong, W., Tran, K., Back, S., Ulissi, Z. W. Toward predicting intermetallics surface properties with high-throughput DFT and convolutional neural networks. *J. Chem. Inf. Model.* **59**, 4742–4749 (2019).
35. Keith, J. A. et al. Combining machine learning and computational chemistry for predictive insights into chemical systems. Preprint at <https://arxiv.org/abs/2102.06321> (2021).
36. MatDeepLearn. <https://github.com/vxfung/MatDeepLearn>. Accessed 1/4/2021 (2021).
37. Hestness, J. et al. Deep learning scaling is predictable, empirically. Preprint at <https://arxiv.org/abs/1712.00409> (2017).
38. You, J., Ying, Z. & Leskovec, J. Design space for graph neural networks. *NeurIPS* **33**, (2020).
39. Anthony, W., Steven, K., Ryan, M. & Taylor, S. Compositionally-restricted attention-based network for materials property prediction. Preprint at https://chemrxiv.org/articles/preprint/Compositionally-Restricted_Attention-Based_Network_for_Materials_Property_Prediction/11869026 (2020).
40. Peterson, G. G. C. & Brgoch, J. Materials discovery through machine learning formation energy. *J. Phys. Energy* **3**, 022002 (2021).
41. Goodall, R. E. A. & Lee, A. A. Predicting materials properties without crystal structure: deep representation learning from stoichiometry. *Nat. Commun.* **11**, 6280 (2020).
42. Qiao, Z., Welborn, M., Anandkumar, A., Manby, F. R. & MillerIII, T. F. OrbNet: deep learning for quantum chemistry using symmetry-adapted atomic-orbital features. *J. Chem. Phys.* **153**, 124111 (2020).
43. Chen, C. et al. Learning properties of ordered and disordered materials from multi-fidelity data. *Nat. Comput. Sci.* **1**, 46–53 (2021).
44. Mamun, O., Winther, K. T., Boes, J. R. & Bligaard, T. High-throughput calculations of catalytic properties of bimetallic alloy surfaces. *Sci. Data* **6**, 76 (2019).
45. Rosen, A. S. et al. Machine learning the quantum-chemical properties of metal–organic frameworks for accelerated materials discovery. *Matter* **4**, 1578–1597 (2021).
46. Hastrup, S. et al. The Computational 2D Materials Database: high-throughput modeling and discovery of atomically thin crystals. *2D Mater.* **5**, 042002 (2018).
47. Fung, V. & D-e, J. Exploring structural diversity and fluxionality of Ptn (n = 10–13) clusters from first-principles. *J. Phys. Chem. C* **121**, 10796–10802 (2017).
48. Unke, O. T. et al. Machine learning force fields. Preprint at arXiv: 201007067, (2020).
49. Behler, J. Constructing high-dimensional neural network potentials: a tutorial review. *Int. J. Quantum Chem.* **115**, 1032–1050 (2015).
50. Kipf, T. N. & Welling, M. Semi-supervised classification with graph convolutional networks. Preprint at <https://arxiv.org/abs/1609.02907> (2016).
51. Bergstra, J., Yamins, D. & Cox, D. Making a science of model search: hyperparameter optimization in hundreds of dimensions for vision architectures. *ICML* 115–123 (2013).
52. Faber, F., Lindmaa, A., von Lilienfeld, O. A. & Armiento, R. Crystal structure representations for machine learning models of formation energies. *Int. J. Quantum Chem.* **115**, 1094–1101 (2015).
53. Fey, M. & Lenssen, J. E. Fast graph representation learning with PyTorch Geometric. Preprint at arXiv:190302428 (2019).
54. Himanen, L. et al. DScribe: library of descriptors for machine learning in materials science. *Comput. Phys. Commun.* **247**, 106949 (2020).
55. Liaw, R. et al. Tune: a research platform for distributed model selection and training. Preprint at <https://arxiv.org/abs/1807.05118> (2018).

ACKNOWLEDGEMENTS

This work was supported by the Center for Understanding and Control of Acid Gas-Induced Evolution of Materials for Energy (UNCAGE-ME), an Energy Frontier Research Center funded by the U.S. Department of Energy, Office of Science, Basic Energy Sciences. Work was performed at the Center for Nanophase Materials Sciences, which is a US Department of Energy Office of Science User Facility. V.F. was also supported by a Eugene P. Wigner Fellowship at Oak Ridge National Laboratory. J.Z. was supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Applied Mathematics program; and by the Artificial Intelligence Initiative at the Oak Ridge National Laboratory (ORNL). ORNL is operated by UT-Battelle, LLC., for the U.S. Department of Energy under Contract DEAC05-00OR22725. This research used resources of the National Energy Research Scientific Computing Center, supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231.

AUTHOR CONTRIBUTIONS

V.F. conceived the project, conducted the calculations, and analyzed the results. V.F., J.Z., E.J., and B.G.S. contributed to the discussion of results and the writing of the manuscript.

COMPETING INTERESTS

The authors declare no competing interests.

ADDITIONAL INFORMATION

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41524-021-00554-0>.

Correspondence and requests for materials should be addressed to V.F.

Reprints and permission information is available at <http://www.nature.com/reprints>

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



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021