# Multi-task learning for molecular electronic structure approaching coupled-cluster accuracy

Hao Tang,[1] Brian Xiao,[2] Wenhao He,[3] Pero Subasic,[4] Avetik R. Harutyunyan,[4] Yao Wang,[5] Fang Liu,[5] Haowei Xu,[6, *] and Ju Li[1, 6, †]

[1]*Department of Materials Science and Engineering,*
*Massachusetts Institute of Technology, MA 02139, USA*
[2]*Department of Physics, Massachusetts Institute of Technology, MA 02139, USA*
[3]*The Center for Computational Science and Engineering,*
*Massachusetts Institute of Technology, Cambridge, MA 02139, USA*
[4]*Honda Research Institute USA, Inc., San Jose, CA 95134, USA*
[5]*Department of Chemistry, Emory University, Atlanta, GA 30322, USA*
[6]*Department of Nuclear Science and Engineering,*
*Massachusetts Institute of Technology, Cambridge, MA 02139, USA*
(Dated: June 26, 2024)

Machine learning (ML) plays an important role in quantum chemistry, providing fast-to-evaluate predictive models for various properties of molecules. However, most existing ML models for molecular electronic properties use density functional theory (DFT) databases as ground truth in training, and their prediction accuracy cannot surpass that of DFT. In this work, we developed a unified ML method for electronic structures of organic molecules using the gold-standard CCSD(T) calculations as training data. Tested on hydrocarbon molecules, our model outperforms DFT with the widely-used hybrid and double hybrid functionals in computational costs and prediction accuracy of various quantum chemical properties. As case studies, we apply the model to aromatic compounds and semi-conducting polymers on both ground state and excited state properties, demonstrating its accuracy and generalization capability to complex systems that are hard to calculate using CCSD(T)-level methods.

## I. INTRODUCTION

Computational methods for molecular and condensed matter systems play essential roles in physics, chemistry, and materials science, which can reveal underlying mechanisms of diverse physical phenomena and accelerate materials design [1, 2]. Among various types of computational methods, quantum chemistry calculations of electronic structure are usually the bottleneck, limiting the computational speed and scalability [3]. In recent years, machine learning (ML) methods have been successfully applied to accelerate molecular dynamics simulations and to improve their accuracy in many application scenarios [4–6]. Particularly, ML inter-atomic potential can predict energy and force of molecular systems with significantly lower computational costs compared to quantum chemistry methods [5–8]. Indeed, recent advances in universal ML potential enable large-scale molecular dynamics simulation with the complexity of realistic physical systems [9–13]. Besides ML inter-atomic potential, rapid advances also appear in another promising direction, namely the ML density functional, which focuses on further improving the energy prediction towards the chemical accuracy (1 kcal/mol) [14–17].

Besides energy and force, other electronic properties that explicitly involve the electron degrees of freedom are also essential in molecular simulations [18]. In the past few years, ML methods have also been extended to electronic structure of molecules, predicting various electronic properties such as electric multipole moments [19–21], electron population [22], excited states properties [23, 24], as well as the electronic band structure of condensed matter [25, 26]. Most of these methods take the density functional theory (DFT) results as the training data, using neural networks (NN) to fit the single-configurational representation (either the Kohn-Sham Hamiltonian or molecular orbitals) of the DFT calculations [19, 23, 25, 27]. Along with the rapid advances of ML techniques, the NN predictions match the DFT results increasingly well, approaching the chemical accuracy [9, 20]. However, as a mean-field level theory, DFT calculations themselves induce a systematic error, which is usually several times larger than the chemical accuracy [28], limiting the overall accuracy of the ML model trained on DFT datasets.

In comparison, the correlated wavefunction method CCSD(T) is considered the gold-standard in quantum chemistry [29]. It provides high accuracy predictions on various molecular properties. Unfortunately, the computational cost of CCSD(T) calculations has a rather unfavorable scaling relationship with system size. Hence, it can only handle small molecules with up to hundreds of electrons. This urges the combination of CCSD(T) and ML methods, which can potentially have both high accuracy and low computational cost. However, above-mentioned ML methods that directly fit the single-configurational representation of the DFT calculations cannot be applied with the CCSD(T) training data. This
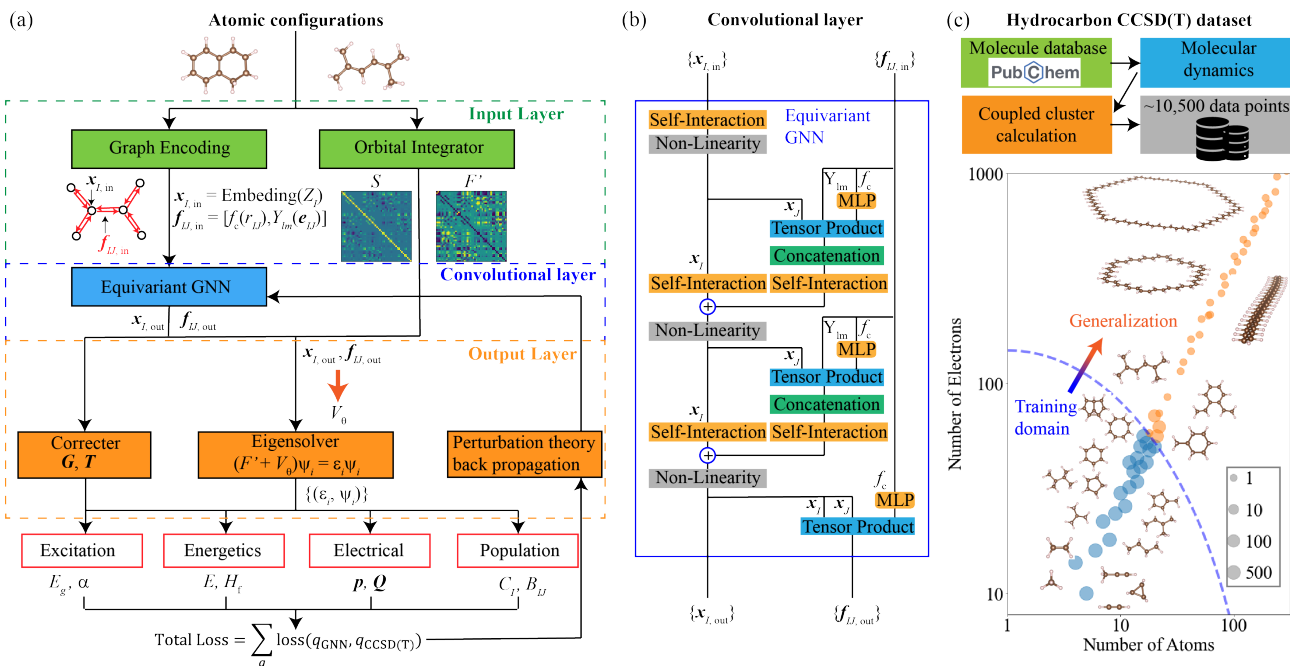
FIG. 1. Schematic of the EGNN electronic structure workflow. (a) Computation graph of the EGNN method that calculate multiple quantum chemical properties from atomic configurations inputs. The computational graph consists of input layer (green blocks), convolutional layer (blue block), and output layer (orange blocks). (b) Model architecture of the EGNN that consists of two layers of graph convolution and output both node feature $\mathbf{x}_{I,\text{out}}$ and edge feature $\mathbf{f}_{IJ,\text{out}}$. (c) Training and testing dataset generation. About 10,500 atomic configurations of 85 different hydrocarbon molecules are sampled from molecular dynamics trajectories. Data points are plot in the map of number of electrons and atoms, and the dot size reflects the number of training data with the same chemical formula.

is because CCSD(T) does not provide either a Kohn-Sham (KS) Hamiltonian or single-body molecular orbitals due to the many-body quantum entanglement nature of its representation.

In this work, we develop a unified multi-task ML method for molecular electronic structures. Instead of focusing solely on energy, our method also provides accurate predictions for various electronic properties; compared with ML models trained on DFT datasets, our method learns from CCSD(T)-accuracy training data. The method incorporates the equivariant graph neural network (EGNN) [5–7, 30], where vectors and tensors are involved in the message-passing step. Using hydrocarbon organic molecules as a testbed, our method predicts molecular energy within chemical accuracy as compared with both CCSD(T) calculation and experiments, and predicts various properties, including electric dipole and quadrupole moments, atomic charge, bond order, energy gap, and electric polarizability with better accuracy than B3LYP, one of the most widely used hybrid DFT functional [31]. Our trained model shows robust generalization capability from small molecules in the training dataset (molecular weight $< 100$) to larger molecules and even semiconducting polymers (molecular weight up to several thousands). Systematically predicting multiple electronic properties using a single model with local DFT computational speed, the method provides a

high-performance tool for computational chemistry and a promising framework for ML electronic structure calculations.

## II. RESULTS

### A. Theory and Model Architecture

In this section, we briefly describe the theoretical background and model architecture of our EGNN method. Basically, we use a NN to simulate the non-local exchange-correlation interactions of a many-body system. Then, a physics-informed approach is used to predict multiple properties from the output of a single NN.

**Computational Workflow.** Given an input atomic configuration, our methods output an effective single-body Hamiltonian matrix that predicts quantum chemical properties, as shown in Fig. 1a. The workflow consists of the input layer, the convolutional layer, and the output layer.

The input layer takes atomic configurations as input, including the information of atomic numbers $(Z_1, Z_2, \cdots, Z_n)$ and atomic coordinates $(\vec{r}_1, \vec{r}_2, \cdots, \vec{r}_n)$ of a $n$-atom system. A molecular graph is constructed, where atoms are mapped to graph nodes, while bonds be-

tween atoms (neighboring atoms within a cut-off radius $r_{\mathrm{cut}} = 2$ Å) are mapped to graph edges. The atomic configuration is then encoded into the node features $\mathbf{x}_{I,\mathrm{in}}$ for atom information and edge features $\mathbf{f}_{IJ,\mathrm{in}}$ for bond information (see Methods A for details). The electron wavefunction is represented using an atomic orbital basis set $\{|\phi_{I,\mu}\rangle\}$ [32], where $I$ is the index of atom and $\mu$ is the index of atomic orbital basis function. Then, a quantum chemistry calculation [33] (the Orbital Integrator block in Fig. 1a) is used to evaluate single-body effective Hamiltonian $F_{I\mu,J\nu}$ and the overlap matrix $S_{I\mu,J\nu} \equiv \langle \phi_{I,\mu} | \phi_{J,\nu} \rangle$ in the non-orthogonal atomic-orbital representation (where $I\mu$ is the row index and $J\nu$ is the column index). The $F_{I\mu,J\nu}$ matrix is obtained by a fast-to-evaluate single-configurational method such as local DFT [34]. The Lowdin-symmetrized KS Hamiltonian [35] is then obtained as

$$\mathbf{F}' \equiv \mathbf{S}^{-1/2}\mathbf{F}\mathbf{S}^{-1/2} + \frac{E_{\mathrm{MB}}}{n_e}\mathbf{I}, \qquad (1)$$

where the last term is an identity shift to account for the many-body energy term (see Methods A for details). Note that $\mathbf{F}'$ is a local DFT-level effective Hamiltonian, meaning that it is easy and fast to compute, but its accuracy is relatively low. We will use $\mathbf{F}'$ as the starting point of our ML model, and the total effective Hamiltonian of the system $\mathbf{H}^{\mathrm{eff}} = \mathbf{F}' + \mathbf{V}^{\theta}$ is obtained by adding the ML correction term $\mathbf{V}^{\theta}$. As $\mathbf{F}'$ is obtained from a local DFT calculation, it only contains local exchange-correlation contribution, and the correction term $\mathbf{V}^{\theta}$ would account for the non-local exchange-correlation effects. Generally, the non-local exchange-correlation effects can be incorporated in CCSD(T) calculations. However, as mentioned before, the computational costs of CCSD(T) methods are formidably high for large system. The essence of our ML method is to obtain the non-local exchange-correlation effects from a NN, whose computational cost scales only *linearly* with system size.

To obtain the ML correction term, we build a neural network model to predict $\mathbf{V}^{\theta}$. The EGNN framework is employed for the convolutional layer because of its outstanding performance in predicting molecular properties [7]. The convolutional layer transforms the input node features $\mathbf{x}_{I,\mathrm{in}}$ and edge features $\mathbf{f}_{IJ,\mathrm{in}}$ to the output node features $\mathbf{x}_{I,\mathrm{out}}$ and edge features $\mathbf{f}_{IJ,\mathrm{out}}$. The EGNN includes a series of linear transformation (Self-Interaction block), activation function (Non-Linearity block), and graph convolution (Tensor Product and Concatenation blocks) layers, as shown in Fig. 1b. Details on the numerical form and dimension of each block are elaborated in Methods B. The output $\mathbf{f}_{IJ,\mathrm{out}}$ and $\mathbf{x}_{I,\mathrm{out}}$ encode equivariant features of atoms and bonds as well as their atomic environment.

Then, we construct an equivariant ML correction Hamiltonian $\mathbf{V}^{\theta}$ from the output features $\mathbf{x}_{I,\mathrm{out}}$ and

$\mathbf{f}_{IJ,\mathrm{out}}$ as follow:

$$V_{I\mu,J\nu}^{\theta} = \begin{cases} [V_{\mathrm{node}}(\mathbf{x}_{I,\mathrm{out}})]_{\mu,\nu} & \text{if } I = J \\ \frac{1}{2}[V_{\mathrm{edge}}(\mathbf{f}_{IJ,\mathrm{out}})]_{\mu,\nu} + \frac{1}{2}[V_{\mathrm{edge}}(\mathbf{f}_{JI,\mathrm{out}})]_{\nu,\mu} & \text{if } I \neq J \end{cases}$$
$$(2)$$

where $V_{\mathrm{node}}(\mathbf{x}_{I,\mathrm{out}})$ is a $N_I \times N_I$ symmetric matrix rearranged from node features $\mathbf{x}_{I,\mathrm{out}}$, while $V_{\mathrm{edge}}(\mathbf{f}_{IJ,\mathrm{out}})$ is a $N_I \times N_J$ matrix obtained from edge features $\mathbf{f}_{JI,\mathrm{out}}$. Here $N_I, N_J$ are the numbers of basis functions of the atom $I, J$. Note that the output matrices $\mathbf{V}^{\theta}$ are Hermitian and equivariant under rotation according to the transformation rule of the basis set $\{|\phi_{I,\mu}\rangle\}$ (see Methods B for details).

An effective electronic structure of the molecule is obtained by solving the eigenvalue equations of the total Hamiltonian $H_{I\mu,J\nu}^{\mathrm{eff}}$:

$$\sum_{J,\nu} H_{I\mu,J\nu}^{\mathrm{eff}} c_{J,\nu}^i = \epsilon_i c_{I,\mu}^i, \qquad (3)$$

where $\epsilon_i$ gives the $i$-th energy levels, and $c_{I,\mu}^i$ gives the corresponding molecular orbitals through basis expansion $|\psi_i\rangle = \sum_{I,\mu} \tilde{c}_{I,\mu}^i |\phi_{I,\mu}\rangle$ with $\tilde{c}^i = S^{-1/2}c^i$.

**Multiple Learning Tasks.** Our scheme aims to predict multiple observable molecular properties (more than just energy). In order to achieve reduced computational costs, we do not include information about the entire electronic Hilbert space as learning targets. The energy levels and molecular orbitals are used to evaluate a series of ground state properties $O_g$ according to the rules of quantum mechanics:

$$O_g^{\mathrm{EGNN}} = f_{O_g}(\{\epsilon_i\}, \{\mathbf{c}^i\}), \quad O_g = E, \vec{p}, \mathbf{Q}, C_I, B_{IJ}, \quad (4)$$

where properties $O_g$ goes through the ground state energy ($E$), electric dipole ($\vec{p}$) and quadrupole ($\mathbf{Q}$) moments, Mulliken atomic charge [36] of each atom $C_I$, and Mayer bond order [37] of each pair of atoms $B_{IJ}$. We also evaluate the energy gap (first excitation energy, $E_g$) and static electric polarizability $\alpha$:

$$\begin{aligned} E_g^{\mathrm{EGNN}} &= f_{E_g}(\{\epsilon_i\}, \{\mathbf{c}^i\}, \mathbf{G}), \\ \alpha^{\mathrm{EGNN}} &= f_\alpha(\{\epsilon_i\}, \{\mathbf{c}^i\}, \mathbf{T}). \end{aligned} \qquad (5)$$

In principle, the ground state electronic structure does not contain the information of the energy gap and electric polarizability. Therefore, we use the EGNN-output correction terms $\mathbf{G}$ (energy gap correction) and $\mathbf{T}$ (dielectric screening matrix) to account for the excited states information and the linear response information, respectively. The function forms of $f_{O_g}$, $f_{E_g}$, and $f_\alpha$ are elaborated in Methods C. Note that these properties are all derived from the underlying electronic structure, so they are internally related. Therefore, multi-task learning methods can utilize these relations to enhance the model's generalization capability.

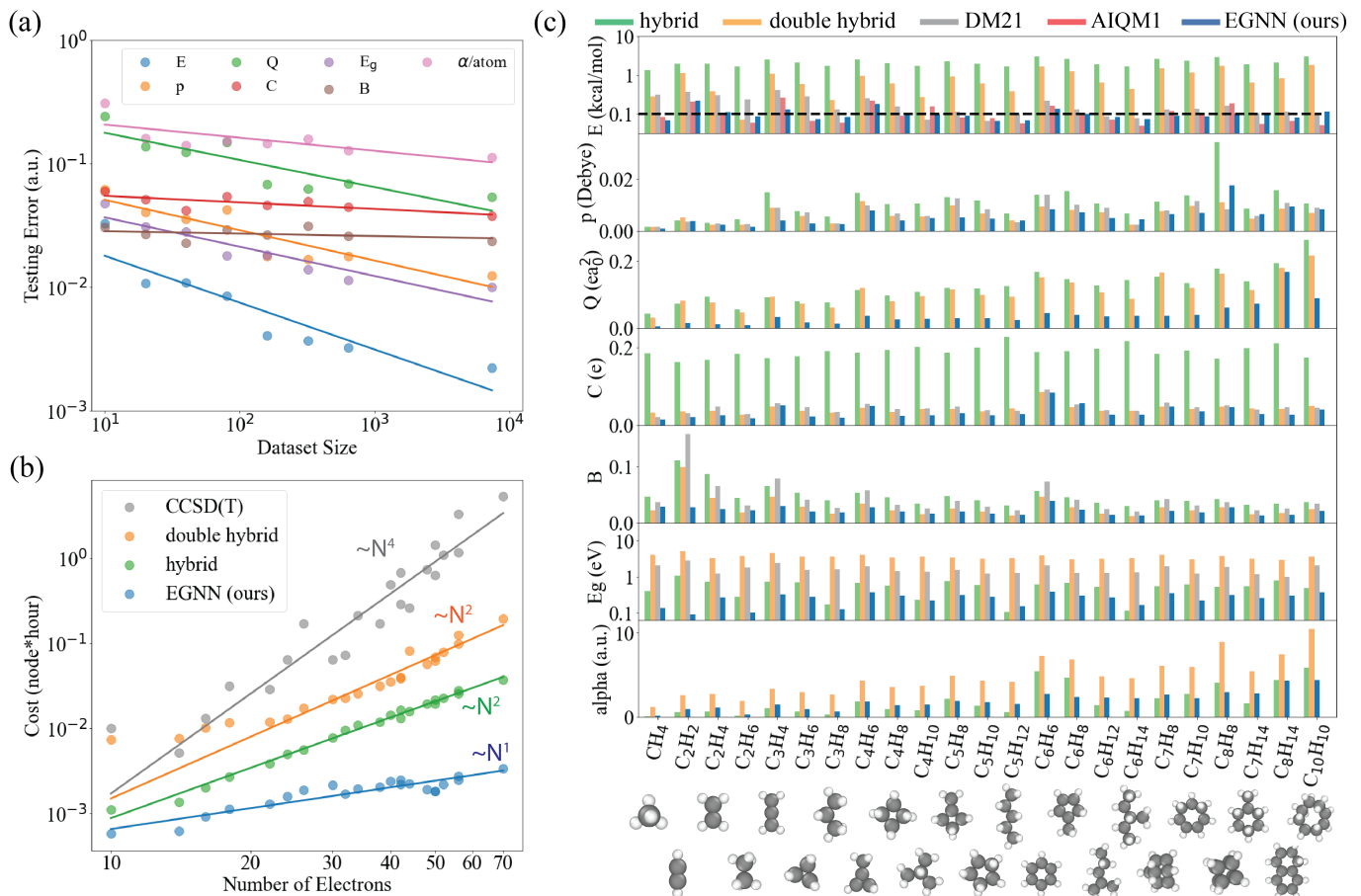The multi-task learning is implemented by minimizing

FIG. 2. Benchmark of the model performance on the testing dataset. (a) Testing root-mean-square errors (RMSE) of different quantities as a function of training dataset size. (b) Computational costs of different methods plot against number of electrons. The computational costs is measured as the calculation time (node hour) on a single Intel Xeon Platinum 8260 CPU node with 48 cores on the MIT SuperCloud [38] with sufficient memory for all calculations. The scaling deviates from the theoretical asymptotic scaling, e.g., $N^7$ for CCSD(T), because the parallelization efficiency is higher for larger molecules. In principle, the $N^7$ scaling for CCSD(T) would appear in the large $N$ limit. (c) Prediction RMSE of the energy ($E$ per atom, reference to separate atoms), electric dipole moment ($\vec{p}$), electric quadrupole moment ($\mathbf{Q}$), Mulliken atomic charge ($C$), Mayer bond order ($B$), energy gap (1$^{\text{st}}$ excitation energy, $E_{\text{g}}$), and static electric polarizability ($\alpha$, a.u. means atomic unit). Our EGNN method is compared with the B3LYP hybrid functional, DSD-PBEP86 double hybrid functional [39], DM21 ML functional [14], and AIQM1 ML potential [13, 40]. A representative atomic configuration of each chemical formula is plotted for illustration.

the total loss function $L_{\text{Total}}$ constructed as follows:

$$L_{\text{Total}} = \sum_{O \in \{O_g, E_g, \alpha\}} l_O + l_V,$$

$$l_O = w_O \times \text{MSEloss}(O^{\text{EGNN}}, O^{\text{label}}), \quad (6)$$

$$l_V = \frac{w_V}{N_{\text{basis}}^2} \sum_{I\mu, J\nu} |V_{I\mu, J\nu}^{\theta}|^2.$$

Here for each property $O$, $l_O$ is the the mean-square error (MSE) loss between $O^{\text{EGNN}}$ and $O^{\text{label}}$, the EGNN predictions (Eq. (4)(5)) and coupled-cluster labels in the training dataset, respectively. Meanwhile, $l_V$ is a regularization that penalizes large correction matrix $\mathbf{V}_\theta$, and $N_{\text{basis}}$ is the total number of basis functions in the molecule. The weights $w_V$ and $w_O$ are hyperparameters whose values are listed in supplementary informa-

tion (SI) section I. The weights are chosen to balance the training tasks so that the training errors of all tasks decrease to satisfactory levels. Minimizing $L_{\text{Total}}$ requires the back-propagation through Eq. (3) (i.e., calculating $\partial \epsilon_i / \partial \mathbf{H}^{\text{eff}}$ and $\partial \mathbf{c}^i / \partial \mathbf{H}^{\text{eff}}$), which is numerically unstable [41]. To overcome this issue, we derive customized back-propagation schemes for each property using perturbation theory in quantum mechanics (see SI section S1 for details), giving

$$\nabla_\theta \epsilon_i = (\mathbf{c}^i)^\dagger (\nabla_\theta \mathbf{V}^\theta) \mathbf{c}^i$$

$$\nabla_\theta \mathbf{c}^i = \sum_{p \neq i} \frac{(\mathbf{c}^p)^\dagger (\nabla_\theta \mathbf{V}^\theta) \mathbf{c}^i}{\epsilon_i - \epsilon_p} \mathbf{c}^p. \quad (7)$$

When evaluating the gradients of properties in Eqs. (4) and (5) using the chain rule, terms that analytically can-

TABLE I. Benchmark of EGNN model's RMSE in predicting different quantum chemical properties on the in-domain (ID) testing dataset and out-of-domain (OOD) validation dataset (written as ID RMSE/OOD RMSE in the table).

| RMSE (ID/OOD) | Unit | Hybrid | Double hybrid | DM21 | AIQM1 | EGNN (ours) |
|---|---|---|---|---|---|---|
| Energy (per atom) | kcal/mol | 2.20/2.41 | 0.94/1.20 | 0.22/0.11 | 0.13/0.06 | **0.11/0.10** |
| Electric Dipole | $10^{-2}$ Debye | 1.27/1.20 | 0.71/0.70 | 0.78/0.88 | – | **0.63/0.81** |
| Electric Quadrupole | $ea_0^2$ | 0.12/0.21 | 0.11/0.18 | – | – | **0.03/0.12** |
| Mulliken Atomic Charge | e | 0.19/0.20 | 0.04/0.05 | 0.05/0.04 | – | **0.04/0.03** |
| Mayer Bond Order | – | 0.05/0.03 | 0.04/0.02 | 0.06/0.03 | – | **0.02/0.02** |
| $1^{st}$-Excitation Energy | eV | 0.59/0.63 | 3.71/3.26 | 1.71/1.47 | – | **0.26/0.31** |
| Polarizability | a.u. | 2.22/4.32 | 4.74/8.05 | – | – | **1.85/3.91** |

cel each other are removed in numerical evaluation, making the scheme numerically stable.

**Software Realization.** Atomic configurations of molecules in our training dataset are generated by the workflow shown in Fig. 1c. First, 85 hydrocarbon molecule structures are collected from the PubChem database [42]. Molecular dynamics (MD) simulation with TeaNet interatomic potential [5, 6] is then performed for each molecule structure to sample an ensemble of atomic configurations. Then, the CCSD(T) calculations are used to calculate the ground-state properties of sampled configurations, and the EOM-CCSD [43] calculations are used to calculate the excited-state properties. The **S** and **F** matrices are evaluated by the ORCA quantum chemistry program package [33] with the fast-to-evaluate BP86 local density functional [34] and the medium-sized cc-pVDZ basis set [32]. The EGNN model based on the e3nn package [30] is trained on small-molecules training dataset (training domain, Fig. 1c) and tested on both small molecules in the training dataset (in-domain validation) and larger molecules outside the training dataset (out-of-domain validation). Details about training and testing dataset and training hyperparameters are elaborated in SI section S2.

### B. Model Performance and Applications

In this section, we benchmark the performance of our EGNN model and display potential applications of the model in systems of practical importance. We will make direct comparisons with experimental results wherever possible, and the results suggest the outstanding generalization capability and predicting power of our approach.

**Benchmark of Model Performance.** The model's generalization capability from small molecules to large molecules is essential for its usefulness on complex systems where coupled cluster calculations cannot be implemented on current computational platforms, due to their formidable computational costs. To test the generalization capability and data efficiency of our model, we train the model with varied training dataset size $N_{train}$, rang-

ing from 10 to 7440 atomic configurations. The testing root-mean-square error (RMSE, absolute error in atomic units) of different trained properties exhibit a decreasing trend when the training dataset size increases, as shown in Fig. S1a, indicating effective model generalization. Notably, the energy error has the fastest drops with a slope of −0.38 (that means the testing error $\propto N_{train}^{-0.38}$). In comparison, some of the recently developed advanced ML potentials (that directly learn energies and their derivatives) exhibits lower slopes of about -0.25 [7, 9]. This implies potential advantage of our multi-task method: as our multi-task method learns different molecular properties through a shared representation (the electronic structure), the domain information learned from one property can help the model's generalization on predicting other properties [44], providing outstanding data efficiency.

Then, we benchmark the computational costs and prediction accuracy of our model trained on 7440 atomic configurations with 70 different molecules, which will be used in the rest of this paper. Our EGNN method exhibits significantly smaller computational cost and slower scaling with system size, as compared with the hybrid functional, double hybrid functional [39], and CCSD(T) method (Fig. S1b). Compared to the hybrid functional, our method avoids the expensive calculation of the exact exchange thus requiring much smaller computational costs [31]. Using the gold-standard CCSD(T) calculation as a reference, the prediction accuracy of our EGNN method on various molecular properties is compared with that of the several theory functionals and existing ML methods, as shown in Fig. S1c and Table I. The comparison is implemented on both an in-domain (ID) and an out-of-domain (OOD) testing dataset of hydrocarbon molecules. The EGNN predictions exhibit smaller RMSE than the hybrid, double hybrid, and DM21 [14] functional on most compared molecular properties (except electric dipole moment on the OOD dataset, where double hybrid gives the smallest RMSE). Remarkably, the RMSE of the combination energy predicted by the EGNN is about 0.1 kcal/mol (about 4 meV) per atom in both the ID and OOD dataset. Our method exhibits similar RMSE of combination energy compared to the AIQM1 ML po-
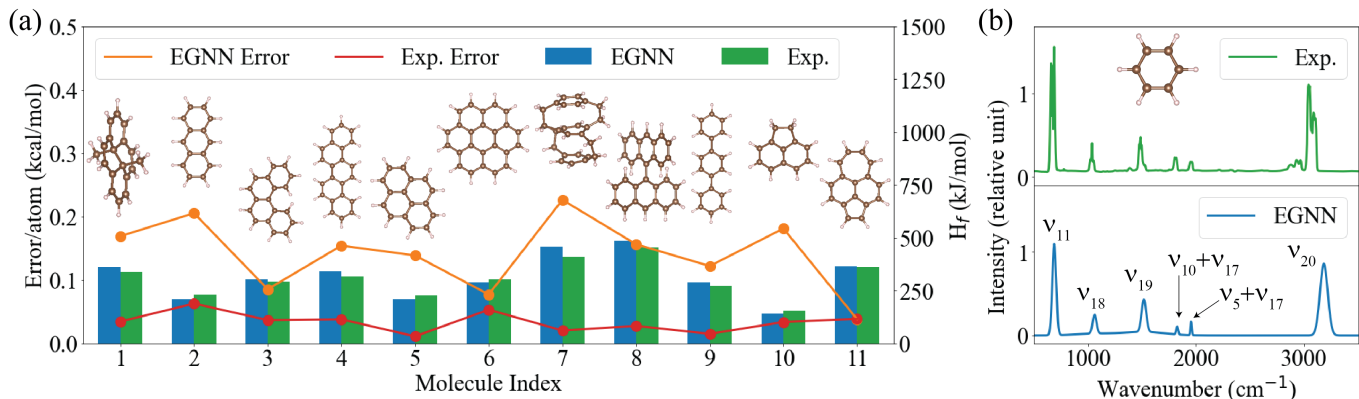
FIG. 3. Validation of the EGNN's predictions on gas phase aromatic hydrocarbon molecules, as compared with experimental results. (a) Standard enthalpy of formation. The EGNN predictions and experimental values from Ref. [45] (right axis) are compared for 11 molecules (see SI Table II for details). The difference between the EGNN method and experimental values are shown by the orange line, and the experimental uncertainty is shown by the red line (left axis). (b) Infrared spectrum of benzene. The experimental data is from the NIST Chemistry WebBook [46]. Vibration modes corresponding to the peaks are labeled following the convention in Ref. [47].

tential that features its chemical accuracy energy predictions. These results confirm that our EGNN's predictions on reaction energies can approach the quantum chemical accuracy (assuming that on average 1 mol molecules in reactants contain ~10 mol atoms). Note that besides the ground-state properties, our EGNN also provides excited-state property $E_g$ and linear response property $\alpha$ with better overall accuracy than other methods in comparison (Table I), though there are certain molecules where the hybrid functional performs better. The statistical distribution of the prediction errors of the EGNN and B3LYP hybrid functional are shown in Fig. S1 in SI. Although the hybrid functional gives smaller median error for $\alpha$, the EGNN is less likely to make large errors (more robust), leading to better root-mean-square error.

**Aromatic Molecules.** Hydrocarbon molecules have a vast structural space, including various types of local atomic environments (see a few examples below Fig. S1c). Our EGNN method provides a single model that exhibits generalization capability among the vast hydrocarbon structural space (cf. Fig. S1). To further examine the model's generalization capability in more complex structures, we apply the EGNN model to a series of aromatic hydrocarbon molecules synthesized in experiments [45]. The gas phase standard enthalpy of formation $H_f$ is an essential thermochemical property of molecules that can be accurately measured in experiments. In this regard, we use the EGNN model to predict $H_f$ of various aromatic molecules in a comprehensive experimental review paper Ref. [45]. Among the 71 molecules lists in Ref. [45], we calculate $H_f$ of those molecules with serial numbers (defined in Ref. [45]) dividable by 4 and compare them with experiments, as shown in Fig. 3a. The selected molecules cover various classes of aromatic molecules, including polycyclic aromatic hydrocarbons, cyclophanes, polyphenyl, and nonalternant hydrocarbons. The EGNN

predictions on $H_f$ are well consistent with experiments on all molecules, and their difference is only around $0.1 \sim 0.2$ kcal/mol per atom. Note that the EGNN prediction error is on the same order of magnitude as the experimental error bar (though numerically larger), indicating high prediction accuracy.

Besides thermochemical properties, our EGNN model can also predict spectral properties, as shown in Fig. 3b. Especially, infrared (IR) spectrum reflects essential information of molecular vibrational modes and their interaction with light. In previous work of ML electronic structure [20], although the predicted peak positions of the IR spectrum are well consistent with the experiment, the predicted peak intensity are inconsistent with the experiment. In comparison, our EGNN model predicts both the peak positions and intensity well consistent with the experiment, providing both the fundamental bands and combination bands known as "benzene fingers" in the IR spectrum. The good consistency of peak intensity is attributed to accurate predictions on the transition dipole moments that determine the intensity of light-matter interaction. Details on calculating the IR spectrum is elaborated in SI section S3.

**Large-scale Semiconducting Polymers.** Besides small molecules, we also apply the EGNN model to semiconducting polymers consisting of hundreds of atoms, which are difficult to calculate by rigorous correlated methods such as CCSD(T). Semiconducting polymers are organic macromolecules with small energy gap and high electrical conductivity compared to insulating polymers. Because of these electronic features, semiconducting polymers attract broad research interests both for the fundamental understanding of quantum chemistry [49] and for applications in semiconductor industry [52]. The essential electronic properties of semiconducting polymers originates from the $\pi$-bonds with delocalized molec-

header

## (a) Molecular orbitals

trans-polyacetylene (t-PA)

HOMO

LUMO

cyclic polyacetylene (c-PA)

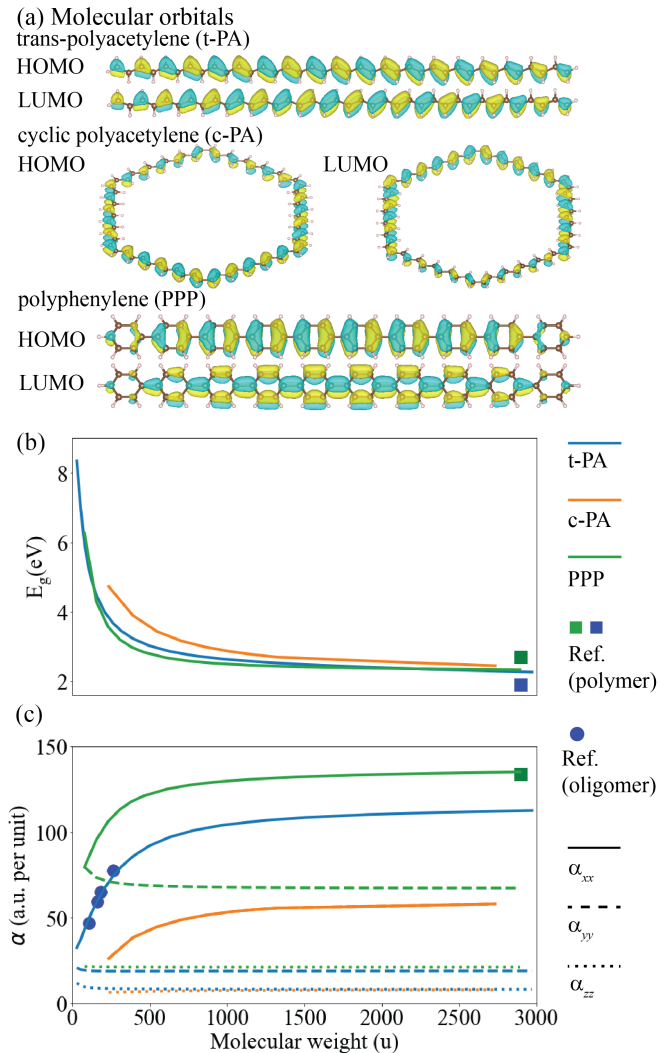HOMO                    LUMO

polyphenylene (PPP)

HOMO

LUMO



FIG. 4. EGNN predictions for the electronic proerpties of semiconducting polymers. (a) Atomic structure and HOMO wavefunctions of t-PA, polyphenylene PPP, and c-PA. The HOMO wavefunctions are visualized by isosurfaces at the level of $\pm 0.01$ Å$^{-2/3}$ (positive isosurface colored blue and negative isosurface colored yellow). (b) Energy gap and (c) static electric polarizability of t-PA (blue lines), PPP (green lines), and c-PA (orange lines) with different polymer chain length. Longitudinal polarizability $\alpha_{xx}$, horizontal polarizability $\alpha_{yy}$, and vertical polarizability $\alpha_{zz}$ are shown as solid, dashed, and dotted lines, respectively. Squares (blue for t-PA and green for PPP) represent literature values for polymers in experiments [48, 49] and correlated calculations [50], and blue dots represent literature values for t-PA oligomers from the MP2 correlated calculations [51].

ular orbitals. As the delocalized molecular orbitals extend through the whole macroscopic molecule (Fig. 4a), the polymers' electronic properties also involve long-range correlation, making it challenging for ML methods. Therefore, it is important to examine whether the EGNN model can capture semiconducting polymers' electronic properties involving delocalized molecular orbitals.

Three kinds of semiconducting polymers, trans-polyacetylene (t-PA), cyclic polyacetylene (c-PA), and polyphenylene (PPP), are studied using our EGNN model. The model correctly captures the delocalized $\pi$-bond feature of frontier orbitals (highest occupied molecular orbital (HOMO) and lowest unoccupied molecular orbital (LUMO)), as shown in Fig. 4a. In the t-PA, the HOMO and LUMO are on the carbon-carbon double and single bonds, respectively, consistent with the predictions from the renowned Su–Schrieffer–Heeger (SSH) model [49]. In the c-PA, the left/right mirror symmetry of the molecule is correctly reflected in the predicted frontier orbitals, which exhibits odd parity upon mirror reflection. Both the HOMO and LUMO have two anti-phase boundaries where the wavefunctions shift from double bonds to single bonds. The same symmetry is also reflected in the PPP with odd-parity HOMO and even-parity LUMO. It should be emphasized that the molecular orbitals are not included as fitting targets in our model training. However, as the model learned relevant properties such as electric moments, atomic charge, and bond order, it gives at least qualitatively correct predictions for frontier orbitals in these semiconducting polymers.

Various important electronic properties of semiconducting polymers depend on the chain length, including the energy gap $E_g$ and polarizability $\alpha$. We calculate such chain-length dependence (up to more than 400 atoms) using our EGNN method, as shown in Fig. 4b,c. One can see that $E_g$ is larger for oligomers with shorter chain length due to the quantum confinement effect and converges to a smaller value for long polymer. The converged energy gap for long t-PA and PPP polymers calculated by our EGNN are in reasonable agreement with the experimental values [48, 49] (relative errors within 10%), which are shown as squares in Fig. 4b. The longitudinal static electric polarizability $\alpha_{xx}$ (per monomer) is positively related to the polymer chain length. This is because in longer chains, more delocalized electron distributions can have larger displacements under external electric field. The predicted $\alpha_{xx}$ for t-PA oligomers and PPP polymers are in perfect agreement with previous correlated calculations using the high-accuracy MP2 method [50, 51], as shown in Fig. 4c. The chain length-dependent $E_g$ and $\alpha$ of cyclic PA, to the best of our knowledge, have not been reported. We provide their values as a prediction to be examined by future work.

## III. CONCLUSION AND OUTLOOK

In this work, we developed a multi-task learning method for predicting molecular electronic structure with coupled-cluster-level accuracy. In our method, an EGNN is trained on a series of molecular properties in order to capture their shared underlying representation, the effective electronic Hamiltonian. Our EGNN model shows significantly lower computational costs and higher pre-

diction accuracy on various molecular properties than the widely-used B3LYP functional.

Compared to other works that take the electronic structure as a direct fitting target [19, 25–27], our method takes the electronic structure as a shared representation to predict various molecular properties. This approach, on the one hand, enables training on data beyond the DFT accuracy; on the other hand, relieve the burden of fitting all matrix elements of the electronic Hamiltonian, which are not direct quantum observable. The later enables us to go beyond the minimal basis used in previous works by a relatively small NN with only 511,589 parameters. With the electronic structure as a physics-informed representation, our method exhibits generalization capability for challenging systems with delocalized molecular orbitals. In such systems, atomic structure at one position has long-range influence to electronic features (such as electron density) at another distant position. Although the direct output of our EGNN, the matrix elements of $\mathbf{V}_\theta$, only depends on local atomic environments, the long-range influence is captured through the output-layer calculations based on the rules of quantum mechanics. In comparison, such long-range influence is unlikely to be captured by directly using graph neural network or local atomic descriptors to predict molecular properties without a physics-informed representation.

Although the current EGNN model is limited to hydrocarbon molecules, our method can be readily applied to systems with different elements. Applying the method to datasets with more elements can produce a general-purpose electronic structure predictor, which is left to future work.

## IV. METHODS

### A. Graph encoding of atomic configuration

The atomic numbers $Z_I$ of input elements are encoded as node features $\mathbf{x}_{I,\text{in}}$ by one-hot embedding. In our case, hydrogen and carbon atoms are encoded as scalar arrays $[0, 1]$ and $[1, 0]$, respectively. The atomic coordinates are encoded as edge features $\mathbf{f}_{IJ,\text{in}} \equiv [f_c(r_{IJ}), Y_{lm}(\vec{e}_{IJ})]$, where $f_c(r) \equiv \frac{1}{2}\left[\cos\left(\pi \frac{r}{r_{\text{cut}}}\right) + 1\right]$ is a smooth cut-off function reflecting the bond length $r_{IJ} \equiv |\vec{r}_I - \vec{r}_J|$ [53], and $Y_{lm}(\vec{e}_{IJ})$ is the spherical harmonic functions acting on the unit vector $\vec{e}_{IJ} \equiv \frac{\vec{r}_I - \vec{r}_J}{|\vec{r}_I - \vec{r}_J|}$ representing the bond orientation [30]. We include $Y_{lm}$ tensors up to $l = 2$.

As the hydrocarbon molecules we study are all close-shell molecules, we use spin-restricted DFT calculations to obtain $\mathbf{F}$ and we assume $\mathbf{V}_\theta$ is spin-independent throughout this paper. Namely, the spin-up and spin-down molecular orbitals and energy levels are the same, and all molecular orbitals are either doubly occupied or vacant. The electronic energy of the BP86 DFT calculation $E_{\text{BP86}}$ equals the band structure energy $2\sum_{i=1}^{n_e/2} \epsilon_i$ (where $\epsilon_i$ is the $i$-th molecular orbital energy level and

$n_e$ is the number of electrons) plus a many-body energy $E_{\text{MB}}$:

$$E_{\text{BP86}} = 2\sum_{i=1}^{n_e/2} \epsilon_i + E_{\text{MB}} \qquad (8)$$

$E_{\text{MB}}$ originates from the double-counting of electron-electron interaction in the band structure energy and is obtained from the output of the ORCA BP86 DFT calculation. In order to incorporate the many-body energy into the KS effective Hamiltonian $\mathbf{F}$, we construct $\mathbf{F}'$ as in Eq. (1), so the direct summation of band structure energies given by $F$ equals:

$$\begin{aligned}
2\sum_{i=1}^{n_e/2} \text{eig}_i(\mathbf{F}') &= 2\sum_{i=1}^{n_e/2} \text{eig}_i(\mathbf{S}^{-1/2}\mathbf{F}\mathbf{S}^{-1/2} + \frac{E_{\text{MB}}}{n_e}\mathbf{I}) \\
&= 2\sum_{i=1}^{n_e/2} \left[\text{eig}_i(\mathbf{S}^{-1/2}\mathbf{F}\mathbf{S}^{-1/2}) + \frac{E_{\text{MB}}}{n_e}\right] \\
&= 2\sum_{i=1}^{n_e/2} \epsilon_i + E_{\text{MB}},
\end{aligned}$$
(9)

where $\text{eig}_i$ is a function that returns the $i$-th lowest eigenvalue of a matrix, and we use the fact that the energy level $\epsilon_i$ is the eigenvalue of the Lowdin-symmetrized Hamiltonian $\text{eig}_i(S^{-1/2}FS^{-1/2})$ [35]. After this transformation, the KS effective Hamiltonian (both before and after the $\mathbf{V}_\theta$ correction) already includes the many-body energy, and the total electronic energy is just the summation of band structure energies. Adding the $E_{\text{MB}}$ term does not change the eigenfunction and relative energy levels, so that all other properties are not changed by adding the $E_{\text{MB}}$ term.

### B. Architecture of the convolutional layer

In the following technical description, we use terminologies defined in the e3nn documentation. One can refer to Ref. [30] for further information. In the convolutional layer, the input feature first goes through a $N_{\text{species}} \times N_{\text{species}}$ linear transformation (the first Self-Interaction block, $N_{\text{species}}$ is the number of different elements in the system) and an activation layer (the first non-Linearity block). All activation layers in our EGNN are realized by tanh function acting on scalar features.

Then, the input features go through the first-step convolution: a fully connected tensor product (the first Tensor Product block) of node feature $\mathbf{x}_J$ and the spherical Harmonic components of all connected edge feature $\mathbf{f}_{IJ}$ mapping to an irreducible representation "8x0e + 8x1o + 8x2e" (denoted as Irreps1), meaning 8 even scalar, 8 odd vector, and 8 even rank-2 tensor. Weights in the fully connected tensor product are from a multilayer perceptron (the first MLP block) taking $f_c(r_{IJ})$ as input. All MLP blocks in Fig. 1b has a $1 \times 16 \times 16 \times 16 \times N_w$ structure

and tanh activation function, where $N_w$ is the number of weights in the tensor product. Then, in the Concatenation block, tensor products from different edges $\mathbf{f}_{IJ}$ connected to the node $I$ are summed to a new node feature on $I$. The new node features then go through a linear transformation (Self-Interaction block) that outputs the same data type Irreps1. In all Self-Interaction layers, linear combinations are only applied to features with the same tensor order. The new node features are added to the original node features undergoes a linear transformation, complete the first-step convolution.

The second step convolution has the same architecture. The only difference is that the second Tensor Product block takes input node features of Irreps1 and output features of "8x0e + 8x0o + 8x1e + 8x1o + 8x2e + 8x2o" (denoted as Irreps2), meaning 8 even and odd scalar, vector, and tensor, respectively. The output of both Self-Interaction blocks are also Irreps2.

After another activation function, the node features are output as $\mathbf{x}_{I,\text{out}}$. Another tensor product acting on the node features of two endpoints of each edge is applied to get the output bond feature $\mathbf{f}_{IJ,\text{out}}$, also with a dimension of Irreps2 and weight parameters from the MLP taking $f_c(r_{IJ})$ as input.

Finally, the output features are used to construct the correction matrix $\mathbf{V}_\theta$ through Eq. (2). $V_{\text{node}}(\mathbf{x}_{I,\text{out}})$ first apply a linear layer from input dimension of Irreps2 to output dimension Irreps$(I)^{\otimes 2}$, where:

$$\text{Irreps}(I) = \begin{cases} (2 \times 0e + 1 \times 1o) & \text{if } I \text{ is H} \\ (3 \times 0e + 2 \times 1o + 1 \times 2e) & \text{if } I \text{ is C} \end{cases} \quad (10)$$

The output dimension corresponds to the irreducible representation of the block diagonal terms of the Hamiltonian (as cc-pVDZ basis of hydrogen includes two s orbitals (0e) and one group of p orbitals (1o); while that of carbon includes three s orbitals (0e), two groups of p orbitals (1o), and one group of $d$ orbitals (2e)). The output is then arranged into the $N_I \times N_I$ matrix form, $V_{I,\text{out}}$, according to the Wigner-Eckart theorem [26], and symmetrized to obtain $V_{\text{node}}(\mathbf{x}_{I,\text{out}}) = \frac{\lambda_V}{2}(V_{I,\text{out}} + V_{I,\text{out}}^T)$. $\lambda$ is a constant hyperparameter set as 0.2 for our model. Similarly, the off-diagonal term $V_{\text{edge}}(\mathbf{f}_{IJ,\text{out}})$ in Eq. (2) applies a linear layer from input dimension of Irreps2 to output dimension Irreps$(I, J)$ that equals:

$$\text{Irreps}(I, J) = \text{Irreps}(I) \otimes \text{Irreps}(J) \quad (11)$$

The output are then arranged into the $N_I \times N_J$ matrix and multiplied by $\lambda_V$, giving $V_{\text{edge}}(\mathbf{f}_{IJ,\text{out}})$.

In addition, the energy gap correction term $\mathbf{G}$ is obtained from a $8 \times 32 \times 3$ MLP that takes the even scalars of $\mathbf{x}_{I,\text{out}}$ as input and output a 3-component scalar array, $\mathbf{g}_{I;0,1,2}$, with tanh activation. The first component is for attention pooling:

$$\mathbf{G}_K = \sum_I \frac{e^{g_{I,0}}}{\sum_J e^{g_{J,0}}} g_{I,K}, \qquad K = 1, 2 \quad (12)$$

Giving the two-component bandgap correction array $\mathbf{G}$. The polarizability correction term, the screening matrix $\mathbf{T}$ is obtained from the edge features $\mathbf{f}_{IJ,\text{out}}$ going through a Irreps2 to $32 \times 0e + 1 \times 2e$ linear layer, an tanh activation layer, and a $32 \times 0e + 1 \times 2e$ to $1 \times 0e + 1 \times 2e$ linear layer. The $1 \times 0e + 1 \times 2e$ array is then multiplied by a factor $\lambda_T$ (set as 0.01 in our case) arranged into the symmetric matrix $\mathbf{T}$'s 6 independent components.

## C. Evaluating molecular properties

The ground state properties in Eq. (4) is evaluated by the electronic structure by the following equations [36, 37]:

$$E^{\text{EGNN}} = E_{\text{NN}} + 2 \sum_{i=1}^{n_e/2} \epsilon_i$$

$$\vec{p}^{EGNN} = -2e \sum_{i=1}^{n_e/2} \sum_{I\mu,J\nu} (\tilde{c}_{I,\mu}^i)^* \tilde{c}_{J,\nu}^i \langle \phi_{I,\mu}|\hat{\vec{r}}|\phi_{J,\nu}\rangle$$

$$\mathbf{Q}^{\text{EGNN}} = 2e^2 \sum_{i=1}^{n_e/2} \sum_{I\mu,J\nu} (\tilde{c}_{I,\mu}^i)^* \tilde{c}_{J,\nu}^i \langle \phi_{I,\mu}|\hat{\vec{r}}\hat{\vec{r}}|\phi_{J,\nu}\rangle$$

$$C_I^{\text{EGNN}} = e\left[Z_I - 2 \sum_{i=1}^{n_e/2} \sum_{J\mu\nu}(\tilde{c}_{I,\mu}^i)^* \tilde{c}_{J,\nu}^i S_{I\mu,J\nu}\right]$$

$$B_{IJ}^{\text{EGNN}} = 4 \sum_{i,j=1}^{n_e/2} \sum_{KL\mu\nu\lambda\sigma} (\tilde{c}_{K,\lambda}^i)^* \tilde{c}_{I,\mu}^i S_{K\lambda,J\nu}(\tilde{c}_{L,\sigma}^j)^* \tilde{c}_{J,\nu}^j S_{L\sigma,I\mu}$$

$$(13)$$

where $E_{\text{NN}}$ is the Coulomb repulsion energy between nuclei and nuclei, and $e$ and $\hat{\vec{r}}$ are the electron charge and position operator, respectively.

Besides, Using the ground state electronic structure obtained from Eq. (3), $E_{\text{g}}$ can be roughly estimated as $\epsilon_{n_e/2+1} - \epsilon_{n_e/2}$, the energy difference between the HOMO and LUMO (abbreviated as the HOMO-LUMO gap). However, in principle, the ground state electronic structure $(\epsilon_n, \mathbf{c}^n)$ does not contain the information of excited states (once a electron is excited, $\epsilon_n$ and $\mathbf{c}^n$ undergo relaxation and become different). Therefore, we use the EGNN to output two correction terms $G_1$ and $G_2$. $E_{\text{g}}$ is then evaluated as a linear transformation of the HOMO-LUMO gap using $G_1$ and $G_2$ as the coefficients:

$$E_{\text{g}}^{\text{EGNN}} = (1 + G_1)(\epsilon_{n_e/2+1} - \epsilon_{n_e/2}) + G_2 \quad (14)$$

Evaluation of the static electric polarizability is done in two steps. First, we evaluate the single-particle polarizability $\alpha_0$ using perturbation theory:

$$\alpha_0 = 2e^2 \sum_{a=n_e/2+1}^{N_{\text{basis}}} \sum_{i=1}^{n_e/2} \frac{\vec{r}_{ai}\vec{r}_{ia}}{\epsilon_a - \epsilon_i} \quad (15)$$

where $N_{\text{basis}}$ is the number of basis functions of the molecule, and $\vec{r}_{ai} \equiv \sum_{I\mu,J\nu}(\tilde{c}_{I,\mu}^a)^* \tilde{c}_{J,\nu}^i \langle \phi_{I,\mu}|\hat{\vec{r}}|\phi_{J,\nu}\rangle$.

However, the single-particle approximation used in Eq. (15) does not consider the electric screening effect from electron-electron interaction. We use the EGNN to output a screening matrix $T$ and evaluate the corrected polarizability $\alpha$ as follow:

$$\alpha^{\text{EGNN}} = (\mathbf{I} + \alpha_0 \mathbf{T})^{-1} \alpha_0. \tag{16}$$

We evaluate the gas phase standard enthalpy of formation of molecules in Fig. 3 using atomic configurations built by Avogadro [54] and relaxed by the BP86 functional with cc-pVDZ basis set. The total energy at the relaxed atomic configuration is then calculated by our EGNN model. The zero-point energy (ZPE) and thermal vibration, rotation, and translation energy at $T = 298.15$ K are also calculated by the BP86 functional with cc-pVDZ basis set implemented in ORCA. The ZPE is corrected by the optimal scaling factor of 1.0393 according to Ref. [55]. Summing all energy terms give the inner energy $U$, and the enthalpy is evaluated as $H \simeq U + N k_{\text{B}} T$ ($N$ is the number of molecules and $k_{\text{B}}$ is the Boltzmann constant), where we use the ideal gas law. To obtain the standard enthalpy of formation, we subtract the reference state enthalpy of graphite and hydrogen gas at standard condition. The reference enthalpy for each carbon and hydrogen atom are determined as -38.04639 a.u. and -0.57550 a.u., respectively, using CCSD(T) calculation with cc-pVTZ basis set combined with measured standard enthalpy of formation of atomic carbon, atomic hydrogen, and benzene. Atomic configurations of semiconducting polymers in Fig. 4 are relaxed using the Pre-Ferred Potential v5.0.0 [5, 6].

## V.   DATA AVAILABILITY

Detailed data of our benchmark test results (Fig. 2, Fig. S1, Table 1) and the calculation results of aromatic molecules (Fig. 3) and semiconducting polymers (Fig. 4) will be made available through figshare. The training and testing dataset is available upon reasonable requests to the corresponding authors.

## VI.   CODE AVAILABILITY

The source code to generate the training dataset, train the EGNN model, and apply the trained EGNN model to hydrocarbon molecules has been deposited into a publicly available GitHub repository https://github.com/htang113/Multi-task-electronic.

## VII.   ACKNOWLEDGEMENTS

[1] S. Yip, Introduction, in *Handbook of Materials Modeling: Methods*, edited by S. Yip (Springer Netherlands, Dordrecht, 2005) pp. 1–5.

[2] E. A. Carter, Challenges in modeling materials properties without experimental input, Science **321**, 800 (2008).

[3] H. J. Kulik, T. Hammerschmidt, J. Schmidt, S. Botti, M. A. Marques, M. Boley, M. Scheffler, M. Todorović, P. Rinke, C. Oses, *et al.*, Roadmap on machine learning in electronic structure, Electronic Structure **4**, 023004 (2022).

[4] P. O. Dral, *Quantum Chemistry in the Age of Machine Learning* (Elsevier, 2022).

[5] S. Takamoto, S. Izumi, and J. Li, Teanet: Universal neural network interatomic potential inspired by iterative electronic relaxations, Comput. Mater. Sci. **207**, 111280 (2022).

[6] S. Takamoto, D. Okanohara, Q. Li, and J. Li, Towards universal neural network interatomic potential, J. Materiomics **9**, 447 (2023).

[7] S. Batzner, A. Musaelian, L. Sun, M. Geiger, J. P. Mailoa, M. Kornbluth, N. Molinari, T. E. Smidt, and B. Kozinsky, E (3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials, Nature communications **13**, 2453 (2022).

[8] L. Zhang, J. Han, H. Wang, R. Car, and E. Weinan, Deep potential molecular dynamics: a scalable model with the accuracy of quantum mechanics, Physical review letters **120**, 143001 (2018).

[9] A. Merchant, S. Batzner, S. S. Schoenholz, M. Aykol, G. Cheon, and E. D. Cubuk, Scaling deep learning for materials discovery, Nature **624**, 80 (2023).

[10] C. Chen and S. P. Ong, A universal graph deep learning interatomic potential for the periodic table, Nature Computational Science **2**, 718 (2022).

[11] S. Takamoto, C. Shinagawa, D. Motoki, K. Nakago, W. Li, I. Kurata, T. Watanabe, Y. Yayama, H. Iriguchi, Y. Asano, *et al.*, Towards universal neural network potential for material discovery applicable to arbitrary combination of 45 elements, Nature Communications **13**, 2991 (2022).

[12] J. S. Smith, B. T. Nebgen, R. Zubatyuk, N. Lubbers, C. Devereux, K. Barros, S. Tretiak, O. Isayev, and A. E. Roitberg, Approaching coupled cluster accuracy with a general-purpose neural network potential through transfer learning, Nature communications **10**, 2903 (2019).

[13] P. Zheng, R. Zubatyuk, W. Wu, O. Isayev, and P. O. Dral, Artificial intelligence-enhanced quantum chemical method with broad applicability, Nature communications

**12**, 7022 (2021).

[14] J. Kirkpatrick, B. McMorrow, D. H. Turban, A. L. Gaunt, J. S. Spencer, A. G. Matthews, A. Obika, L. Thiry, M. Fortunato, D. Pfau, *et al.*, Pushing the frontiers of density functionals by solving the fractional electron problem, Science **374**, 1385 (2021).

[15] R. Pederson, B. Kalita, and K. Burke, Machine learning and density functional theory, Nature Reviews Physics **4**, 357 (2022).

[16] M. Bogojeski, L. Vogt-Maranto, M. E. Tuckerman, K.-R. Müller, and K. Burke, Quantum chemical accuracy from density functional approximations via machine learning, Nature communications **11**, 5223 (2020).

[17] K. Bystrom and B. Kozinsky, Cider: An expressive, non-local feature set for machine learning density functionals with exact constraints, Journal of Chemical Theory and Computation **18**, 2180 (2022).

[18] T. Helgaker, P. Jorgensen, and J. Olsen, *Molecular electronic-structure theory* (John Wiley & Sons, 2013).

[19] K. T. Schütt, M. Gastegger, A. Tkatchenko, K.-R. Müller, and R. J. Maurer, Unifying machine learning and quantum chemistry with a deep neural network for molecular wavefunctions, Nature communications **10**, 5024 (2019).

[20] X. Shao, L. Paetow, M. E. Tuckerman, and M. Pavanello, Machine learning electronic structure methods based on the one-electron reduced density matrix, Nature communications **14**, 6281 (2023).

[21] C. Feng, J. Xi, Y. Zhang, B. Jiang, and Y. Zhou, Accurate and interpretable dipole interaction model-based machine learning for molecular polarizability, Journal of Chemical Theory and Computation **19**, 1207 (2023).

[22] G. Fan, A. McSloy, B. Aradi, C.-Y. Yam, and T. Frauenheim, Obtaining electronic properties of molecules through combining density functional tight binding with machine learning, The Journal of Physical Chemistry Letters **13**, 10132 (2022).

[23] E. Cignoni, D. Suman, J. Nigam, L. Cupellini, B. Mennucci, and M. Ceriotti, Electronic excited states from physically constrained machine learning, ACS Central Science (2023).

[24] P. O. Dral and M. Barbatti, Molecular excited states through a machine learning lens, Nature Reviews Chemistry **5**, 388 (2021).

[25] H. Li, Z. Wang, N. Zou, M. Ye, R. Xu, X. Gong, W. Duan, and Y. Xu, Deep-learning density functional theory hamiltonian for efficient ab initio electronic-structure calculation, Nature Computational Science **2**, 367 (2022).

[26] X. Gong, H. Li, N. Zou, R. Xu, W. Duan, and Y. Xu, General framework for e (3)-equivariant neural network representation of density functional theory hamiltonian, Nature Communications **14**, 2848 (2023).

[27] O. Unke, M. Bogojeski, M. Gastegger, M. Geiger, T. Smidt, and K.-R. Müller, Se (3)-equivariant prediction of molecular wavefunctions and electronic densities, Advances in Neural Information Processing Systems **34**, 14434 (2021).

[28] N. Mardirossian and M. Head-Gordon, Thirty years of density functional theory in computational chemistry: an overview and extensive assessment of 200 density functionals, Molecular physics **115**, 2315 (2017).

[29] R. J. Bartlett and M. Musiał, Coupled-cluster theory in quantum chemistry, Reviews of Modern Physics **79**, 291

[30] (2007).

[30] M. Geiger and T. Smidt, e3nn: Euclidean neural networks, arXiv preprint arXiv:2207.09453 (2022).

[31] J. Tirado-Rives and W. L. Jorgensen, Performance of b3lyp density functional methods for a large set of organic molecules, Journal of chemical theory and computation **4**, 297 (2008).

[32] T. H. Dunning Jr, Gaussian basis sets for use in correlated molecular calculations. i. the atoms boron through neon and hydrogen, The Journal of chemical physics **90**, 1007 (1989).

[33] F. Neese, F. Wennmohs, U. Becker, and C. Riplinger, The orca quantum chemistry program package, The Journal of chemical physics **152** (2020).

[34] J. P. Perdew, Density-functional approximation for the correlation energy of the inhomogeneous electron gas, Physical review B **33**, 8822 (1986).

[35] P.-O. Löwdin, On the non-orthogonality problem connected with the use of atomic wave functions in the theory of molecules and crystals, The Journal of Chemical Physics **18**, 365 (1950).

[36] R. S. Mulliken, Electronic population analysis on lcao–mo molecular wave functions. i, The Journal of chemical physics **23**, 1833 (1955).

[37] I. Mayer, Bond order and valence indices: A personal account, Journal of computational chemistry **28**, 204 (2007).

[38] A. Reuther, J. Kepner, C. Byun, S. Samsi, W. Arcand, D. Bestor, B. Bergeron, V. Gadepally, M. Houle, M. Hubbell, M. Jones, A. Klein, L. Milechin, J. Mullen, A. Prout, A. Rosa, C. Yee, and P. Michaleas, Interactive supercomputing on 40,000 cores for machine learning and data analysis, in *2018 IEEE High Performance extreme Computing Conference (HPEC)* (IEEE, 2018) pp. 1–6.

[39] S. Kozuch and J. M. Martin, Spin-component-scaled double hybrids: an extensive search for the best fifth-rung functionals blending dft and perturbation theory, Journal of Computational Chemistry **34**, 2327 (2013).

[40] P. O. Dral, F. Ge, Y.-F. Hou, P. Zheng, Y. Chen, M. Barbatti, O. Isayev, C. Wang, B.-X. Xue, M. Pinheiro Jr, *et al.*, Mlatom 3: A platform for machine learning-enhanced computational chemistry simulations and workflows, Journal of Chemical Theory and Computation (2024).

[41] https://pytorch.org/docs/stable/generated/torch.linalg.eigh.html, (2023).

[42] S. Kim, J. Chen, T. Cheng, A. Gindulyte, J. He, S. He, Q. Li, B. A. Shoemaker, P. A. Thiessen, B. Yu, L. Zaslavsky, J. Zhang, and E. E. Bolton, PubChem 2023 update, Nucleic Acids Research **51**, D1373 (2022), https://academic.oup.com/nar/article-pdf/51/D1/D1373/48441598/gkac956.pdf.

[43] A. I. Krylov, Equation-of-motion coupled-cluster methods for open-shell and electronically excited species: The hitchhiker's guide to fock space, Annu. Rev. Phys. Chem. **59**, 433 (2008).

[44] R. Caruana, Multitask learning, Machine learning **28**, 41 (1997).

[45] S. W. Slayden and J. F. Liebman, The energetics of aromatic hydrocarbons: an experimental thermochemical perspective, Chemical reviews **101**, 1541 (2001).

[46] P. Linstrom and E. W.G. Mallard, Nist chemistry webbook, nist standard reference database number 69 (National Institute of Standards and Technology, Gaithers-

burg MD, 20899, 2024) Chap. Infrared Spectra.

[47] E. B. Wilson Jr, The normal modes and frequencies of vibration of the regular plane hexagon model of the benzene molecule, Physical Review **45**, 706 (1934).

[48] G. Grem, G. Leditzky, B. Ullrich, and G. Leising, Realization of a blue-light-emitting device using poly (p-phenylene), Advanced Materials **4**, 36 (1992).

[49] A. J. Heeger, Nobel lecture: Semiconducting and metallic polymers: The fourth generation of polymeric materials, Reviews of Modern Physics **73**, 681 (2001).

[50] P. Otto, M. Piris, A. Martinez, and J. Ladik, Dynamic (hyper) polarizability calculations for polymers with linear and cyclic $\pi$-conjugated elementary cells, Synthetic metals **141**, 277 (2004).

[51] B. Champagne, E. A. Perpete, S. J. Van Gisbergen, E.-J. Baerends, J. G. Snijders, C. Soubra-Ghaoui, K. A. Robins, and B. Kirtman, Assessment of conventional density functional schemes for computing the polarizabilities and hyperpolarizabilities of conjugated oligomers: An ab initio investigation of polyacetylene chains, The Journal of chemical physics **109**, 10489 (1998).

[52] B. Geffroy, P. Le Roy, and C. Prat, Organic light-emitting diode (oled) technology: materials, devices and display technologies, Polymer international **55**, 572 (2006).

[53] L. Zhang, J. Han, H. Wang, W. Saidi, R. Car, *et al.*, End-to-end symmetry preserving inter-atomic potential energy model for finite and extended systems, Advances in neural information processing systems **31** (2018).

[54] M. D. Hanwell, D. E. Curtis, D. C. Lonie, T. Vandermeersch, E. Zurek, and G. R. Hutchison, Avogadro: an advanced semantic chemical editor, visualization, and analysis platform, Journal of cheminformatics **4**, 1 (2012).

[55] M. K. Kesharwani, B. Brauer, and J. M. Martin, Frequency and zero-point vibrational energy scale factors for double-hybrid density functionals (and other selected methods): can anharmonic force fields be avoided?, The Journal of Physical Chemistry A **119**, 1701 (2015).

[56] Q. Sun, X. Zhang, S. Banerjee, P. Bao, M. Barbry, N. S. Blunt, N. A. Bogdanov, G. H. Booth, J. Chen, Z.-H. Cui, J. J. Eriksen, Y. Gao, S. Guo, J. Hermann, M. R. Hermes, K. Koh, P. Koval, S. Lehtola, Z. Li, J. Liu, N. Mardirossian, J. D. McClain, M. Motta, B. Mussard, H. Q. Pham, A. Pulkin, W. Purwanto, P. J. Robinson, E. Ronca, E. R. Sayfutyarova, M. Scheurer, H. F. Schurkus, J. E. T. Smith, C. Sun, S.-N. Sun, S. Upadhyay, L. K. Wagner, X. Wang, A. White, J. D. Whitfield, M. J. Williamson, S. Wouters, J. Yang, J. M. Yu, T. Zhu, T. C. Berkelbach, S. Sharma, A. Y. Sokolov, and G. K.-L. Chan, Recent developments in the PySCF program package, The Journal of Chemical Physics **153**, 024109 (2020), https://pubs.aip.org/aip/jcp/article-pdf/doi/10.1063/5.0006074/16722275/024109_1_online.pdf.

[57] B. A. Hess Jr, L. J. Schaad, P. Carsky, and R. Zahradnik, Ab initio calculations of vibrational spectra and their use in the identification of unusual molecules, Chemical Reviews **86**, 709 (1986).

[58] P. E. Maslen, N. C. Handy, R. D. Amos, and D. Jayatilaka, Higher analytic derivatives. iv. anharmonic effects in the benzene spectrum, The Journal of chemical physics **97**, 4233 (1992).

## S1. PERTURBATION THEORY-BASED BACK-PROPAGATION

In the EGNN training, gradient of the loss function to the model parameters needs to be calculated. Gradient back-propagation schemes are well-developed for all computation steps except solving the Schrodinger equation Eq. 3 in the main text. The gradients are numerically unstable when there are near-degenerate energy levels, which is usually the case in molecules. Here, we first use quantum perturbation theory to obtain the first-order change of energy levels and molecular orbitals:

$$\delta\epsilon_i = (\mathbf{c}^i)^\dagger \delta H^{\text{eff}} \mathbf{c}^i$$
$$\delta\mathbf{c}^i = \sum_{p \neq i} \frac{(\mathbf{c}^p)^\dagger \delta H^{\text{eff}} \mathbf{c}^i}{\epsilon_i - \epsilon_p} \mathbf{c}^p \tag{S1}$$

Then, we have the gradients to model parameters as:

$$\nabla_\theta \epsilon_i = (\mathbf{c}^i)^\dagger (\nabla_\theta \mathbf{V}^\theta) \mathbf{c}^i$$
$$\nabla_\theta \mathbf{c}^i = \sum_{p \neq i} \frac{(\mathbf{c}^p)^\dagger (\nabla_\theta \mathbf{V}^\theta) \mathbf{c}^i}{\epsilon_i - \epsilon_p} \mathbf{c}^p \tag{S2}$$

Using these equations, we derive the gradients of each molecule properties in Eq. 4 as follow:

$$\nabla_\theta f_E = 2 \sum_{i=1}^{n_e/2} \nabla V_{ii}$$

$$\nabla_\theta f_{\vec{p}} = -4e \sum_{i=1}^{n_e/2} \sum_{a=n_e/2+1}^{N_{\text{basis}}} \text{Re} \frac{\nabla V_{ai}}{\epsilon_i - \epsilon_a} \langle \psi_i | \hat{\vec{r}} | \psi_a \rangle$$

$$\nabla_\theta f_{\mathbf{Q}} = 4e^2 \sum_{i=1}^{n_e/2} \sum_{a=n_e/2+1}^{N_{\text{basis}}} \text{Re} \frac{\nabla V_{ai}}{\epsilon_i - \epsilon_a} \langle \psi_i | \hat{\vec{r}}\hat{\vec{r}} | \psi_a \rangle \tag{S3}$$

$$\nabla_\theta f_{C_I} = -4e \sum_{i=1}^{n_e/2} \sum_{a=n_e/2+1}^{N_{\text{basis}}} \text{Re} \frac{\nabla V_{ai} (\tilde{\mathbf{c}}^i)^\dagger (I_I S) \tilde{\mathbf{c}}^a}{\epsilon_i - \epsilon_a}$$

$$\nabla_\theta f_{B_{IJ}} = 4 \sum_{i=1}^{n_e/2} \sum_{a=n_e/2+1}^{N_{\text{basis}}} \text{Re} \frac{\nabla V_{ai}}{\epsilon_i - \epsilon_a}$$
$$\times (\tilde{\mathbf{c}}^i)^\dagger (S I_J P S I_I + S I_I P S I_J) \tilde{\mathbf{c}}^a$$

where $\nabla V_{ai} \equiv (\tilde{\mathbf{c}}^a)^\dagger (\nabla_\theta \mathbf{V}^\theta) \tilde{\mathbf{c}}^i$, the $N_{\text{basis}} \times N_{\text{basis}}$ matrix $I_J$ is identity in the block diagonal part of atom $J$ and zero elsewhere. Meanwhile, we define $P \equiv 2 \sum_{i=1}^{n_e/2} \tilde{\mathbf{c}}^i (\tilde{\mathbf{c}}^i)^\dagger$. The essential method to avoid numerical instability is to remove terms that can be proved to cancel each other. Taking $\nabla_\theta f_{\vec{p}}$ as an example: in Eq. (S2), the summation over $m$ goes through all states except $n$. But as the summed formula is antisymmetric to $m$ and $n$, the terms that $m$ goes from 1 to $n_e/2$ cancel each other. Only terms that $m$ goes from $n_e/2 + 1$ to $N_{\text{basis}}$ have a non-zero contribution to the final gradient. Therefore, $n$ is always occupied, and $m$ is always unoccupied in the summation. As close-shell molecules always

have a finite bandgap, $\epsilon_n$ and $\epsilon_m$ are not close to each other in any term of the summation, so evaluating Eq. S3 is numerically stable.

Similarly, the gradients of $E_{\text{g}}$ and $\alpha$ are as follow:

$$\nabla_\theta f_{E_{\text{g}}} = (1 + G_1) \left[ \nabla V_{n_e/2+1, n_e/2+1} - \nabla V_{n_e/2, n_e/2} \right]$$
$$+ (\epsilon_{n_e/2+1} - \epsilon_{n_e/2}) \nabla_\theta G_1 + \nabla_\theta G_2 \tag{S4}$$

To calculate the gradient of $\alpha$, we first evaluate the gradient of $\alpha_0$, and then derive $\nabla_\theta f_\alpha$ using the chain rule:

$$\nabla_\theta \alpha_0 = 2e^2 \sum_{a=n_e/2+1}^{N_{\text{basis}}} \sum_{i=1}^{n_e/2} \text{Re} \left\{ \frac{\vec{r}_{ai}\vec{r}_{ia}(\nabla V_{ii} - \nabla V_{aa})}{(\epsilon_a - \epsilon_i)^2} \right.$$
$$\left. -2 \sum_{p \neq a,i} \frac{\vec{r}_{ai}}{(\epsilon_a - \epsilon_i)} \left[ \frac{\vec{r}_{ip}\nabla V_{pa}}{(\epsilon_p - \epsilon_a)} + \frac{\vec{r}_{pa}\nabla V_{ai}}{(\epsilon_p - \epsilon_i)} \right] \right\}$$
$$\nabla_\theta f_\alpha = (I + \alpha_0 T)^{-1} (\nabla_\theta \alpha_0)(I - T\alpha)$$
$$- \alpha_0 (\nabla_\theta T)(I + \alpha_0 T)^{-1} \alpha \tag{S5}$$

The above equations give gradients of all terms in the loss function expressed by gradients to the direct outputs of the EGNN, $\nabla V$, $\nabla_\theta \mathbf{G}$, and $\nabla_\theta T$.

## S2. DATASET AND TRAINING PARAMETERS

The training domain and out-of-distribution testing testing dataset include 20 and 3 different chemical formula shown as the horizontal axis labels of the first 20 and last 3 columns in Fig. 1c in the main text, respectively. Each chemical formula includes up to 5 different molecules (conformers) taken from the PubChem database. The total number of molecules in the training domain and out-of-distribution testing dataset is 70 and 15, respectively. The full list of molecules and the number of atomic configurations for each molecule are listed in Table I in this file.

Each molecule structure is then set as the initial configuration of a MD simulation. The MD simulation uses PreFerred Potential v4.0.0 [11] at a temperature of 2000 K that enables large bond distortion but does not break the bonds. Initial velocity is set as Maxwell Boltzmann distribution with the same temperature of 2000 K. Langevin NVT dynamics is used with the friction factor of 0.001 fs$^{-1}$ and timestep of 2 fs, and one atomic configuration is sampled every 200 timesteps in the MD trajectory. 500 configurations (including the inital equilibrium configuration) are sampled for each chemical formula in the training domain, where 3/4 of the 10,000 configurations are sampled to form the training dataset, and the left 1/4 forms the in-domain testing dataset. The out-of-distribution testing dataset contains 500 configurations.

A CCSD(T) calculation with cc-pVTZ basis set is implemented in ORCA [33] for each selected configuration, giving the training labels of total energy, electric dipole and quadrupole moment, Mulliken atomic charge, and
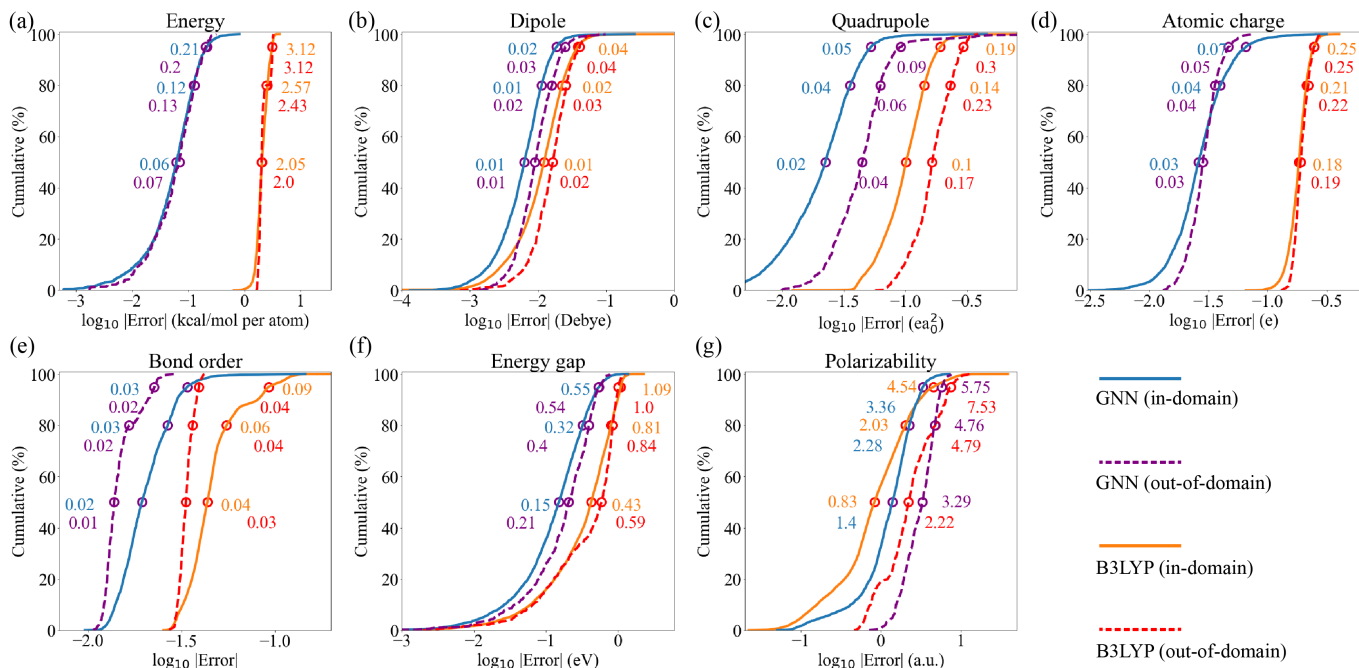
FIG. S1. The distribution of model prediction accuracy on the test dataset compared to the B3LYP DFT calculations using the CCSD(T) results as the ground truth. a-g Cumulative distribution of prediction errors for the (a) energy, (b) electric dipole moment, (c) electric quadrupole moment, (d) Mulliken atomic charge, (e) Mayer bond order, (f) energy gap (1st excitation energy), and (g) static electric polarizability (a.u. means atomic unit). The blue and orange solid lines represent EGNN and B3LYP results on the in-domain testing dataset, and the purple and red dashed lines represent GNN and B3LYP results on the out-of-domain testing dataset, respectively. We denote the model errors at 50%, 80%, and 95% percentile from the bottom to the top by hollow circles.

Mayer bond order. An EOM-CCSD calculation with cc-pVDZ basis set is then implemented to obtain the first excitation energy (bandgap). Finally, we conduct a polarizability calculation with CCSD and cc-pVDZ basis set. The overlap matrix $S$ and Effective Hamiltonian $F$ is obtained from a BP86 DFT calculations with cc-pVDZ basis set.

For comparison, B3LYP hybrid DFT calculations are implemented with def2-SVP basis set in ORCA. DSD-PBEP86 double-hybrid DFT calculations are implemented with the def2-TZVP basis set in ORCA. The DM21 double-hybrid DFT calculations are implemented with the def2-TZVP basis set in PySCF program package [56] The AIQM1 calculations are implemented in the MLatom platform [40]. For DM21 and AIQM1, the isolated-atom energies of carbon and hydrogen are recalibrated according to the least-mean-square criterion to give optimal combination energies in our dataset.

The weight parameters in the loss function is listed as follow: $w_V = 0.1$, $w_E = 1$, $w_{\vec{p}} = 0.2$, $w_Q = 0.01$, $w_C = 0.01$, $w_B = 0.02$, $w_{E_g} = 0.2$, $w_\alpha = 3 \times 10^{-5}$. All quantities are in atomic unit. The model training is implemented by full gradient descend (FGD) with Adam optimizer. For the finally deployed model (7,440 training data points), it is first trained on 1240 data points sampled from the whole dataset for 5000 FGD steps with initial learning rate of 0.01. The learning rate is decayed by a constant factor $\gamma_1 = 10^{-1/10}$ per 500 steps. Then, the model is trained on the whole dataset with 7,440 data points for 6,000 FGD steps with a initial learning rate of 0.001. For other models trained on smaller dataset in Fig. 2a in the main text, the model is trained for 3,000 FGD steps with initial learning rate of 0.01 decayed by $\gamma_2 = 10^{-1/6}$ per 500 steps. As the model trained on 640 data points do not converge in the 3000-step training, we implement 10,000-step training, initial learning rate of 0.01 decays by $\gamma_1$ every 500 steps in the first 5,000 steps and keeps constant at the last 5,000 steps.

Aromatic molecules in the main text Fig. 3 include all molecules with serial numbers dividable by 4 and enthalpy of formation provided in Ref. [45]. Details of these aromatic molecules are listed in Table II.

## S3. INFRARED SPECTRUM

In order to evaluate the infrared spectrum, we first implement a B3LYP hybrid DFT calculation with def2-TZVP basis set to obtain the vibrational modes and frequency of a benzene molecule. Then, we generate atomic configurations displaced from the equilibrium configuration along each vibrational mode by a displacement of -0.1, -0.05, 0.05, and 0.1 Å. Our EGNN model is used to evaluate the electric dipole moment at each config-

TABLE I. List of molecule names and number of atomic configurations (labelled in the superscript) for each molecule in the training and testing dataset.

| Chemical formula | Molecule 1 | Molecule 2 | Molecule 3 | Molecule 4 | Molecule 5 |
|---|---|---|---|---|---|
| $CH_4$ | Methane[500] | – | – | – | – |
| $C_2H_2$ | Acetylene[500] | – | – | – | – |
| $C_2H_4$ | Ethylene[500] | – | – | – | – |
| $C_2H_6$ | Ethane[500] | – | – | – | – |
| $C_3H_4$ | Propyne[300] | Allene[100] | Cyclopropene[100] | – | – |
| $C_3H_6$ | Propylene[260] | Cyclopropane[250] | – | – | – |
| $C_3H_8$ | Propane[500] | – | – | – | – |
| $C_4H_6$ | 1,2-Butadiene[100] | 1,3-Butadiene[100] | 1-Butyne[100] | 2-Butyne[100] | 1-Methylcyclopropene[100] |
| $C_4H_8$ | Isobutylene[100] | Cyclobutane[100] | 1-Butene[100] | 2-Butene[100] | Methylcyclopropane[100] |
| $C_4H_{10}$ | Butane[250] | Isobutane[250] | – | – | – |
| $C_5H_8$ | Isoprene[100] | Cyclopentene[100] | 1-Pentyne[100] | Methylene-cyclobutane[100] | 1,3-Pentadiene[100] |
| $C_5H_{10}$ | Cyclopentane[100] | 1-Pentene[100] | 2-Methyl-1-Butene[100] | 2-Methyl-2-Butene[100] | 3-Methyl-1-Butene[100] |
| $C_5H_{12}$ | Neopentane[200] | Isopentane[200] | Pentane[100] | – | – |
| $C_6H_6$ | Benzene[100] | 1,5-Hexadiyne[100] | 2,4-Hexadiyne[100] | Divinylacetylene[100] | 3,4-Dimethylene-cyclobut-1-ene[100] |
| $C_6H_8$ | 1,3-Cyclohexadiene[100] | 1,4-Cyclohexadiene[100] | Hexa-1,3,5-triene[100] | Methyl-cyclopentadiene[100] | Divinylethylene[100] |
| $C_6H_{12}$ | Methyl-cyclopentane[100] | Cyclohexane[100] | 1-Hexene[100] | cis-4-Methyl-2-pentene[100] | 2-Methyl-1-Pentene[100] |
| $C_6H_{14}$ | 2,2-Dimethylbutane[100] | 2,3-Dimethylbutane[100] | 3-Methylpentane[100] | 2-Methylpentane[100] | Hexane[100] |
| $C_7H_8$ | Toluene[100] | 2,5-Norbornadiene[100] | Quadricyclane[100] | 1,6-Heptadiyne[100] | Cycloheptatriene[100] |
| $C_7H_{10}$ | Norbornene[100] | 1,3-Cycloheptadiene[100] | 1-Methyl-1,3-cyclohexadiene[100] | 2-Methyl-1,3-cyclohexadiene[100] | 3-Methylenecyclohexene[100] |
| $C_7H_{14}$ | Methyl-cyclohexane[50] | Cycloheptane[50] | 1-Heptene[50] | (E)-4,4-Dimethyl-2-pentene[50] | trans-3-Heptene[50] |
| $C_8H_8$ | Styrene[100] | Benzocyclobutene[100] | Cubane[100] | Semibullvalene[100] | Cyclooctatetraene[100] |
| $C_8H_{14}$ | Bimethallyl[25] | Diisocrotyl[50] | 1,7-Octadiene[25] | CYCLOOCTENE[25] | (4E)-2,3-dimethylhexa-1,4-diene[25] |
| $C_{10}H_{10}$ | 1,3-Divinylbenzene[20] | Diolin[20] | 1,4-Divinylbenzene[20] | Divinylbenzene[20] | 4-Phenyl-1-butyne[20] |

uration, and the dipole-moment derivative with respect to each normal coordinate is evaluated by linear regression. The infrared band intensity of fundamental bands are then evaluated following the method in Ref. [57].

As the two combination bands at 1800 - 2000 $cm^{-1}$ are mainly contributed by $\nu_{10} + \nu_{17}$ and $\nu_5 + \nu_{17}$ [58], we generate atomic configurations displaced from the equilibrium configuration by displacement vectors of $0.1(\vec{e}_i + \vec{e}_j), 0.1(\vec{e}_i - \vec{e}_j), 0.1(-\vec{e}_i + \vec{e}_j)$, and $0.1(-\vec{e}_i - \vec{e}_j)$ Å, where $(\vec{e}_i, \vec{e}_j)$ are the pair of vibrational modes contribution to each combination band. The second-order dipole-

moment derivatives with respect to each pair of normal coordinates $\frac{\partial^2 \vec{p}}{\partial Q_i \partial Q_j}$ are then obtained by finite difference method. The leading-order anharmonic constants are also evaluated by finite difference method. These parameters are then used to calculate intensity of the combination bands by Fermi's golden rule. Using the calculated infrared spectrum peak positions and intensity, we add Gaussian broadening to each peak and fit their bandwidth to the experimental spectrum.

TABLE II. List of the names and serial numbers (superscript, defined in Ref. [45]) of aromatic molecules in the main text Fig. 3.

| Molecule index | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Name[serial number] | trans-10b,10c-dimethyl-10b,10c-dihydropyrene[12] | anthracene[20] | benzo[c]phenanthrene[24] | 5-ring phenacene, picene[28] |
| Molecule index | 5 | 6 | 7 | 8 |
| Name[serial number] | Pyrene[32] | Coronene[36] | 1,4:2,5-[2.2.2]cyclophane[40] | 9,9'-bianthryl[44] |
| Molecule index | 9 | 10 | 11 | – |
| Name[serial number] | p-terphenyls[48] | acenaphthene[64] | Aceplaidylene[68] | |